



**HAL**  
open science

# Synergie infrarouge et micro-onde pour la restitution atmosphérique

Maxime Paul

► **To cite this version:**

Maxime Paul. Synergie infrarouge et micro-onde pour la restitution atmosphérique. Géophysique [physics.geo-ph]. Université Pierre et Marie Curie - Paris VI, 2013. Français. NNT: . tel-00918775

**HAL Id: tel-00918775**

**<https://theses.hal.science/tel-00918775>**

Submitted on 17 Dec 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# THÈSE

PRÉSENTÉE À

L'UNIVERSITÉ PIERRE ET MARIE CURIE

ÉCOLE DOCTORALE : ED 129

Par Maxime PAUL

POUR OBTENIR LE GRADE DE :

Docteur

SPÉCIALITÉ : Terre, environnement, biodiversité

# SYNERGIE INFRAROUGE ET MICRO-ONDE POUR LA RESTITUTION ATMOSPHERIQUE

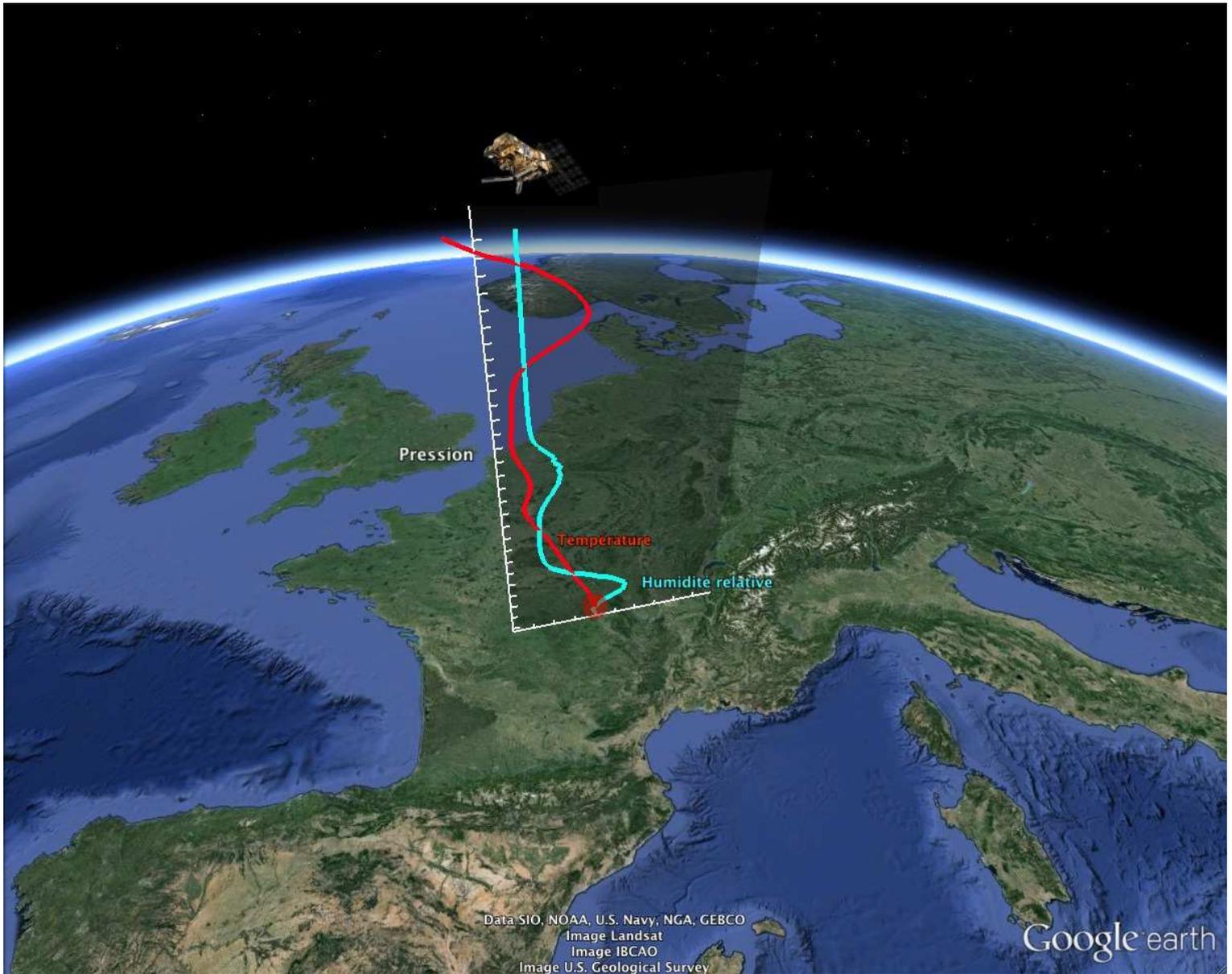
Directeur de recherche : Filipe AIRES, Estellus

Co-directeur de recherche : Catherine PRIGENT, LERMA

Soutenue le 30 septembre 2013

**Devant la commission d'examen formée de :**

Claudia STUBENRAUCH	LMD	Présidente du jury
Sid BOUKABARA	NOAA	Rapporteur
Dieter KLAES	Eumetsat	Rapporteur
Stephen ENGLISH	ECMWF	Examinateur
Hervé LE TREUT	IPSL	Examinateur
Cathy CLERBAUX	LATMOS	Invitée



# Synergie infrarouge et micro-onde pour la restitution atmosphérique

---

## Résumé

---

L'étude du climat et la météorologie nécessitent des modèles, mais également des bases de données indépendantes, issues d'observations *in situ* ou satellites. L'importance de ces dernières grandit. Nous proposons, dans cette thèse, d'optimiser leur utilisation pour restituer, à l'échelle globale, des profils atmosphériques de température et de vapeur d'eau. La contribution des surfaces terrestres sur le rayonnement conditionne la qualité des restitutions au-dessus des continents. Un schéma d'inversion bayésienne de l'équation de transfert radiatif a donc été mis au point. Il permet de restituer simultanément la température et l'émissivité hyper-spectrale infrarouge de la surface, à partir des mesures de l'instrument IASI. Une chaîne opérationnelle a été construite, entraînant la création d'une base de données d'émissivités infrarouges et de températures de surface depuis 2007. Ces informations de surface sont ensuite intégrées dans un algorithme de restitution de profils atmosphériques. Elles permettent une diminution notable de l'erreur, notamment dans les basses couches de l'atmosphère, cruciales en météorologie. Les réseaux de neurones utilisés pour les restitutions nécessitent des bases d'apprentissage. Nous avons donc mis au point une méthode d'échantillonnage de variables hétérogènes et de grande dimension. Enfin, nous avons montré que l'utilisation conjointe des observations infrarouges et micro-ondes est une source prometteuse d'amélioration de la télédétection satellite. La synergie entre des instruments tels IASI, AMSU-A et MHS sur la plateforme MetOp permet de mieux restituer les profils atmosphériques.

**Mots clés :** télédétection, synergie, émissivité de surface, analyse statistique de données, réseau de neurones, échantillonnage

---

## Infrared and microwave synergy applied to atmospheric retrievals

---

### Abstract

---

Climatology and meteorology are mainly based on numerical models, but they also need independent data from *in situ* measurements or satellite observations. In this thesis, we attempt to optimize the use of the satellite observations in order to globally retrieve atmospheric profiles of temperature and water vapor. Knowledge about the impact of land surfaces on the radiation measured by a satellite is crucial to be able to determine the quality of the retrieved profiles. A Bayesian estimator has been used to invert the radiative equation, leading to a simultaneous retrieval of surface temperature and emissivity in the infrared, based on IASI measurements. An operational algorithm has been built. It has allowed the creation of a surface emissivity and temperature database from 2007 to today. Those surface retrievals have been used in an atmospheric inversion scheme, which led to a global decrease in the error on the retrieval of temperature and water vapor profiles, especially in the troposphere, which is the most important in meteorology. The neural network-based algorithm used for the retrievals needs a representative learning database. To build such datasets, we created a multi-variate sampling method able to compute numerous non-homogeneous variables. Finally, we have shown that the simultaneous use of infrared and microwave observations is a promising way to increase the quality of the satellite retrievals. The synergy between instruments like IASI, AMSU-A and MHS on board MetOp decreases the error of the retrieved atmospheric profiles of temperature and water vapor.

**Keywords :** remote sensing, synergy, surface emissivity, statistical data analysis, neural network, sampling



# REMERCIEMENTS

Ce travail a été financé en collaboration entre l'Université Pierre et Marie Curie et le Collège Des Ingénieurs. Je les remercie de m'avoir permis de mener cette étude. Je remercie également l'Observatoire de Paris et plus particulièrement le Laboratoire d'Étude du Rayonnement en Astrophysique de m'avoir accueilli. Le cadre splendide m'a aidé à m'épanouir et à travailler sainement et sereinement.

Je remercie aussi mes directeurs de thèse : Filipe AIRES et Catherine PRIGENT. Leurs conseils avisés m'ont aidé à avancer tout au long du noir cheminement du doctorat. Je remercie également Cathy CLERBAUX et Claudia STUBENRAUCH qui ont suivi mon travail tout au long de la thèse et m'ont apporté leur soutien. Je remercie mes voisins de bureau avec qui nous avons passé trois riches années : Frédéric, Walter et Victoria, ainsi que les divers intermittents qui se sont succédés. Je remercie également mes collègues du LERMA : Laurent, Éric, Jana, Carlos, Valérie, Annick, Vivianne... La bonne ambiance au sein du laboratoire assure un cadre de travail agréable.

Je remercie les enseignants du MBA du Collège des Ingénieurs qui, à travers des formations diverses, m'ont ouvert de nouveaux horizons et permis de me détacher par moment de mon travail de thèse. Ils m'ont ainsi fait prendre du recul sur mon travail afin de mieux l'approfondir.

Je tiens aussi à remercier ma compagne Chloé qui m'a soutenu au cours de ces trois ans. Je remercie toute ma famille qu'ils soient plus ou moins éloignés, ils ont tous à leur façon contribué à ce travail. Dans un autre domaine, je remercie mon chien Macaron dont la présence rassurante m'a porté tout au long de la route. Je remercie enfin le club de l'ASM Clermont Auvergne qui a montré année après année que l'on peut toujours grandir de ses échecs.



---

# SOMMAIRE

---

<b>I</b>	Le sondage atmosphérique	5
<b>2</b>	Restitution de la température de surface et de l'émissivité infrarouge	37
<b>3</b>	Chaîne opérationnelle de restitution de surface	83
<b>4</b>	Synergie : Principes généraux	107
<b>5</b>	Synergie au-dessus des océans, en ciel clair	121
<b>6</b>	Échantillonnage par entropie	155
<b>7</b>	Synergie au-dessus des continents, en ciel clair	183
<b>A</b>	Acronymes	207
<b>B</b>	Mathématiques	211
<b>C</b>	Re-analyses de l'ECMWF	219

---



# INTRODUCTION

La société contemporaine est toujours en quête de plus de rapidité, de précision et de fiabilité. Les exigences de nos systèmes d'information augmentent chaque jour. La technologie doit suivre l'évolution constante de la science et des besoins sociétaux toujours plus forts. Les domaines météorologique et climatique n'échappent pas à la règle, et la moindre erreur ou le moindre retard est parfois catastrophique (dans le cas d'évènements climatiques extrêmes, ou l'évolution du climat en général par exemple). L'apport du sondage atmosphérique par satellite est un élément de réponse face à ces exigences. Il permet de fournir, de façon continue, des données précises et globales sur l'état de l'atmosphère. Ces données sont nécessaires aux prévisions météorologiques mais aussi à un suivi des changements climatiques. La précision recherchée dans la météorologie ne sert pas seulement à fournir au public un indice sur la tenue vestimentaire à adopter ou sur la destination à privilégier pour partir en week-end. Ces données permettent également d'anticiper le mouvement des différentes tempêtes et d'évacuer à temps les zones à risques, de guider les secours lors de violents incendies, ou encore d'optimiser les pratiques agricoles. D'un point de vue climatique, le débat sur l'origine du réchauffement climatique est d'actualité. L'obtention d'une longue série de mesures globales et cohérentes, indépendantes des modèles de climat, permet de mieux étudier et comprendre les facteurs qui sont à l'origine de ce réchauffement.

Les nombreuses plates-formes satellites d'observation de l'atmosphère qui orbitent autour de la terre disposent, pour la plupart, de plusieurs capteurs sondant à différentes longueurs d'onde (visible, infrarouge, micro-onde). En dépit de la disponibilité de ces observations multi-fréquences, peu d'efforts ont été investis jusqu'à présent afin de concevoir des algorithmes de restitution qui exploiteraient au mieux les synergies potentielles. Plusieurs instruments peuvent être sensibles aux mêmes variables atmosphériques : les utiliser conjointement peut alors améliorer la précision de la restitution de cette variable. De telles méthodes sont déjà utilisées au sein des centres de prévision, mais nous cherchons à obtenir ici des données indépendantes des modèles.

La plupart du temps, les restitutions des différents instruments sont combinées *a posteriori*, ce qui peut introduire des incohérences et la synergie potentielle entre les observations n'est pas exploitée. Un algorithme synergique utilisera simultanément ou hiérarchiquement les observations de deux, ou plus, gammes spectrales dans le but d'obtenir une restitution plus précise que les restitutions indépendantes réunies. Comment utiliser efficacement différentes sources d'information pour un meilleur produit final? Nous verrons, dans cette étude, que l'algorithme d'inversion doit pouvoir fusionner intelligemment l'information afin d'optimiser leur utilisation. Les modèles d'assimilation numérique de prévision du temps sont capables d'utiliser simultanément des observations multi-spectrales. Toutefois, l'assimilation mélange les observations satellites avec un modèle de circulation générale, or il est

essentiel d'obtenir des bases de données atmosphériques qui soient purement observationnelles. De plus, il est difficile, avec cette approche, de réellement comprendre et d'évaluer le potentiel synergique des observations et donc d'identifier de nouvelles pistes d'amélioration des algorithmes actuels. Il est important de noter que si la synergie est exploitée au mieux, il existe encore un potentiel d'amélioration de la restitution des variables atmosphériques clés, comme les profils de température, de vapeur d'eau, ou d'ozone, sur lesquelles les centres de prévision se sont focalisés depuis la naissance de la météorologie. La synergie peut également avoir des avantages au delà des paramètres traditionnels, tels que la pluie, la couverture nuageuse, ou d'autres caractéristiques atmosphériques.

L'utilisation des données satellite à des fins de sondage atmosphérique souffre de nombreuses complications. Par exemple, le mélange des contributions de la surface et de l'atmosphère dans le signal mesuré par le satellite peut être un véritable problème. Ainsi, une méconnaissance de ces variables de surface peut entraîner une mauvaise caractérisation de l'atmosphère, et *vice versa*. Une longue étude sera donc présentée dans ce manuscrit afin de restituer la température de la surface et son émissivité aux différentes longueurs d'onde considérées. Nous nous concentrons, ici, sur l'infrarouge car une grande quantité de travail a déjà été accomplie sur les micro-ondes par les membres de notre équipe.

Cette étude est axée sur les données météorologiques et l'utilisation de la synergie entre les différentes longueurs d'onde pour restituer des variables atmosphériques. Une longue présentation du sondage atmosphérique permet aux néophytes de comprendre les tenants et les aboutissants de ce domaine de recherche et de pouvoir ainsi comprendre l'intérêt de la synergie. Les principes ainsi que les outils statistiques présentés sont généraux, ils peuvent être utilisés dans des contextes très différents (astronomie, écologie, biologie, éthologie et même économétrie ou finance). Le lecteur pourra retrouver, au cours de cette thèse, une présentation de différents outils statistiques ainsi que leurs avantages et inconvénients respectifs. Une partie de l'étude est notamment consacrée à une méthode d'échantillonnage de bases hétérogènes et de grande dimension dont l'intérêt est général et qui est facilement réutilisable pour d'autres applications.

L'utilisation de MetOp présentée dans cette étude est une application de grande envergure mais une méthodologie similaire à la notre peut être utilisée pour d'autres plates-formes avec d'autres instruments. L'avantage de cette plate-forme est que ses capteurs micro-ondes (AMSU-A et MHS) sont en phase avec le capteur infrarouge (IASI). Ainsi, ces capteurs sondent la même colonne d'atmosphère au même moment. Il est donc plus facile de combiner leurs mesures pour effectuer une restitution. Cependant, une méthodologie utilisant la synergie est envisageable même avec plusieurs plates-formes, dans la mesure où l'algorithme mis en place prend en compte les écarts temporels ou spatiaux entre les mesures.

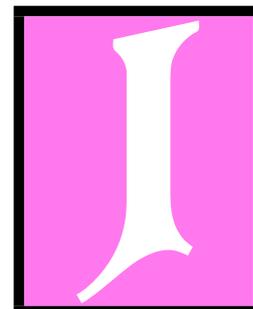
Cette thèse est composée de plusieurs volets. Dans une première partie, nous présenterons les bases du sondage atmosphérique, en particulier l'impact de la surface sur le rayonnement reçu par le satellite. Nous présenterons ensuite l'algorithme que nous avons mis en place

afin de restituer simultanément l'émissivité infrarouge et la température de surface, nous le validerons en l'exploitant opérationnellement sur une grande base de données provenant d'EUMETSAT. La suite de l'étude portera sur la synergie elle-même, avec une présentation des principes mathématiques ainsi que des différentes méthodes statistiques pouvant être utilisées. Cette partie a été séparée volontairement du reste du manuscrit afin de ne pas noyer les conclusions obtenues sous un flot d'équations qui pourraient sembler rébarbatives pour le lecteur non averti. Enfin, l'étude de la synergie en elle-même sera présentée, tout d'abord au-dessus des océans. Une présentation de l'échantillonnage mis en place pour créer des bases d'apprentissage sera effectuée, et pour conclure, nous étudierons la synergie infrarouge et micro-onde pour la restitution atmosphérique au-dessus des continents, en prenant en compte l'émissivité de la surface précédemment restituée.



---

# LE SONDAGE ATMOSPHERIQUE



## Sommaire

---

<b>1.1 Le transfert radiatif</b> . . . . .	<b>7</b>
1.1.1 Le rayonnement électromagnétique . . . . .	7
1.1.1.1 L'émission . . . . .	8
1.1.1.2 L'absorption . . . . .	9
1.1.1.3 La diffusion . . . . .	11
1.1.2 L'équation de transfert radiatif . . . . .	12
<b>1.2 Le sondage satellite : principe</b> . . . . .	<b>18</b>
1.2.1 Le spectre électromagnétique . . . . .	21
1.2.2 Le sondage dans le micro-onde . . . . .	24
1.2.3 Le sondage dans l'infrarouge . . . . .	25
<b>1.3 La plate-forme MetOp</b> . . . . .	<b>25</b>
1.3.1 IASI . . . . .	29
1.3.2 AMSU-A . . . . .	31
1.3.3 MHS . . . . .	32
<b>1.4 Le modèle de transfert radiatif RTTOV</b> . . . . .	<b>33</b>
<b>1.5 Les centres de prévision météorologique</b> . . . . .	<b>34</b>
1.5.1 Dénomination des données satellites . . . . .	34
1.5.2 Les centres de prévision numérique : NWP . . . . .	34

---

L'observation de l'atmosphère est utile à deux points de vue :

**En météorologie** pour les prévisions à plus ou moins court terme de l'état de l'atmosphère.

Une connaissance précise de l'état de l'atmosphère en temps réel et de façon globale permet de mieux contraindre les modèles de prévision et ainsi de mieux anticiper les événements météorologiques ;

**En climatologie** pour l'étude de l'évolution de l'atmosphère passée et future. Cela permet de prédire le devenir de l'atmosphère sur le long terme.

La météorologie a été étudiée dès l'antiquité. Cependant, c'est après l'invention du baromètre par Evangelista Torricelli au XVII<sup>ème</sup> siècle, qui cherchait à prouver l'existence du vide, contraire aux croyances religieuses de l'époque, que les mesures atmosphériques ont

réellement commencé. Le 19 septembre 1648, Blaise Pascal chargea son beau-frère Florin Périer d'effectuer des mesures à l'aide d'un tube rempli de mercure (le fameux baromètre de Torricelli) depuis la place de Jaude à Clermont-Ferrand jusqu'au sommet du Puy de Dôme (voir Figure 1.1). Ces mesures permirent d'établir l'existence de la pression atmosphérique et sa variation avec l'altitude. Des observations *in situ* de l'atmosphère, à l'aide de capteurs au sol ou de ballons sondes, ont été effectuées dès les années 1850 (Welsh 1853; Barral 1852), voire même dès 1804 avec les ascensions en ballon de Gay-Lussac (voir Figure 1.1). Cependant, ces mesures restent, même de nos jours, très limitées spatialement et sont dépendantes des différents instruments utilisés. Leur utilisation permet une détermination de l'atmosphère continue mais pas globale.



FIGURE 1.1 – Mesures météorologiques par Florin Périer et Gay-Lussac : Sur l'image de gauche, Florin Périer mesurant la pression atmosphérique au long de son ascension du Puy de Dôme le 19 septembre 1648 ; Sur l'image de droite Louis Joseph Gay-Lussac et Jean-Baptiste Biot effectuant des mesures en ballon le 20 septembre 1804. Source : Météo-France.

La couverture spatiale des mesures par satellite est nettement plus importante que les campagnes de mesures *in situ*. En effet, ces dernières sont limitées à un seul point (instruments au sol, ballons sondes...) voire en une zone (campagnes par avions) mais ne peuvent figurer la Terre dans sa globalité. Certains programmes (par exemple le “global drifter program” qui recueille les données des bouées météorologiques ou le GOS (Global Observing System) qui réunit des mesures au sol et par satellite) rassemblent plusieurs sources d'informations pour avoir une vue globale, mais là encore cela reste dépendant des erreurs liées à chacun des instruments de mesure (ici chaque bouée ne fait pas la même erreur, il est donc difficile d'intégrer tous les signaux). Une couverture plus complète est obtenue grâce aux mesures par avion, néanmoins, ces campagnes ne couvrent pas la Terre dans sa globalité. L'utilisation des satellites permet d'avoir une couverture globale du globe, notamment dans

les régions peu instrumentées car peu accessibles.

Les problèmes liés à la calibration des instruments sont des problèmes auxquels font face toutes les mesures (par satellite ou *in situ*). Cependant, les centres de prévision météorologique parviennent à les exploiter et à s'affranchir de ces problèmes. L'avantage majeur, que présentent les mesures par satellite, est leur couverture spatiale. Cela permet d'intégrer de plus en plus de données dans les centres météorologiques, surtout pour l'hémisphère sud nettement moins équipé pour des observations *in situ*. Il est intéressant de pouvoir disposer de mesures par satellites pour mieux caractériser l'atmosphère. Ces données sont incorporées en direct dans les centres de prévision opérationnels afin d'affiner les modèles.

D'un point de vue climatologique, l'apport des données *in situ* est plus important car la durée de vie d'un instrument au sol est plus élevée. De plus, le sondage atmosphérique par satellite étant encore assez "récent", il n'existe pas encore de longue série temporelle. Cependant, le développement de la restitution des propriétés atmosphériques à partir des mesures satellites permettra dans un avenir proche de créer des bases de données globales sur de longues séries temporelles<sup>1</sup>. La plupart des missions spatiales actuelles prévoient une succession de satellites afin de palier l'usure des instruments dans l'espace.

Ainsi, le sondage atmosphérique par satellite révolutionne les prévisions opérationnelles actuelles qui intègrent de plus en plus ces nouvelles données. A l'heure actuelle, près de 90% des données utilisées dans les modèles des centres de prévision météorologiques viennent des satellites.

Ce chapitre détaille le principe du sondage atmosphérique par satellite, décrit les instruments qui seront utilisés dans cette étude, l'utilisation qui est faite de ces données par les différents acteurs (météorologues, climatologues...) ainsi que le modèle de transfert radiatif RTTOV qui sera utilisé pour simuler les différents instruments.

## 1.1 Le transfert radiatif

### 1.1.1 Le rayonnement électromagnétique

Une molécule peut interagir de trois façons différentes avec le rayonnement électromagnétique : absorber ou diffuser le rayonnement incident, ou émettre son propre rayonnement. Chacune des molécules de l'atmosphère ou de la surface terrestre interagit avec le rayonnement suivant ces trois modes. C'est grâce à cette interaction entre l'atmosphère et le rayonnement électromagnétique que l'observation satellite peut être utilisée pour caractériser l'atmosphère. Les sections suivantes détaillent ces différentes interactions.

---

1. De telles bases de données ont déjà été rassemblées, notamment pour SSM/I qui vole depuis près de trente ans (on peut également citer les satellites de la NOAA ou Meteosat). Cependant, il est nécessaire de continuer à développer de telles bases car les plus étendues ne couvrent que quelques dizaines années.

### 1.1.1.1 L'émission

Un corps noir est un corps qui absorbe la totalité du rayonnement incident et émet à son tour un rayonnement de sorte à équilibrer son bilan radiatif (*i.e.*, l'énergie incidente est la même que celle émise). Le rayonnement ainsi émis suit une distribution spectrale conforme à la loi de Planck :

$$B_{\lambda}(T) = \frac{2hc^2}{\lambda^5} \cdot \frac{1}{e^{\left(\frac{hc}{\lambda kT}\right)} - 1}$$

où  $\lambda$  est la longueur d'onde du rayonnement considéré,  $h$  est la constante de Planck ( $6,62617 \cdot 10^{-34}$  J·s),  $k$  est la constante de Boltzmann ( $1,38066 \cdot 10^{-23}$  J·K<sup>-1</sup>),  $c$  est la vitesse de la lumière ( $299792458$  m·s<sup>-1</sup>) et  $T$  est la température absolue du corps.

Le Soleil peut être considéré en première approximation comme un corps noir à une température de 5800 K auquel il faut ajouter des raies d'absorption dues aux différents éléments chimiques qui composent son atmosphère. Ce rayonnement intervient, en majeure partie, de l'ultra-violet au proche infrarouge. Le maximum d'émission se situe dans le visible. Le rayonnement terrestre quant à lui est semblable à celui d'un corps noir à 288 K. Il est maximum dans l'infrarouge thermique. Les spectres d'émission du Soleil et de la Terre sont présentés sur la Figure 1.2. L'ordonnée (en représentation logarithmique) représente la puissance du rayonnement par mètres carrés par unité d'angle solide et par longueur d'onde. En abscisse, la longueur d'onde en  $\mu\text{m}$  est indiquée en bas et le nombre d'onde en  $\text{cm}^{-1}$  en haut. Les valeurs indiquées sont des valeurs qui seront reprises par la suite : la longueur d'onde du rayonnement visible (0,3 à 0,7  $\mu\text{m}$ ) et les canaux extrêmes de IASI (entre 645 et 2760  $\text{cm}^{-1}$ , voir Section 1.3.1, page 29).

Sur cette figure on peut remarquer que le rayonnement solaire est nettement supérieur à celui de la Terre. C'est pourquoi dans l'espace, les satellites (en orbite au sommet de l'atmosphère) doivent être conçus pour résister au rayonnement solaire très fort et plus contraignant que celui émis par la Terre. Du point de vue d'un satellite visant la Terre, le rayonnement solaire traverse d'abord l'atmosphère, se réfléchit sur la surface de la Terre et retransverse l'atmosphère avant d'atteindre le satellite. Si le satellite est sur le cheminement de ce rayonnement (c'est-à-dire exactement dans l'axe du reflet du Soleil sur la Terre, appelé "sunglint" en anglais), le rayonnement mesuré est faussé par le rayonnement solaire (même si ce dernier est en partie absorbé lors de sa traversée de l'atmosphère). Ces situations sont donc systématiquement filtrées pour éviter des erreurs. Elles sont suffisamment rares pour être éliminées. La réflexion du Soleil sur la Terre est très directionnelle. Si le satellite n'est pas dans la direction du reflet, il ne recevra quasiment pas de rayonnement solaire. Le rayonnement solaire sera donc négligé par la suite<sup>2</sup>.

---

2. La bande 3 de IASI est particulièrement sensible au rayonnement solaire, elle ne sera donc pas prise en compte, voir Section 1.3.1, page 29

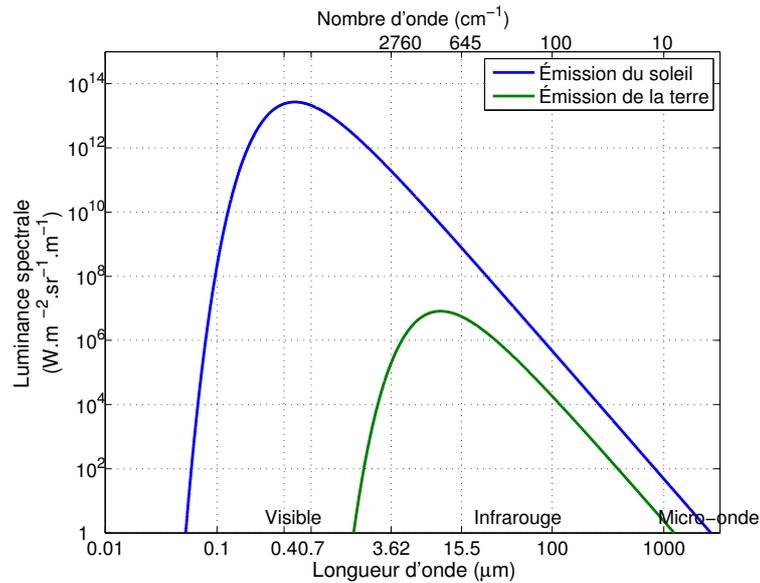


FIGURE 1.2 – Spectre d'émission du Soleil et de la Terre considérés ici comme deux corps noirs à respectivement 5800 K et 288 K.

### 1.1.1.2 L'absorption

Chaque molécule peut se trouver sous différents états électroniques. Une molécule peut quitter son état stable, appelé fondamental, si elle reçoit une énergie suffisante. Chaque état est quantifié par un niveau d'énergie. La molécule peut atteindre différents états d'excitation suivant la quantité d'énergie qu'elle reçoit. Chaque état est à son tour distingué en plusieurs états de transition, appelés états de vibration ou de rotation. Leur dénomination décrit directement le mouvement de la molécule légèrement excitée. Une quantité précise d'énergie est nécessaire à la molécule pour passer d'un état à un autre. La Figure 1.3 présente les différents niveaux d'excitation d'une molécule. On voit que l'état avec le moins d'énergie est l'état fondamental. Un faible apport d'énergie peut faire changer l'état de cette molécule jusque vers les états de transition (rotations, vibrations...). L'état électronique excité (ici un seul est représenté mais il peut y en avoir plusieurs) nécessite plus d'énergie pour être atteint. Chacun de ces "paliers" représente une quantité d'énergie fixe pour chaque molécule. Les molécules ont toutes leurs propres modes. Leurs vibrations ou leurs rotations dépendent de leur structure. Elles peuvent tourner ou vibrer selon différents axes. Elles peuvent également avoir différents nombres de niveaux électroniques.

L'énergie contenue dans le rayonnement électromagnétique peut entraîner ces processus d'excitation, de dissociation ou d'ionisation au sein des molécules qu'il rencontre. Cet échange d'énergie se traduit par une absorption de photon. Le processus qui est engendré dépend de l'énergie du photon incident. L'énergie  $E$  d'un photon se propageant à la fréquence  $f$  vaut  $E = h \times f$ , où  $h$  est la constante de Planck. Le rayonnement dans le micro-onde contient donc une faible énergie, il ne peut qu'engendrer des rotations de la molécule. Dans

l'infrarouge, l'énergie transportée est plus importante, le rayonnement incident peut donc engendrer des vibrations en plus des rotations. Les rayonnements dans le domaine visible et ultra-violet, non traités dans cette étude, transportent plus encore d'énergie et peuvent engendrer des changements de l'état électronique des molécules, voire des dissociations ou des ionisations.

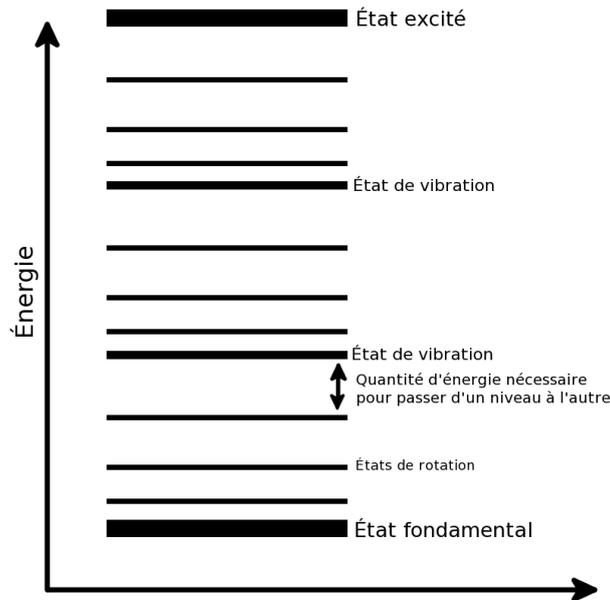


FIGURE 1.3 – Schéma des différents niveaux d'énergie pour une molécule donnée. Un seul niveau électronique supplémentaire est représenté mais il peut y en avoir plus. De même, le nombre de niveaux de rotations et de vibrations sont donnés à titre d'exemple.

Chacun des différents états de la molécule (vibration, rotation...) correspond à un niveau d'énergie donné. Il faut que le photon incident contienne exactement l'énergie entre deux états de la molécule pour que celle-ci l'absorbe et change d'état. Ainsi seuls certains photons munis de la bonne quantité d'énergie seront absorbés. Ceci explique que chaque molécule absorbe le rayonnement dans différentes zones spectrales (uniquement aux longueurs d'onde qui correspondent à l'énergie entre leurs modes).

Une molécule se trouvant dans un état excité peut à son tour émettre un rayonnement afin de revenir à un état plus stable. La température de l'atmosphère est telle que les molécules se trouvent majoritairement dans leur état fondamental. Seuls les molécules situées au sommet de l'atmosphère peuvent être dans des états électroniques excités. Ces molécules sont soumises à un rayonnement solaire plus important et la faible densité de l'atmosphère implique aussi que les chocs entre molécules sont plus rares et il y a donc moins d'échanges d'énergie entre les molécules. La transition de ces molécules d'un état excité à l'état fondamental entraîne une importante émission de photons dans l'ultra-violet ou le visible. Ces

émissions sont à l'origine des aurores boréales et à la luminescence du ciel nocturne<sup>3</sup>.

Dans le reste de l'atmosphère, la température et le plus grand nombre de chocs impliquent que les molécules soient dans leur état fondamental. Elles passent simplement de différents niveaux de vibration et/ou de rotation à l'état fondamental. Ces transitions impliquent l'émission d'un rayonnement de faible énergie (*i.e.*, dans le domaine micro-onde ou infrarouge). L'atmosphère émet donc un rayonnement pour équilibrer l'énergie absorbée. D'après la loi de Kirchhoff, les coefficients d'absorption et d'émission sont égaux lorsque l'équilibre thermodynamique local est atteint (approximation faite dans la Section 1.1.2, page 12). Le rayonnement qui est absorbé par les molécules est entièrement réémis. L'émission aura donc lieu dans les mêmes bandes spectrales que l'absorption. Cependant, les molécules n'émettent dans la même direction que le rayonnement incident. On est alors en présence de diffusion présentée dans la section qui suit.

### 1.1.1.3 La diffusion

On parle de diffusion lorsque le rayonnement est dévié dans diverses directions par des particules. Si la diffusion est équirépartie entre toutes les directions, on parle de diffusion isotrope. À l'inverse, si le rayonnement est dévié de façon orientée, on parle de diffusion anisotrope ou directionnelle. La partie de l'onde qui est diffusée dans le sens contraire à l'onde incidente est appelée rétrodiffusée (voir la Figure 1.4 pour mieux comprendre ces notions). On distingue encore deux types de diffusion : la diffusion élastique où il n'y a pas de perte d'énergie (*i.e.*, la longueur d'onde du rayonnement diffusé est la même que celle du rayonnement incident) et la diffusion inélastique avec changement d'énergie et donc de longueur d'onde du rayonnement.

La diffusion de Rayleigh correspond à la diffusion d'une onde par une particule largement inférieure en taille à la longueur d'onde du rayonnement incident. Cela correspond au phénomène d'absorption et d'émission décrit précédemment. C'est donc une diffusion élastique car sans changement de longueur d'onde. De façon nettement plus faible en intensité, la molécule diffusante peut émettre un rayonnement légèrement différent en longueur d'onde. On parle alors de diffusion Raman. Celle-ci n'est pas traitée dans cette étude.

La diffusion de Rayleigh est une diffusion considérée comme isotrope (les particules émettent dans toutes les directions de façon égale, voir Figure 1.4). L'intensité lumineuse diffusée est proportionnelle à l'inverse de la longueur d'onde puissance quatre ( $\frac{1}{\lambda^4}$ ). Un rayonnement dans le bleu sera plus diffusé (par les fines particules) qu'un rayonnement dans le rouge car  $\lambda_{bleu} < \lambda_{rouge}$ . Comme présenté sur la Figure 1.2, le Soleil émet un rayonnement dans tout le domaine visible. Le ciel paraît bleu, car la partie bleue du rayonnement solaire est plus diffusée par l'atmosphère.

---

3. Le meilleur exemple d'émission de photons est le feu d'artifice : les molécules excitées par la chaleur de l'explosion changent d'état électronique. En revenant à leur état fondamental elles émettent un rayonnement à différentes longueurs d'onde suivant les molécules, qui peuvent correspondre à différentes couleurs dans le visible.

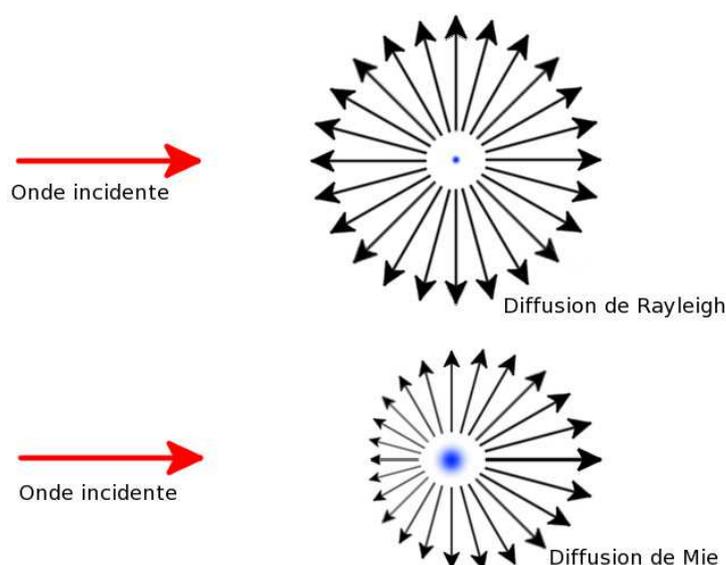


FIGURE 1.4 – Schéma des différents types de diffusion. Le rayonnement arrive de la gauche. La molécule du dessus est de taille très inférieure à la longueur d’onde du rayonnement, c’est la diffusion de Rayleigh. La molécule au-dessous est de taille comparable à la longueur d’onde du rayonnement incident, c’est la diffusion directionnelle de Mie.

Lorsque les particules diffusantes sont de taille comparable à la longueur d’onde, c’est alors la diffusion de Mie (Mie 1908; van de Hulst 1981). L’intensité du rayonnement diffusé varie alors, de façon inversement proportionnelle à la longueur d’onde ( $\frac{1}{\lambda}$ ). Plus la taille des particules augmente, plus la diffusion sera directionnelle et vers l’avant. On peut remarquer la directionnalité de cette diffusion sur la Figure 1.4. Cette diffusion a lieu, entre autres, au niveau des nuages. Le fait qu’elle varie nettement moins en fonction de la longueur d’onde du rayonnement fait que, dans le visible, toutes les longueurs d’onde sont diffusées de la même manière. On voit donc les nuages blancs (mélange de tous les rayonnements aux différentes couleurs du Soleil). La directionnalité de cette diffusion explique les dégradés de gris de certains nuages car dans certaines directions, le rayonnement est moins diffus et donc moins fort<sup>4</sup>. On le voit alors moins brillant et plus terne.

### 1.1.2 L’équation de transfert radiatif

Comme présenté précédemment, l’atmosphère terrestre peut absorber, émettre ou diffuser le rayonnement, mais ses caractéristiques sont hétérogènes. Il est donc nécessaire d’étudier précisément la propagation d’une onde en son sein. Celle-ci dépend de la longueur d’onde considérée  $\lambda$ , de la température de l’atmosphère  $T$  et de la concentration des différents constituants de l’atmosphère.

La loi de Beer-Lambert définit un coefficient d’extinction  $K_\lambda = k_\lambda + \sigma_\lambda$ , avec  $k_\lambda$  le coeffi-

4. Leur couleur est également liée à leur opacité.

cient d'absorption et  $\sigma_\lambda$  le coefficient de diffusion (incluant toutes les différentes diffusions), tel que la variation  $dI_\lambda$  de la luminance  $I_\lambda$  (*i.e.*, intensité spectrale) qui traverse une couche atmosphérique infinitésimale de longueur  $dl$  vaut (voir Schéma 1.5) :

$$dI_\lambda = -K_\lambda \cdot I_\lambda \cdot dl$$

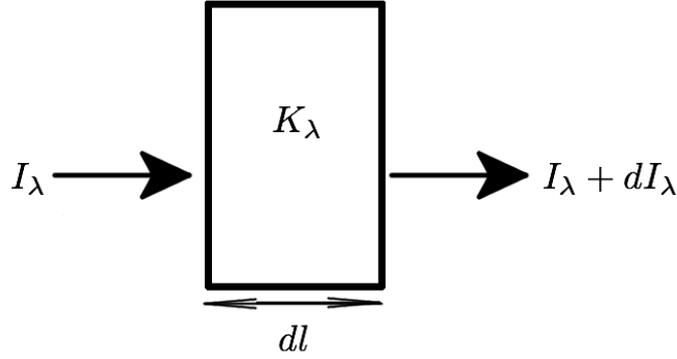


FIGURE 1.5 – Schéma de la variation de la luminance spectrale traversant une couche simple de longueur  $dl$ .  $I_\lambda$  est la luminance incidente,  $I_\lambda + dI_\lambda$  est la luminance qui en résulte.  $K_\lambda$  est le coefficient d'extinction du milieu traversé.

On a donc en intégrant sur l'épaisseur de l'atmosphère traversée depuis la surface jusqu'à l'altitude  $z$ , la luminance  $I_\lambda(z)$  qui vaut :

$$I_\lambda(z) = I_\lambda(0) \cdot e^{-\int_0^z K_\lambda dl}$$

On définit alors le facteur de transmission atmosphérique  $\tau = e^{-\int_0^z K_\lambda dl}$ . Cette relation s'appelle la loi de Beer-Lambert.

Dans le cas de l'atmosphère, il faut prendre en compte l'émission propre de l'atmosphère traversée :  $J_\lambda$  appelé fonction source. On a alors la variation de la luminance en un point  $M$  se propageant dans la direction  $\vec{s}$  :

$$dI_\lambda(M, \vec{s}) = -K_\lambda(M) \cdot (I_\lambda(M, \vec{s}) - J_\lambda(M, \vec{s})) \cdot d\vec{s} \quad (1.1)$$

Afin de simplifier cette équation on fait plusieurs hypothèses sur l'état de l'atmosphère :

- L'atmosphère est plan-parallèle, c'est-à-dire que le rayon de courbure de la Terre est très supérieur à l'épaisseur de l'atmosphère. On peut alors considérer l'atmosphère en un point donné comme une succession de couches horizontales ;
- L'atmosphère est à l'équilibre thermodynamique local. On peut donc appliquer localement la loi de Kirchhoff, c'est-à-dire que les coefficients locaux d'émission et d'absorption sont égaux.

Du point de vue d'un satellite, si on considère uniquement le rayonnement montant vertical, la luminance spectrale reçue par le satellite  $I_{\lambda}^{sat}$  peut être exprimée en fonction de la luminance spectrale perçue à une altitude  $z$ , d'après la loi de Beer-Lambert présentée précédemment.

S'il n'y a pas de source entre le niveau de l'atmosphère d'altitude  $z$  et le satellite, on a, en caractérisant la luminance par son facteur de transmission  $\tau$  :

$$I_{\lambda}^{sat} = I_{\lambda}(\tau) \cdot \tau$$

Dans le cas d'un satellite, le facteur de transmission  $\tau$  se calcule entre l'altitude  $z$  d'où part le rayonnement et le sommet de l'atmosphère (voir Figure 1.6). Comme indiqué sur le schéma, une altitude peut être référencée par son altitude  $z$  ou par son facteur de transmission  $\tau$ . Dans la suite nous utiliserons le facteur de transmission pour référencer les altitudes afin de simplifier l'équation finale.

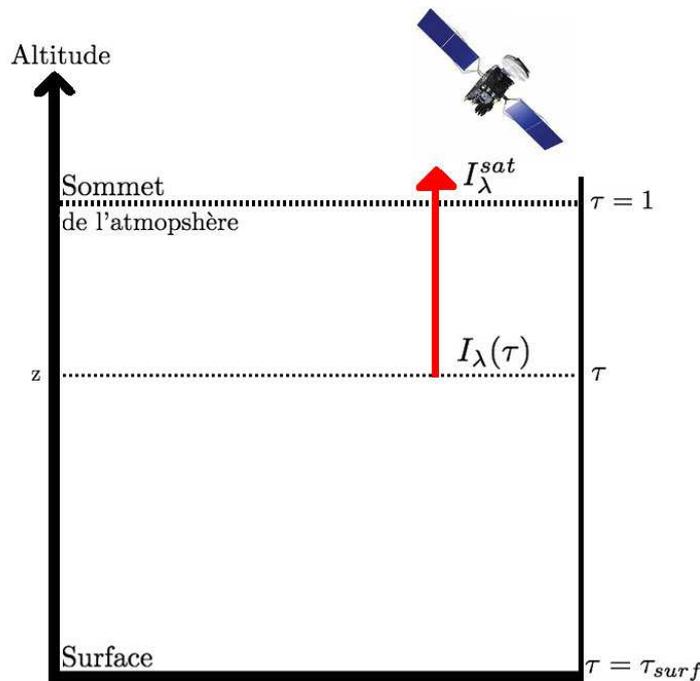


FIGURE 1.6 – Schéma de la luminance spectrale reçue par le satellite en fonction de la luminance spectrale perçue à une altitude  $z$ .

On a donc en différenciant l'équation précédente :

$$\tau \cdot dI_{\lambda}(\tau) + d\tau \cdot I_{\lambda}(\tau) = 0$$

Soit :

$$\frac{dI_{\lambda}}{I_{\lambda}} = -\frac{1}{\tau} d\tau \quad (1.2)$$

On considère maintenant une couche de l'atmosphère située entre les niveaux  $\tau$  et  $\tau + d\tau$ . La loi de Kirchhoff appliquée localement et l'équation (1.2) nous permettent de dire que cette couche de l'atmosphère émet :

$$\frac{1}{\tau} \cdot I(\tau) \cdot d\tau$$

Soit :

$$\frac{1}{\tau} \cdot B_\lambda (T(\tau)) \cdot d\tau$$

où  $T(\tau)$  est la température de la couche considérée et  $B_\lambda$  est la fonction d'émission de Planck à la longueur d'onde  $\lambda$ . En faisant le bilan d'énergie de cette couche on obtient :

$$dI_\lambda(\tau) = -\frac{1}{\tau} (I_\lambda - B_\lambda (T(\tau))) \cdot d\tau$$

on retrouve bien l'équation (1.1) où  $B_\lambda (T(\tau))$  correspond au terme source.

En intégrant cette équation entre la surface et un niveau de facteur de transmission  $\tau$ , on obtient :

$$I_\lambda(\tau) = \frac{1}{\tau} \cdot \int_{\tau_\lambda^{surf}}^{\tau} B_\lambda (T(\tau')) d\tau' + \frac{1}{\tau} \cdot A$$

où  $A$  est une constante à déterminer.

Si on se place au niveau de la surface on a :

$$I_\lambda^{surf}(\tau) = \varepsilon_\lambda^{surf} B_\lambda (T^{surf}) = \frac{1}{\tau_\lambda^{surf}} \cdot A$$

Au niveau du satellite, le facteur de transmission atmosphérique  $\tau$  vaut 1, on a alors la luminance perçue par le satellite qui vaut :

$$I_\lambda^{sat}(\tau) = \tau_\lambda^{surf} \varepsilon_\lambda^{surf} B_\lambda (T^{surf}) + \int_{\tau_\lambda^{surf}}^1 B_\lambda (T(\tau')) d\tau'$$

Dans cette équation, on a considéré uniquement le rayonnement montant reçu par le satellite, mais il faut également prendre en compte le rayonnement descendant émis par l'atmosphère, réfléchi par la surface puis qui atteint le satellite. Un raisonnement similaire nous conduit au terme :

$$(1 - \varepsilon_\lambda^{surf}) \cdot \tau_\lambda^{surf} \cdot \int_1^{\tau_\lambda^{surf}} B_\lambda (T(\tau')) d\tau'$$

Le terme  $(1 - \varepsilon_\lambda^{surf})$  correspond au coefficient de réflexion de la surface.

L'équation de transfert radiatif finale est donc constituée de trois termes :

$$I_\lambda^{sat}(\tau) = \tau_\lambda^{surf} \varepsilon_\lambda^{surf} B_\lambda (T^{surf}) + \int_{\tau_\lambda^{surf}}^1 B_\lambda (T(\tau')) d\tau' + (1 - \varepsilon_\lambda^{surf}) \cdot \tau_\lambda^{surf} \cdot \int_1^{\tau_\lambda^{surf}} B_\lambda (T(\tau')) d\tau'$$

Le premier terme correspond au rayonnement émis par la surface, le deuxième correspond au rayonnement montant émis par chaque couche de l'atmosphère et le troisième correspond au rayonnement descendant émis par chaque couche de l'atmosphère puis réfléchi par la surface (voir Schéma 1.7).

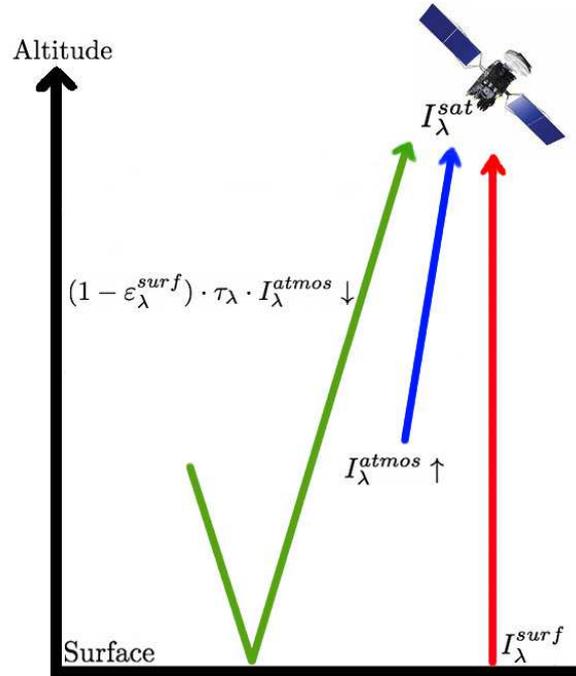


FIGURE 1.7 – Schéma des différents rayonnements reçus par le satellite au sommet de l'atmosphère : l'émission de la surface, le rayonnement montant émis par l'atmosphère et le rayonnement descendant émis par l'atmosphère puis réfléchi par la surface.

Deux cas de figures particuliers peuvent être mis en évidence. Si l'atmosphère est transparente à la longueur d'onde considérée (région "fenêtre"),  $\tau_\lambda = 1$  pour chaque couche, la majeure partie du rayonnement mesuré au sommet de l'atmosphère provient de la surface :

$$I_\lambda^{sat}(\tau) = \tau_\lambda^{surf} \varepsilon_\lambda^{surf} B_\lambda(T^{surf}(\tau))$$

L'intégrale est alors calculée entre deux bornes très proches et le terme à intégrer est fini. Les deux intégrales sont donc nulles. Dans ce cas, la luminance mesurée par le satellite au sommet de l'atmosphère est une fonction de l'émissivité de la surface et de la température de surface.

Si l'atmosphère est très absorbante,  $\tau_\lambda^{surf} = 0$ , la contribution de la surface est faible. L'équation de transfert radiatif se résume donc à :

$$I_\lambda^{sat}(\tau) = \int_{\tau_\lambda^{surf}}^1 B_\lambda(T(\tau')) d\tau'$$

Soit :

$$I_{\lambda}^{sat}(z) = \int_{Z_{surf}}^{Z_{sommet}} B_{\lambda}(T(z')) \frac{\partial \tau_{\lambda}}{\partial z'} dz'$$

Le terme  $\frac{\partial \tau_{\lambda}}{\partial z'}$  qui apparaît dans l'équation est important, il s'appelle la fonction de poids. Comme la transmission atmosphérique peut s'exprimer comme l'exponentielle d'une intégrale en  $z$  du coefficient d'extinction, ce facteur peut s'exprimer comme un produit du coefficient d'extinction par la transmission atmosphérique :

$$\frac{\partial \tau_{\lambda}}{\partial z'} = K_{\lambda} \times e^{-\int_0^z K_{\lambda} dt}$$

Le coefficient d'extinction est une fonction croissante du nombre de molécule par volume d'air et donc décroît avec l'altitude. La transmission atmosphérique croît de 0, à la surface, à 1, au sommet de l'atmosphère. La fonction de poids est donc une multiplication entre un terme qui tend vers 0 au sommet de l'atmosphère ( $K_{\lambda}$ ) et un autre croissant de 0 à 1 ( $\tau_{\lambda}$ ). Elle vaut 0 au sommet de l'atmosphère ainsi qu'à la surface de la Terre et présente un pic à une altitude donnée, qui est une fonction croissante de l'absorption atmosphérique (voir Figure 1.8). Cette fonction de poids permet de savoir quelle longueur d'onde sonde quelle altitude de l'atmosphère<sup>5</sup>.

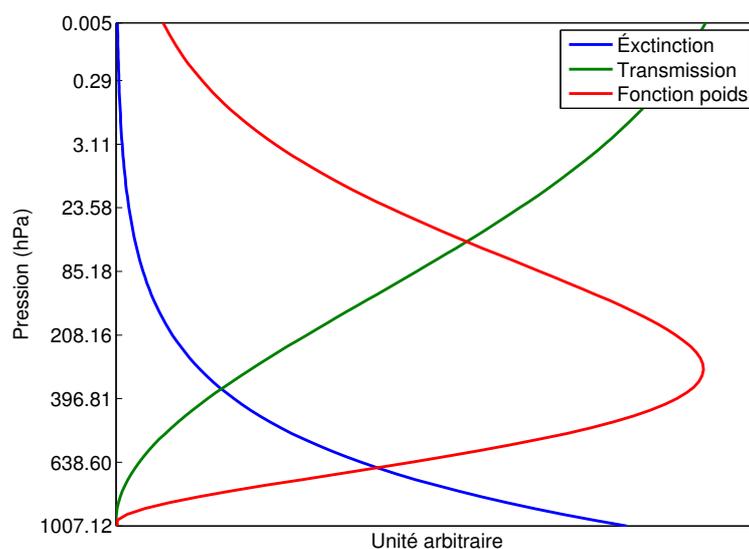


FIGURE 1.8 – Sur le graphique, la variation du coefficient d'extinction atmosphérique est représentée en bleu. La transmission atmosphérique est en vert, et la fonction de poids résultante est représentée en rouge. Ces courbes sont données à titre d'exemple, elles ne représentent pas une longueur d'onde particulière.

L'absorption atmosphérique varie en fonction de la longueur d'onde considérée. Il est

5. Au niveau des canaux fenêtres, c'est-à-dire que la transmission atmosphérique est élevée jusqu'à la surface, le maximum de la fonction de poids peut se trouver au niveau de la surface. Sonder à de telles longueurs d'onde permet d'avoir accès aux informations de surface.

donc possible de restituer des informations sur la température et la composition atmosphérique tout le long du profil. Si l'absorbeur est uniformément réparti avec une concentration connue ( $\frac{\partial \tau_\lambda}{\partial z}$  est connu), le profil de température peut être restitué. Inversement, si le profil de température est connu, la seule inconnue est la concentration de l'absorbeur. En pratique, les profils de température et de concentration de l'absorbeur sont inconnues, ce qui rend la restitution compliquée.

L'ozone en haute altitude absorbe le rayonnement dans l'ultra-violet, le gaz carbonique et surtout la vapeur d'eau absorbent le rayonnement dans l'infrarouge et le micro-onde. Les nombreuses poussières en suspension dans l'atmosphère contribuent également à absorber et diffuser certains rayonnements. Il faut également tenir compte de la température de chaque couche de l'atmosphère. Le problème est donc complexe bien que bien identifié. Cependant, une utilisation simultanée de plusieurs mesures sensibles tant à la température qu'à la composition atmosphérique et des mesures plus sensibles à la température seule peut permettre de dissocier les deux informations et donc de les restituer.

Afin d'avoir des données plus facilement interprétables et dans une gamme de variations comparables d'une longueur d'onde à une autre, les mesures de radiances par satellites sont souvent converties en températures de brillance. En effet, l'énergie contenue dans le rayonnement électromagnétique varie beaucoup d'un domaine de longueur d'onde à un autre, ce qui complique son interprétation. La température de brillance correspond à la température qu'aurait un corps noir qui émettrait cette radiation à cette longueur d'onde.

## 1.2 Le sondage satellite : principe

Il existe deux types majeurs d'observation de l'atmosphère, que ce soit au sommet de l'atmosphère pour un satellite, au sein même de l'atmosphère pour un instrument à bord d'un avion, d'un ballon sonde, ou en surface :

- **Le sondage passif** consiste à mesurer la luminance spectrale reçue en un point. Les sources qui émettent sont alors les différents composants de l'atmosphère, la surface de la Terre et le Soleil ;
- **Le sondage actif** consiste à émettre un rayonnement et à mesurer la luminance spectrale perçue en retour. La source d'émission est alors imposée par le sondeur, les autres sources (émission des molécules de l'atmosphère) sont dans ce cas négligeables par rapport à l'intensité de l'émission du sondeur. La mesure porte uniquement sur la diffusion et la réflexion du rayonnement par les particules (aérosols, hydrométéores).<sup>6</sup>

La télédétection passive requiert moins de puissance que la télédétection active, c'est pourquoi elle est plus souvent utilisée à bord des satellites. En effet, plus le satellite a besoin d'énergie pour faire fonctionner ses différents instruments, plus il devra transporter

---

6. Il existe aussi le sondage GPS qui repose sur la l'occultation radio, ce type de mesures n'est pas traitée dans cette étude.

de batteries, lourdes et encombrantes. Les panneaux solaires qu'il transporte lui permette de récupérer de l'énergie, mais de façon limitée. Cette énergie est utilisée pour faire les différentes mesures grâce à ses capteurs, pour transmettre les données au sol. Les panneaux solaires ne fournissent pas assez d'énergie pour permettre aux instruments actifs d'avoir une très longue durée de vie dans l'espace. Les instruments actifs (LIDAR, radars) ne seront pas étudiés ici.

La contrainte énergétique principale des satellites reste cependant la quantité de carburant dont ils disposent. Ce carburant est nécessaire pour se maintenir sur son orbite car les frottements, même infimes, tendent à le faire dévier de sa trajectoire, de même que la non-sphéricité de la Terre et les différentes anomalies géodésiques. Si la position du satellite n'est pas connue de façon extrêmement précise, il est très compliqué de déterminer la zone visée par ses différents capteurs. La géolocalisation très précise des zones sondées est primordiale pour pouvoir exploiter les mesures. La quantité de carburant à bord du satellite est donc une contrainte vitale puisqu'elle conditionne leur durée de vie.

Cette étude est centrée sur la télédétection passive par satellite. Il s'agit de mesurer la luminance spectrale au sommet de l'atmosphère. Comme expliqué dans la section précédente, cette luminance spectrale résulte de l'interaction de l'atmosphère et de la surface sur le rayonnement. Les différents composants de l'atmosphère (vapeur d'eau, ozone...) influencent le rayonnement électromagnétique à différentes longueurs d'onde. Il est donc important de disposer de mesures de la luminance spectrale à différentes longueurs d'onde afin de pouvoir mesurer ces interactions et ainsi caractériser les différents composants de l'atmosphère.

Il existe trois principales méthodes pour déterminer, à partir des mesures passives par satellite, la composition de l'atmosphère et sa température ainsi que celle de la surface :

- Des méthodes physiques qui utilisent le transfert radiatif pour vérifier la justesse de la caractérisation de l'atmosphère. Ces méthodes fonctionnent souvent par itérations qui utilisent un état supposé de l'atmosphère pour simuler le transfert radiatif puis modifient l'état de l'atmosphère afin de coller au mieux avec les mesures satellites réelles (voir Section 1.5.2, page 34) ;
- Des méthodes statistiques qui utilisent des données antérieures pour paramétrer des régressions ou des réseaux de neurones (Aires et al. 2001, 2011a). Des tables de correspondance (look-up tables en anglais) peuvent également être utilisées pour associer rapidement un sondage satellite à une situation atmosphérique ;
- Des méthodes physico-statistiques basées sur des inversions mathématiques de l'équation de transfert radiatif (Paul et al. 2012; Li et al. 2007).

Cette étude utilise des méthodes physico-statistiques (voir Section 2.3, page 60) ainsi que des méthodes statistiques (voir Section 5.2, page 134, ou 7.3, page 194). L'objectif est de montrer l'existence d'une synergie entre les mesures satellites à différentes longueurs d'onde. En utilisant ces méthodes d'inversion appliquées à différentes mesures (ici micro-onde et infrarouge), des profils atmosphériques seront restitués. La finalité est de montrer

qu'il est préférable d'utiliser toutes les mesures conjointement, pour la restitution, plutôt que de combiner les résultats *a posteriori* (Aires et al. 2012).

Nous nous intéresserons plus particulièrement aux profils atmosphériques de température et de vapeur d'eau. Ces variables atmosphériques sont très importantes et dictent en partie le comportement et l'évolution de l'atmosphère. Elles sont facilement compréhensibles et influent sur la vie quotidienne du lecteur qui comprend alors la problématique mise en place ici. Il sera facile de mettre en avant la synergie car les deux variables sont influentes sur une large gamme de longueur d'ondes et de nombreux instruments satellites dans différentes gammes de longueurs d'onde y sont sensibles. La température et la vapeur d'eau atmosphérique sont des variables majeures dans la caractérisation de l'atmosphère, que ce soit pour les modèles de prévisions ou pour le suivi de l'évolution de l'atmosphère.

La température atmosphérique utilisée sera en Kelvin ( $0^{\circ}\text{C}=273,15^{\circ}\text{K}$ ). La vapeur d'eau, quant à elle, est généralement utilisée en  $\text{kg.kg}^{-1}$  (masse de vapeur d'eau par masse d'air humide). Cette unité de mesure est utile pour les calculs de transfert radiatif, cependant la valeur est moins évidente à comprendre. De plus, la variation de la température modifie la quantité de vapeur d'eau qui peut être contenue dans l'atmosphère. C'est pourquoi, nous utiliserons la plupart du temps une mesure de la vapeur d'eau non pas en humidité spécifique (notée  $q$  en  $\text{kg.kg}^{-1}$  ou en ppmv) mais en humidité relative (notée  $rh$  en %). L'humidité relative correspond au rapport entre la pression partielle en vapeur d'eau et la pression de vapeur saturante, qui est la pression d'équilibre de l'eau entre son état liquide et gazeux (ici la pression partielle en vapeur d'eau à laquelle elle redeviendrait liquide). Ainsi, une humidité relative de 100 % impliquera un phénomène de condensation et donc la formation de nuages voire de précipitations. Ponctuellement, l'humidité relative peut dépasser 100 %, on est alors dans un cas de sursaturation.

La pression de vapeur saturante est donnée par la formule simplifiée de Clausius-Clapeyron (Clausius 1856; Clapeyron 1834), en considérant que la vapeur d'eau est un gaz parfait et que l'enthalpie de vaporisation ne varie pas avec la température :

$$\ln\left(\frac{P_{sat}}{P_0}\right) = \frac{M \cdot L_v}{R} \cdot \left(\frac{1}{T_0} - \frac{1}{T}\right)$$

où  $T_0$  est la température d'ébullition à une pression  $P_0$  donnée,  $P_{sat}$  est la pression de vapeur saturante,  $M$  est la masse molaire de l'eau,  $L_v$  est la chaleur latente de vaporisation de l'eau,  $R$  est la constante des gaz parfaits et  $T$  est la température. Nous utilisons ici l'approximation de Rankine (Rankine 1849) :

$$\ln\left(\frac{P_{sat}}{1013,25}\right) = 13,7 - \frac{5120}{T}$$

avec la pression en hPa et la température en Kelvin.

On peut alors calculer l'humidité relative en fonction de l'humidité spécifique

$$rh = 100 \cdot \frac{q}{Sat_q}$$

où  $Sat_q = \frac{M_{eau}}{M_{air}} \cdot \frac{P_{sat}}{P - P_{sat}}$  est le rapport de vapeur d'eau saturante.

Dans cette unité, la vapeur d'eau atmosphérique est plus facilement compréhensible. L'humidité spécifique prend des valeurs d'un ordre de grandeur inférieur dans les hautes couches de l'atmosphère, il est plus compliqué de représenter un profil. C'est donc l'humidité relative qui sera utilisée lors des restitutions et l'humidité spécifique qui sera utilisée pour les calculs de transfert radiatif puisque ces derniers sont conçus de cette manière.

### 1.2.1 Le spectre électromagnétique

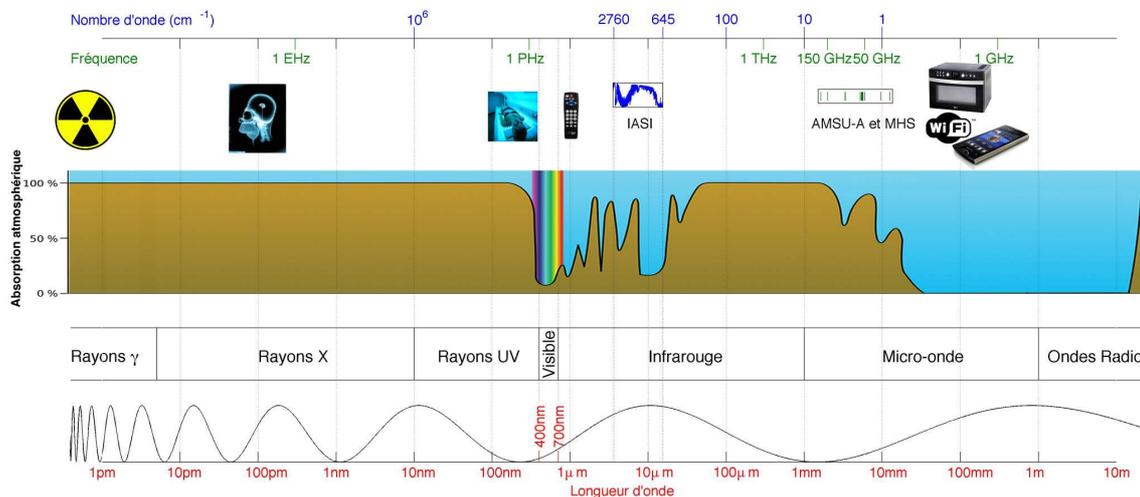


FIGURE 1.9 – Représentation du spectre électromagnétique et correspondance entre les différentes unités. Les utilisations les plus communes des différentes gammes de longueur d’onde sont représentées. Les instruments IASI et AMSU sont également présent.

La Figure 1.9 présente les différentes gammes de rayonnement, leurs principales utilisations et également l’absorption atmosphérique. Une grande partie du rayonnement électromagnétique est absorbé par l’atmosphère. Dans ces régions du spectre il n’est pas possible de sonder l’atmosphère par satellite. En effet, comme on l’a expliqué précédemment, c’est grâce à la capacité du rayonnement à pénétrer l’atmosphère que l’on peut accéder à ses caractéristiques. Dans une région du spectre où l’atmosphère absorbe 100% du rayonnement, du point de vue du satellite, seules les hautes couches de l’atmosphère sont visibles. Tout le rayonnement émis par les couches plus basses sera absorbé par les couches supérieures. C’est pourquoi, comme on le voit sur la figure, les instruments de sondage atmosphérique sont situés dans des régions “fenêtres” du spectre, c’est-à-dire transparentes, ou dans des régions avec une absorption atmosphérique faible. C’est la variation spectrale de l’absorption et de

la transmission atmosphérique qui permet de mesurer des caractéristiques de l’atmosphère à différentes altitudes.

Au cours de cette étude, les fréquences infrarouges seront référencées par leur nombre d’onde (en  $\text{cm}^{-1}$ ) tandis que celle dans le micro-onde seront référencées par leur fréquence (en GHz). La relation entre une fréquence  $f$  en GHz, un nombre d’onde  $\nu$  en  $\text{cm}^{-1}$  et une longueur d’onde  $\lambda$  en  $\mu\text{m}$  (unités utilisées par la suite) est donnée par :

$$\lambda = \frac{10^4}{\nu} = \frac{1000 \cdot c}{f}$$

où  $c$  est la vitesse de la lumière. Le rayonnement électromagnétique se propage à la vitesse de la lumière. La longueur d’onde d’un rayonnement correspond à la distance entre deux maxima du champ propagé. Plus ces maxima sont rapprochés, plus l’énergie véhiculée par le rayonnement est important. Ceci explique que les rayons gamma soient aussi destructeurs. Le rayonnement dans le micro-onde véhicule moins d’énergie que le rayonnement infrarouge. Un rayonnement plus énergétique est plus facile à mesurer, il n’est pas nécessaire d’intégrer sur des grandes zones géographiques (ou dans le temps) pour pouvoir le mesurer. La taille d’une antenne dépend de la longueur d’onde considérée. Plus la longueur d’onde est grande, plus il faut une grande antenne pour avoir un lobe principal équivalent, c’est-à-dire la même résolution. Généralement la résolution spatiale d’un instrument dans le micro-onde (*i.e.*, la zone qu’il sonde) est plus large que celle d’un instrument sondant dans l’infrarouge (pour un satellite volant à la même altitude).

Chaque bande de fréquence a sa propre appellation et sa propre utilisation suivant ses différentes interactions avec le milieu environnant (se reporter à la Figure 1.9). Sur cette figure, le spectre d’absorption atmosphérique est également représenté. Lorsqu’il vaut 100%, l’atmosphère est opaque, inversement, plus il est faible plus l’atmosphère est transparente à ces longueurs d’onde. Le tableau suivant détaille, gamme par gamme, les différentes utilisations du rayonnement électromagnétique.

Longueur d’onde	Fréquence	Nombre d’onde ( $\text{cm}^{-1}$ )	Appellation	Utilisation
Jusqu’à 5 pm	au delà de 60 EHz	moins de $2 \cdot 10^9$	Rayons gammas	Mesures des retombées nucléaires en cas de guerre
De 5 pm à 10 nm	30 à 60000 PHz	$2 \cdot 10^9$ à $10^6$	Rayons X	Imagerie médicale et surveillance aux frontières...

Longueur d'onde	Fréquence	Nombre d'onde ( $\text{cm}^{-1}$ )	Appellation	Utilisation
De 10 nm à 400 nm	750 à 30000 THz	$10^6$ à 25000	Rayons Ultra-Violets	Bronzage, stérilisation en laboratoire, luminothérapie, pièges à insectes
De 400 nm à 700 nm	420 à 750 THz	25000 à 14000	Lumière visible	Vision humaine
De 700 nm à 1 mm	0,3 à 420 THz	14000 à 10	Infrarouge	Caméras thermiques, chauffage, télécommandes
De 1 mm à 1 m	0,3 à 300 GHz	10 à $10^{-2}$	Micro-onde	Wifi, Bluetooth, télévision, téléphonie mobile, four micro-onde, radars météorologiques
Au delà d'1 m	Moins de 300 MHz	Moins de $10^{-2}$	Ondes Radio	Radio FM et AM, communications diverses (VHF...)

Le rayonnement dans le micro-onde est encore divisé en plusieurs bandes suivant sa fréquence : les bandes L (1 à 2 GHz), S (2 à 4 GHz), C (4 à 8 GHz), X (8 à 12 GHz), Ku (12 à 18 GHz), K (18 à 26,5 GHz), Ka (26,5 à 40 GHz), U (40 à 60 GHz). D'autres bandes existent mais sont moins utilisées.

Il existe une autre catégorisation des différentes fréquences établie par l'UIT (Union Internationale des Télécommunications) : TLF (0 à 3 Hz), ELF (3 à 30 Hz), SLF (30 à 300 Hz), ULF (0,3 à 3 kHz), VLF (3 à 30 kHz), LF (30 à 300 kHz), MF (0.3 à 3 MHz), HF (3 à 30 MHz), VHF (30 à 300 MHz), UHF (0,3 à 3 GHz), SHF (3 à 30 GHz), EHF (30 à 300 GHz) et THF (au delà de 300 GHz). L'UIT régule l'utilisation des différentes fréquences afin d'éviter les interférences entre plusieurs utilisations. Cela permet de ne pas brouiller son wifi en appelant avec son portable ou de ne pas influencer les radars météorologiques avec des transmissions aux satellites... Chaque bande de fréquence est dédiée à une utilisation, ou partagée entre plusieurs applications. Cependant, certaines bandes de fréquence sont encore "piratées" et des antennes émettent de façon illicite. Ces émissions entraînent de larges imprécisions dans les mesures de certains instruments satellite, comme par exemple la mission SMOS (Soil Moisture and Ocean Salinity) lancé en 2009, qui sonde en bande L. Ce dernier a souffert de nombreux problèmes de RFI (Radio Frequency Interference). Des méthodes de décontamination doivent être mises en place, mais un bruit résiduel persiste.

## 1.2.2 Le sondage dans le micro-onde

Le spectre d'absorption atmosphérique dans le micro-onde est majoritairement dominé par le continuum d'absorption de l'eau et par les raies d'absorption de l'oxygène. Le spectre présente de nombreuses raies rotationnelles de l'eau. On retrouve également la plus forte raie rotationnelle de l'eau à 556,9361 GHz. De nombreux capteurs sondent autour de la raie d'absorption de H<sub>2</sub>O à 138,3 GHz afin de mesurer la quantité de vapeur d'eau.

La mesure du dioxygène (O<sub>2</sub>) permet d'avoir accès indirectement à des mesures de températures. Le dioxygène présente une bande d'absorption à environ 57 GHz, une autre raie d'absorption isolée à 118,75 GHz, ainsi que de nombreuses autres raies à plus hautes fréquences. La plupart des sondeurs micro-ondes exploitent la raie d'absorption à 57 GHz.

De nombreuses autres molécules présentent des raies d'absorption dans les micro-ondes. On peut citer par exemple l'ozone (O<sub>3</sub>), le dioxyde d'azote (NO<sub>2</sub>), le monoxyde de chlore (ClO), le dioxyde de soufre (SO<sub>2</sub>), le protoxyde d'azote (N<sub>2</sub>O) ou encore le dioxyde de soufre (SO<sub>2</sub>). Cependant, la faible concentration de ces constituants dans l'atmosphère, ajoutée à leur faible absorption, les rend difficilement détectables dans le micro-onde. Une géométrie de sondage aux limbes est nécessaire pour mesurer la teneur de l'atmosphère en ces molécules. Il s'agit dans ce cas de viser l'atmosphère de façon parallèle à la surface de la Terre afin que le rayonnement mesuré ne traverse que l'atmosphère et ne soit pas issu de la surface. De plus, le rayonnement mesuré traverse une longueur supérieure de l'atmosphère, ce qui permet de mesurer des concentrations plus importantes. On obtient une excellente résolution verticale, mais en contre partie la résolution spatiale est moins bonne. Cette géométrie n'est pas utilisée dans cette étude mais les outils utilisés ici peuvent être appliqués à de tels sondeurs.

L'avantage du sondage dans le micro-onde est la faculté qu'a le rayonnement, à cette fréquence, de traverser en partie les nuages. Ainsi, des mesures peuvent être réalisées sous les nuages et donner accès à des zones inaccessibles aux rayonnements visibles ou infrarouges. Certains sondeurs utilisent également des régions spectrales "fenêtres" afin d'avoir accès aux caractéristiques du sol (comme SMOS précédemment cité), mais également des précipitations et des nuages, ou pour mesurer la quantité totale d'eau dans la colonne atmosphérique sondée (comme l'instrument MADRAS, Microwave Analysis and Detection of Rain and Atmospheric Structures, radiomètre imageur avec 5 canaux de 18 à 157 GHz).

Malgré le nombre bien plus important de canaux et donc de volume de données des instruments infrarouges récents, les observations micro-ondes sont encore de nos jours plus importantes dans les centres de météorologie. Ceci s'explique en partie par leur faible sensibilité aux nuages (non précipitants) et aux aérosols et donc la possibilité de les utiliser en permanence, mais également par leur faible sensibilité aux constituants mineurs de l'atmosphère peu connus et donc difficilement quantifiables.

### 1.2.3 Le sondage dans l'infrarouge

L'absorption dans l'infrarouge est majoritairement constituée de bandes de vibration-rotation. Seule la vapeur d'eau présente également une bande rotationnelle pure dans cette gamme de longueur d'onde. Les principaux absorbeurs sont la vapeur d'eau ( $\text{H}_2\text{O}$ ), le dioxyde de carbone ( $\text{CO}_2$ ) et l'ozone ( $\text{O}_3$ ). De nombreux autres gaz, dits gaz traces car moins présents dans l'atmosphère, y sont également actifs. On peut citer le méthane ( $\text{CH}_4$ ), le monoxyde de carbone ( $\text{CO}$ ) ou le protoxyde d'azote ( $\text{N}_2\text{O}$ ). La plupart de ces gaz se trouvent dans la troposphère (partie basse de l'atmosphère).

Les particules diatomiques, majoritaires dans l'atmosphère, sont presque inactives dans l'infrarouge. Cela permet d'avoir accès à la mesure de l'absorption par les composants mineurs<sup>7</sup>. La nouvelle génération de capteurs infrarouges hyperspectraux (CrIs, AIRS, IASI, voir Section 1.3.1, page 29) mesure précisément le rayonnement émis par l'atmosphère et la surface, à des longueurs d'onde précises. Ceci donne accès à la mesure des bandes d'absorption de chaque constituant atmosphérique. La connaissance de ces nombreuses bandes d'absorption permet de déterminer la composition chimique de l'atmosphère.

Le sondage dans l'infrarouge permet également la mesure de la température et de la vapeur d'eau de l'atmosphère. Les bandes d'absorption de la vapeur d'eau donnent directement accès à sa mesure. Celles du  $\text{CO}_2$  permettent de déterminer la température de l'atmosphère. Cependant, l'absorption du rayonnement par les hydrométéores rend les nuages imperméables au rayonnement infrarouge. Il n'est pas possible de mesurer par satellite le rayonnement infrarouge sous un nuage. Même si la complémentarité des rayonnements infrarouges et micro-ondes est évidente dans le cas de situations nuageuses, nous avons décidé de n'utiliser que le cas simple en ciel clair afin de bien mettre en avant la synergie entre les rayonnements. Si la synergie est présente en ciel clair, elle sera encore plus forte pour les cas nuageux ou précipitants, où les mesures dans les différentes gammes de longueurs d'onde donnent accès à des informations sensiblement différentes, car le rayonnement infrarouge ne traverse pas les nuages.

## 1.3 La plate-forme MetOp

La plate-forme MetOp est une plate-forme satellite prévue pour être lancée en trois exemplaires. Elles ont été construites par EADS Astrium pour l'ESA (European Spatial Agency) et EUMETSAT (EUropean organisation for the exploitation of METeorological SATellites). La première a été lancée le 19 octobre 2006 (nommée MetOp-A), la deuxième a été lancée le 17 septembre 2012 (nommée MetOp-B), la troisième devrait être lancée courant 2017 ou 2018. Cette série de satellites identiques munis des mêmes capteurs permet une continuité des mesures. Chaque satellite est lancé de façon à voler, pendant une courte période, de façon simultanée avec le satellite précédent. Il est alors possible de comparer les

7. Plus la résolution de l'instrument est fine plus on a accès à la mesure de fines bandes d'absorption

différentes mesures et de calibrer les différents capteurs afin d’avoir des données similaires d’un satellite à l’autre. Ces plateformes font partie du programme EPS (EUMETSAT Polar System) (Klaes et al. 2007).

Les plateformes à orbite géostationnaire sont situées à 36000 km au-dessus de l’équateur. Elles tournent à la même vitesse que la Terre sur elle-même. Elles sont situées en permanence au-dessus du même point et permettent de surveiller en continu toute une zone du globe (un disque d’environ 60°). Par conséquent, elles ne peuvent couvrir la totalité de la Terre. Il faut pour cela plusieurs satellites, ce qui pose des problèmes d’intercalibration entre les instruments à bord des différentes plates-formes. Ces plates-formes sont situées au-dessus de l’équateur, elles n’observent pas les pôles, où de nombreux phénomènes climatiques importants ont lieu. Les changements de température et les vents au niveau des pôles ont une influence directe sur le climat, notamment en Europe et en Amérique du nord. C’est dans le but de suppléer ces satellites géostationnaires que l’orbite de MetOp a été choisie.

Cette plate-forme orbite à 824 km de la Terre en orbite polaire, c’est-à-dire qu’elle survole les pôles à chaque orbite (voir Figure 1.10). C’est une orbite héliosynchrone, c’est-à-dire en phase avec la rotation de la Terre sur elle-même. MetOp survole donc chaque point du globe à heure fixe. Autrement dit, son plan d’orbite garde la même orientation par rapport au Soleil. À chaque tour, il se décale vers l’ouest à la même vitesse que le Soleil. Il passe donc au-dessus de points à la même latitude à heure fixe localement. Sur la Figure 1.10, il croise l’équateur en phase ascendante à la même heure locale. Il s’agit d’un satellite du “matin”, il franchit l’équateur à 9h30 (heure solaire locale) en orbite descendante, il fait le tour du globe en 101 minutes. En effectuant 14 orbites par jour, il survole chaque point de la Terre tous les cinq jours.

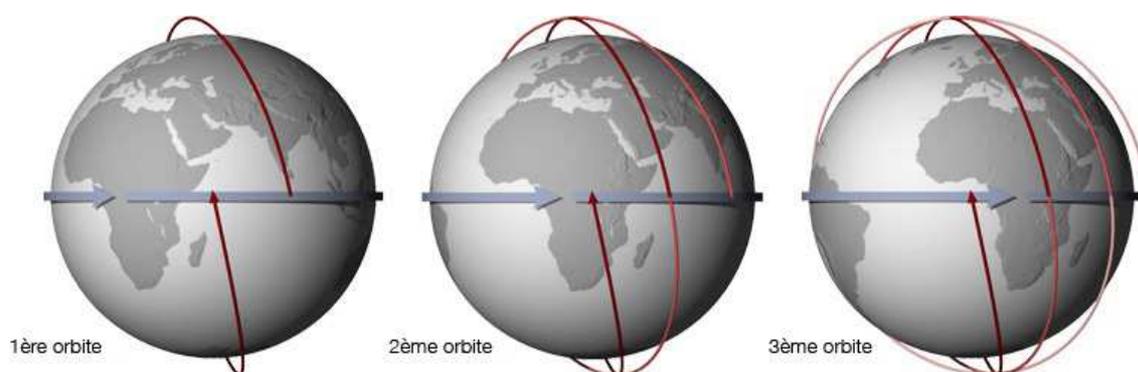


FIGURE 1.10 – Représentation de l’orbite polaire de MetOp (flèche rouge) en prenant en compte la rotation de la Terre sur elle-même (flèche bleue). Source EADS Astrium : <http://www.astrium.eads.net/>.

MetOp a été conçu dans un but météorologique. Sa fonction première est d’acquérir des profils atmosphériques, d’observer les nuages, les vents et la surface de la Terre. Ces données sont ensuite introduites dans les modèles de prévisions météorologiques et de surveillance

climatique. Il est aussi équipé pour des missions à but humanitaire (sauvetage, recherche...). Il transporte à bord douze instruments :

- IASI : il s’agit d’un interféromètre sondant entre 645 et 2760  $\text{cm}^{-1}$  qui permet la mesure de différents profils atmosphériques et de caractéristiques de surface ;
- AMSU-A1 et AMSU-A2 : ce sont des capteurs micro-ondes disposant de 15 canaux entre 23 et 90 GHz qui servent à mesurer la température atmosphérique, mais également les précipitations et les caractéristiques de la surface ;
- MHS : c’est un sondeur micro-onde doté de cinq canaux entre 89 et 190 GHz qui mesurent l’humidité atmosphérique ;
- A-DCS : il s’agit d’un capteur de signal UHF (voir Section 1.2.1, page 21) qui sert à recueillir des données environnementales ;
- ASCAT : c’est un scattéromètre (analyse de rayonnement diffracté) en bande C (voir Section 1.2.1, page 21) qui sert essentiellement à faire des mesures du vent à la surface des océans ;
- AVHRR : ce radiomètre imageur possède 6 canaux dans le visible, le proche infrarouge et l’infrarouge qui permet de faire des images des terres, des océans et des nuages ;
- GOME-2 : il s’agit d’un spectromètre dans le visible et l’ultra-violet qui permet de faire notamment des mesures d’ozone ;
- GRAS : c’est un récepteur radio (mesure de l’affaiblissement du signal GPS) permettant la mesure de profils de température et d’humidité ;
- HIRS : ce sondeur infrarouge possède 19 canaux entre 666 et 2631  $\text{cm}^{-1}$  et un canal dans le domaine visible. Cet instrument “historique” mesure les profils atmosphériques de température mais également la hauteur des nuages ;
- S&R : Ce récepteur, émetteur et processeur de signaux VHF et UHF permet la réception et le traitement des signaux d’urgence des navires et aéronefs ;
- SEM-2 : Il s’agit d’un spectromètre multi-canaux qui permet la mesure de différents flux de particules ionisées.

Dans cette étude, nous nous intéresserons uniquement aux trois premiers instruments cités : IASI, AMSU-A et MHS. Cependant, un raisonnement similaire pourrait être mené avec d’autres instruments.

Afin d’augmenter la couverture du globe par les différents capteurs, certains sont mobiles et “balayent” le sol sous le satellite. Au cours de l’avancée du satellite, ils observent différents points. Il existe deux modes de balayage :

- Le balayage conique, le capteur est orienté avec un certain angle et sa trace au sol correspond à des arcs de cercle de sorte à avoir un angle constant avec le sol ;
- Le balayage “cross-track” (*i.e.*, balayage transverse), la fauchée de l’instrument est perpendiculaire à la trace du satellite, l’angle avec le sol varie.

La Figure 1.11 met en évidence ces deux types de balayage. Sur cette figure on peut voir deux satellites munis chacun d’un instrument de sondage. L’instrument de gauche est à

balayage conique et celui de droite est un “cross-track”. L’orbite du satellite est représentée en pointillés blancs. Sa projection au sol est appelée la “trace” du satellite. Les ellipses rouges au sol représentent la zone que sonde l’instrument. Sur l’image de gauche, on peut constater que l’angle d’incidence (noté  $\alpha$ ) est constant. Ainsi, la zone sondée est toujours la même et la colonne atmosphérique depuis la zone jusqu’au sondeur est également la même. Sur l’image de droite, la variation de l’angle d’incidence fait varier la zone mesurée au sol mais également la colonne de l’atmosphère traversée. Le trajet du rayonnement depuis la surface jusqu’au capteur est représenté par un trait rouge. Plus le sondeur vise sur le côté plus le rayonnement traverse une couche épaisse d’atmosphère avant de parvenir au capteur. Il faut donc prendre en compte l’angle de visée lorsque l’on s’intéresse aux capteurs “cross-track”. De plus, la taille et la forme des pixels au sol est différente : plus petits et circulaires proche de la trace et plus grands et elliptiques en bordure de la fauchée. L’angle qui sera utilisé par la suite est l’angle zénithal. Il correspond à l’angle entre la verticale sous le satellite et le rayonnement reçu par le capteur. Il est noté en vert sur la figure. On appelle nadir le sondage effectué avec un angle zénithal nul.

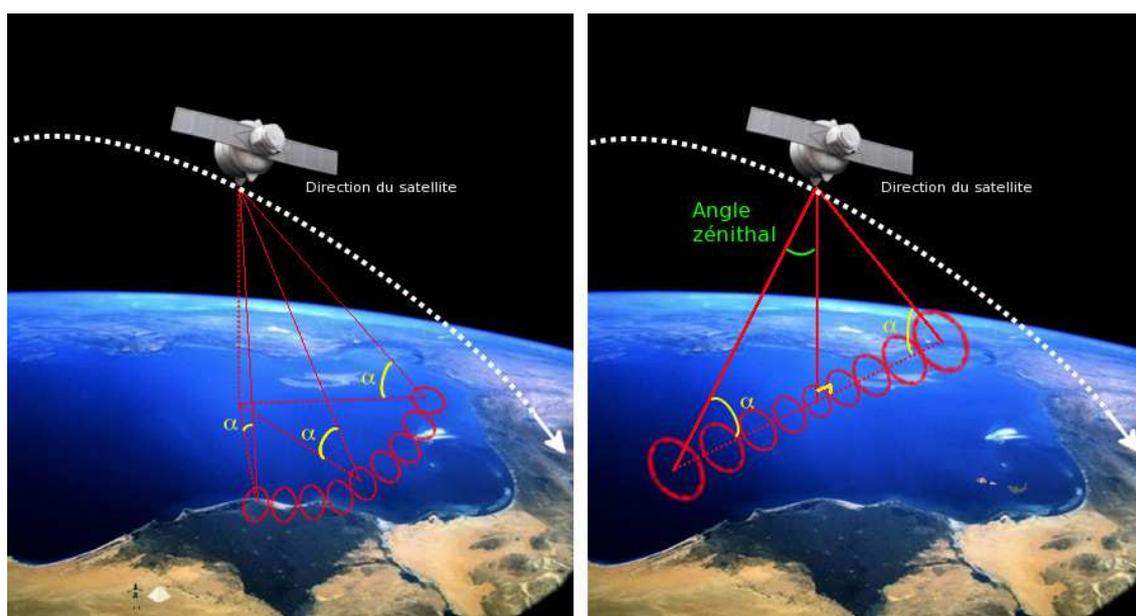


FIGURE 1.11 – Schéma des géométries de balayage pour les capteurs satellite : balayage conique à gauche et cross-track à droite.

L’inconvénient des capteurs à balayage conique est la difficulté de calibration en vol. En effet, la méthode utilisée pour calibrer en vol les capteurs à balayage “cross-track” consiste en la mesure d’une cible froide (dans l’espace) et d’une cible chaude à bord de la plateforme. Une telle méthode n’est pas possible avec un balayage conique. C’est pourquoi la majorité des capteurs sont à balayage “cross-track”.

Les capteurs que nous utiliserons dans cette étude (IASI, AMSU-A et MHS) sont des

capteurs “cross-track”. Du fait de cette fauchée plus large que la trace même du satellite, les capteurs peuvent recouvrir la totalité de la surface terrestre plus rapidement que la trace du satellite. Si MetOp met 5 jours complets pour survoler chaque point de la Terre, le capteur IASI aura vu chaque point du globe au bout de deux jours seulement<sup>8</sup>.

La Figure 1.12 représente la trajectoire de MetOp pendant une demi-journée. Le trait rouge indique le chemin suivi par MetOp en lui-même, les traits orange indiquent la fauchée de IASI et les traits verts représentent la trace du satellite. On constate ici que la largeur des fauchées permet d’étendre la surface couverte par le satellite.

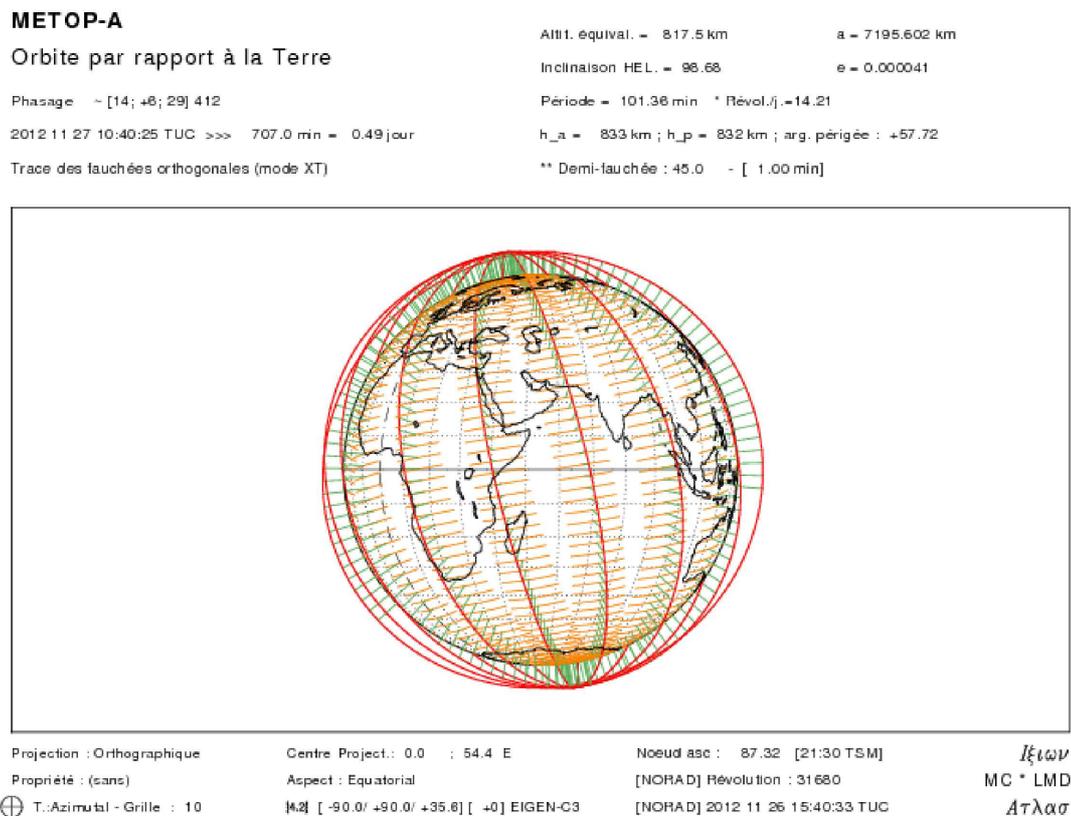


FIGURE 1.12 – Représentation de l’orbite polaire de MetOp (en rouge) ainsi que de la fauchée (en orange) d’un instrument (IASI) grâce au logiciel IXION <http://climserv.ipsl.polytechnique.fr/ixion.html>, pour une demi-journée (partie nuit de l’orbite).

### 1.3.1 IASI

Cette étude porte plus précisément sur les mesures effectuées par IASI (Infrared Atmospheric Sounding Interferometer) (Chalon et al. 2001). Il est constitué de 8461 canaux répartis linéairement entre 645 et 2760  $\text{cm}^{-1}$  (voir Figure 1.9). Le fait que chaque mesure

8. Il reste une petite zone au pôle qui n’est pas survolée par le satellite, du fait de son angle de  $92^\circ$ .

soit constitué de 8461 valeurs à des longueurs d'ondes différentes fait que la quantité de données mesurées par IASI est colossale. Il représente à lui seul quasiment la moitié des données transmises quotidiennement par MetOp. Ceci explique pourquoi nous chercherons des moyens de réduire au mieux cette quantité de données. La mesure s'effectue grâce à un spectromètre infrarouge passif bien calibré à transformée de Fourier. IASI est séparé en trois bandes :

- La bande B1 de 645 à 1210  $\text{cm}^{-1}$  ;
- La bande B2 de 1210 à 2000  $\text{cm}^{-1}$  ;
- La bande B3 de 2000 à 2760  $\text{cm}^{-1}$ .

La bande B3 a un rapport signal sur bruit très faible et est donc difficilement exploitable. Sur la Figure 1.2, on peut remarquer que à 2760  $\text{cm}^{-1}$ , l'émission de la Terre est plus faible. De plus, l'émission du Soleil (même si très faible lorsque le capteur n'est pas dans l'axe du reflet du Soleil sur la Terre) est plus forte. La combinaison de ces deux facteurs fait que le signal à mesurer est plus faible dans la bande B3 et le bruit instrumental y est élevé, il est donc difficile d'extraire des informations de cette bande.

IASI est un capteur "cross-track" qui sonde donc de façon perpendiculaire au déplacement du satellite. La résolution spatiale de IASI est de 12 km au nadir, tandis que lorsque l'angle zénithal est à son maximum ( $\pm 48.3^\circ$ ) la résolution spatiale est de 27 km.

D'après la Section 1.2.3 (page 25), chaque partie de son spectre peut être utile à mesurer l'absorbant qui y est actif. D'après Aires (1999), on a pour IASI :

Nombre d'onde ( $\text{cm}^{-1}$ )	Caractéristiques atmosphériques mesurées
650-770	Température
770-980	Nuages et surface
1000-1070	Ozone
1080-1150	Nuages et surface
1210-1650	Vapeur d'eau, température, $\text{N}_2\text{O}$ , $\text{CH}_4$ et $\text{SO}_2$
2100-2150	Quantité totale de CO
2150-2250	Température et quantité totale de $\text{N}_2\text{O}$
2350-2420	Température
2420-2700	Nuages et surface
2700-2760	Quantité totale de $\text{CH}_4$

Ainsi chaque partie du spectre est sensible à différentes caractéristiques de l'atmosphère. Le fait d'utiliser les différents canaux de façon simultanée pour obtenir des informations sur une des variables atmosphériques (température, vapeur d'eau...) est déjà, en soit, une utilisation de la synergie. Les informations contenues dans chaque canal se couplent les unes aux autres afin de mieux déduire les variables recherchées.

La Figure 1.13 représente les jacobiens des différents canaux de IASI en température (à gauche) et en vapeur d'eau (à droite). Pour chaque canal, la dérivé de la température de

brillance en fonction de la température (à gauche) ou de l'humidité relative (à droite) a été calculée à l'aide de RTTOV (voir Section 1.4, page 33), pour des situations au-dessus des océans, en ciel clair.

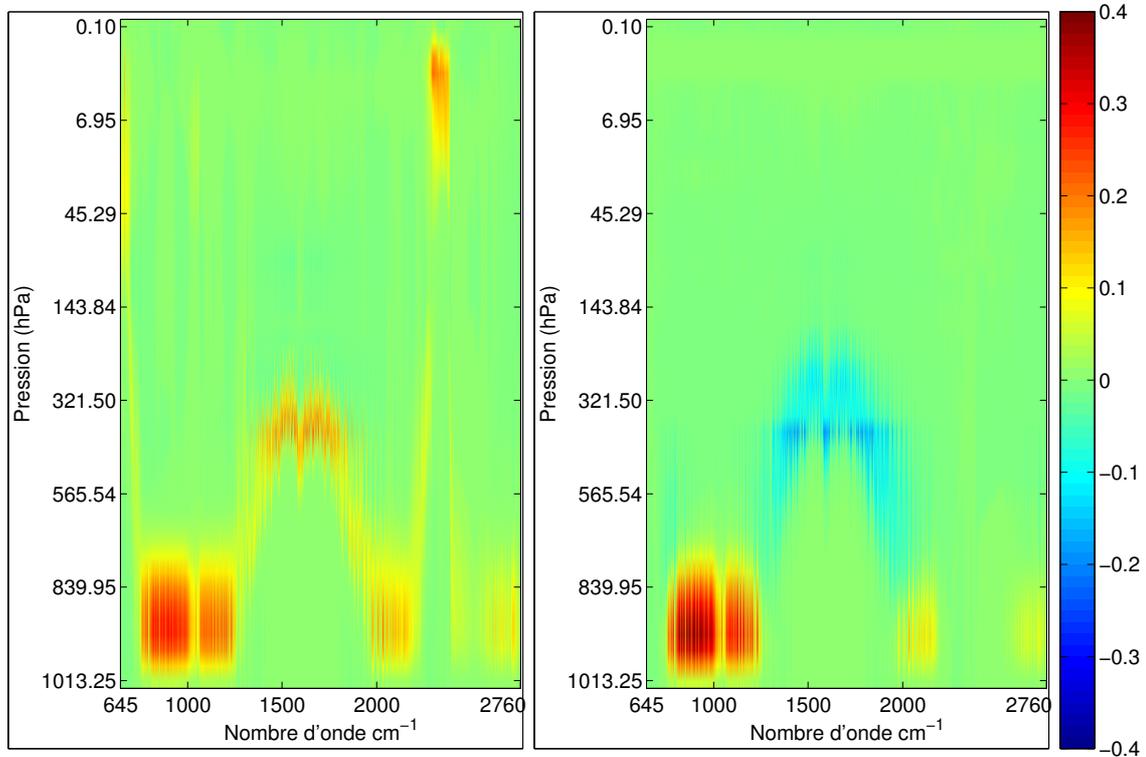


FIGURE 1.13 – Jacobien en température (à gauche) et en humidité relative (à droite) de IASI.

On retrouve sur ces jacobiens la sensibilité des différentes bandes spectrales présentées dans le tableau précédent. On remarque également que les zones sensibles à la vapeur d'eau (1210 à 1650  $\text{cm}^{-1}$ ) sont également sensibles à la température. Ceci est dû à la corrélation entre la température et la vapeur d'eau. En utilisant tous les canaux de IASI, on a une sensibilité à la température tout le long de l'atmosphère. C'est donc un instrument adapté à la restitution de profils atmosphériques de température. On ne retrouve pas de sensibilité à la vapeur d'eau dans les hautes couches de l'atmosphère, car la teneur en vapeur d'eau y est très faible. Cependant, la vapeur d'eau dans les basses couches de l'atmosphère influence de nombreux canaux. IASI est donc également adapté pour restituer la vapeur d'eau.

### 1.3.2 AMSU-A

Nous nous intéresserons également dans cette étude au radiomètre AMSU-A (Advanced Microwave Sounding Unit-A). Il sonde dans la bande de l'oxygène entre 50 et 60 GHz. Il est dédié aux restitutions de profils de température (Mo 1996). Il s'agit également d'un capteur "cross-track" avec un angle allant jusqu'à 48,3°. Chacun des scans est composé de

30 sondages. Il parcourt un scan (de  $-48,3$  à  $48,3^\circ$ ) en 8 secondes. Chaque scan fait une largeur de 1650 km et chaque sondage a une résolution au sol de 48 km au nadir (ce qui représente quatre sondages de IASI).

AMSU-A est composé de deux modules :

- A-1 : il est composé des canaux 3 à 15. Les canaux 3 et 4 (50,3 et 52,8 GHz) servent à mesurer les caractéristiques de surface et la quantité de vapeur d'eau dans l'atmosphère. Les canaux 5 à 8 (53,596 54,4 54,94 et 55,5 GHz) servent à mesurer la température de la troposphère (partie basse de l'atmosphère). Les canaux 9 à 14, tous autour de 57,290 GHz, servent à mesurer la température stratosphérique. Le canal 15 à 89 GHz sonde dans une région "fenêtre" du spectre (l'atmosphère y est transparente) et donne accès aux nuages ou aux caractéristiques de surface ;
- A-2 : il est composé des canaux 1 et 2 qui sondent respectivement à 23,8 et 31,4 GHz. Ces canaux servent essentiellement à mesurer la quantité de vapeur d'eau dans l'atmosphère.

AMSU-A a été conçu pour restituer de façon globale des profils atmosphériques de température et de vapeur d'eau depuis la surface jusqu'à la haute stratosphère ( $\approx 50$  km). Il sert également à mesurer les précipitations et des informations de surface comme la couverture neigeuse, la glace ou l'humidité du sol.

### 1.3.3 MHS

MHS est un instrument à géométrie d'acquisition similaire à celle d'AMSU-A, avec une meilleure résolution au sol ([Hewison and Saunders 1996](#)). Chacun des sondages a une résolution de 15 km au nadir. Quand les deux instruments sont utilisés ensemble (AMSU-A et MHS), au sein de chaque sondage de AMSU-A, une moyenne des 9 sondages de MHS est faite (matrice de 3 par 3) pour revenir à la même résolution (dans l'application que l'on en fait dans cette étude, sinon chacun des sondages MHS est utilisé en tant que tel).

MHS est dédié aux restitutions de l'humidité. Il est composé de 5 canaux (89,9 et 157 GHz et 3 canaux à 183,31 GHz à différentes résolutions spectrales). Les canaux de MHS et AMSU-A sont représentés sur la Figure 1.9.

AMSU-A et MHS forment avec le sondeur infrarouge HIRS (qui a une résolution spatiale de 10 km) ce que l'on appelle ATOVS (Advanced TIROS Operational Vertical Sounder).

Comme nous l'avons expliqué, les différentes résolutions spatiales de AMSU-A et MHS sont traitées en moyennant les sondages de MHS pour obtenir la résolution spatiale de AMSU-A.

La Figure 1.14 présente les jacobiens de AMSU-A et MHS en température (à gauche) et en humidité relative (à droite). De la même façon que pour les jacobiens de IASI, ces dérivés partielles ont été calculées grâce à RTTOV (voir Section 1.4, page 33), pour des situations au-dessus des océans, en ciel clair. Ici encore, les canaux sensibles à la température (AMSU-A notamment) sont également sensibles à la vapeur d'eau de par leur corrélation. L'inverse

est également vrai si on regarde les canaux de MHS, plutôt sensibles à la vapeur d'eau. En combinant ces deux instruments, on a une sensibilité à la température tout le long de l'atmosphère et on peut donc restituer un profil complet. La faible quantité de vapeur d'eau dans les hautes couches de l'atmosphère implique également une très faible sensibilité des instruments. La sensibilité des instruments à l'humidité dans les basses couches permet d'obtenir un profil d'humidité. Ces jacobiens viennent justifier le choix des instruments pour notre application.

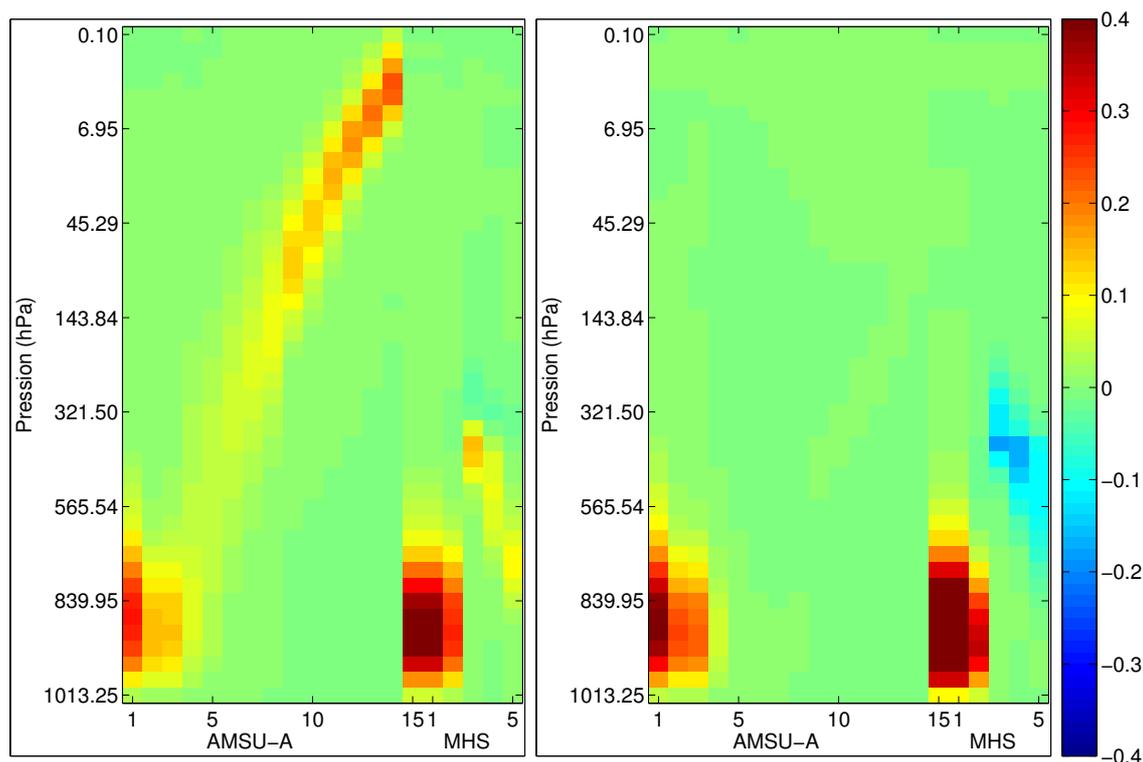


FIGURE 1.14 – Jacobien en température (à gauche) et en humidité relative (à droite) des 15 canaux de AMSU-A et des 5 canaux de MHS.

## 1.4 Le modèle de transfert radiatif RTTOV

Nous utilisons dans cette étude le code de transfert radiatif RTTOV (Radiative Transfer for TOVS). Il s'agit d'un modèle de transfert radiatif rapide développé par l'ECMWF (Eyre 1991) et qui est maintenant soutenu par l'Eumetsat NWP-SAF (Satellite Application Facility) (Saunders et al. 1999; Matricardi et al. 2004). Ce code permet de simuler des luminances rapidement dans l'infrarouge ou le micro-onde à partir de l'information des profils atmosphériques de température et de vapeur d'eau, de diverses concentrations en gaz, de la couverture nuageuse et des propriétés de la surface. Il a été implémenté pour la plupart des instruments existants. Sa rapidité est due à des matrices de coefficients, qui ont été calculées pour chaque instrument afin de ne pas avoir à recalculer le transfert radiatif intégré sur toute

la colonne atmosphérique, pour toute nouvelle situation et pour chaque canal. RTTOV-10 permet d’avoir accès au spectre de luminances ainsi qu’au jacobien associé, c’est-à-dire la matrice de sensibilité de la luminance aux variables géophysiques décrivant la situation atmosphérique. Si certains champs ne sont pas connus, une base de données incluse permet de les estimer. Elle contient notamment un profil moyen d’ozone et de gaz traces. RTTOV contient également un simulateur d’émissivité au-dessus des océans et des terres dans le micro-onde (voir Section 2.1.1, page 41, (Aires et al. 2011b)), et dans l’infrarouge (voir Section 2.1.2.3, page 46, (Seemann et al. 2008)). Ce modèle de transfert radiatif a été validé grâce à des comparaisons avec les mesures réelles des différents instruments (Saunders et al. 2012).

## 1.5 Les centres de prévision météorologique

### 1.5.1 Dénomination des données satellites

Les données satellites sont classées en deux catégories : les données L1 (level 1) et les données L2 (level 2).

Les données L1 correspondent aux mesures des capteurs satellites. Elles sont fournies sur la grille correspondant à la résolution de l’instrument considéré, géolocalisées et éventuellement calibrées. Elles sont principalement utilisées par les NWP.

Les données L2 correspondent aux quantités géophysiques calculées à partir des L1 représentées sur la même grille que les L1. Ce sont des restitutions de paramètres atmosphériques ou de surface. Ces données sont également utilisées par les NWP mais aussi dans les centres de suivi du climat, ou pour l’observation de la Terre et le suivi de certains paramètres. Dans le contexte actuel de réchauffement climatique, ces suivis sont de plus en plus réclamés par le grand public. Toutes sortes de données peuvent être extraites : la quantité de neige, de glace, la température, l’humidité, l’ozone et aussi des gaz moins répandus comme le méthane... Le développement d’algorithmes permettant d’avoir accès de façon plus précise à ces différents paramètres est crucial afin de permettre de réduire les barres d’erreur et de mettre en évidence des tendances d’évolution plus fines.

### 1.5.2 Les centres de prévision numérique : NWP

Les centres de prévision météorologiques (NWP pour “Numerical Weather Prediction” center) sont chargés de prédire le temps. Ils ont pour cela trois sources d’information distinctes : les mesures *in situ* (ballons, stations météorologique, avions, bouées...), les mesures par satellites et les données issues des modèles de prévision (GCM pour “Global Circulation Model”). À partir de ces données, leur objectif est de prévoir le futur. Si seuls les modèles de prévision étaient utilisés, la côté chaotique des approximations utilisées les rendrait instables. Les données *in situ* et issues de mesures satellites permettent alors de contraindre

les sorties des modèles.

La complexité des modèles climatiques fait que certains des plus puissants ordinateurs actuels sont utilisés au sein des NWP. Le domaine météorologique a toujours été à la pointe de la technologie en termes de fusion d'informations (entre les données issues des modèles et des mesures). Aujourd'hui les techniques utilisées reposent sur de l'assimilation variationnelle et des filtres de Kalman. Ces techniques sont très performantes pour fusionner les informations et donc exploiter leur synergie potentielle.

Les modèles sont basés sur une discrétisation du temps et de l'espace. La Terre est découpée en mailles distinctes qui interagissent l'une avec l'autre et l'évolution temporelle se fait par paliers. Une telle méthode est très dépendante des conditions initiales considérées. Il est important de pouvoir caractériser l'atmosphère le plus précisément possible.

L'affinage de l'état de l'atmosphère par rapport aux mesures satellites se fait par assimilation variationnelle aussi appelée "3D-var" (et maintenant "4D-var") (Courtier et al. 1998). Il s'agit de modifier petit à petit chacune des variables de l'atmosphère considérées afin que les simulations de sondages atmosphériques soient le plus semblables possibles aux vraies luminances. Comme nous le verrons dans la suite de cette étude, une telle méthodologie est idéale pour exploiter la synergie entre les différents capteurs. Il est cependant important de définir précisément les différentes matrices de propagation d'erreurs (voir Chapitre 4, page 107). Sans utiliser les modèles, nous cherchons à développer des méthodologies capables d'utiliser de façon efficace les informations des nombreux satellites, qui sillonnent le ciel, sondant notre atmosphère pour mieux le déterminer. Au final, nous aurons des résultats indépendants et que l'on pourra comparer à ceux des NWP.

Les NWP fournissent donc des prévisions météorologiques, mais leur rôle ne se limite pas à cela. Ils fournissent également des analyses en temps réel. Il s'agit de données sur l'état géophysique de la Terre (profils atmosphériques, état de la surface...), maille par maille à un moment donné. Ces données sont actuellement disponibles toutes les 6 heures. Ils fournissent également des re-analyses. Ces dernières consistent en un calcul *a posteriori* de l'état géophysique de la Terre. Le calcul de ces re-analyses est effectué sur une période de temps donnée pendant laquelle toutes les mesures utilisées sont intercalibrées pour assurer leur cohérence. Elles sont également disponibles de façon globale toutes les 6 heures.



# RESTITUTION DE LA TEMPÉRATURE DE SURFACE ET DE L'ÉMISSIVITÉ INFRAROUGE



## Sommaire

<b>2.1</b>	<b>L'émissivité</b>	<b>40</b>
2.1.1	Émissivité de surface dans le micro-onde	41
2.1.2	Émissivité de surface dans l'infrarouge	42
2.1.2.1	Bases de données d'émissivités mesurées en laboratoire	43
2.1.2.2	L'émissivité restituée à partir de MODIS	44
2.1.2.3	La base de données UWIRemis	46
2.1.2.4	L'émissivité IASI de la NASA	48
2.1.2.5	L'émissivité IASI du groupe ARA	48
2.1.3	Mise en coïncidence des bases	49
<b>2.2</b>	<b>Construction d'une base de données d'émissivités infrarouges</b>	<b>49</b>
2.2.1	Les réseaux de neurones	50
2.2.2	Base hyperspectrale d'émissivité	51
2.2.2.1	Une base de première ébauche indépendante de IASI	51
2.2.2.2	Une classification de surface	51
2.2.2.3	Création de la base à partir des observations MODIS	52
2.2.3	Analyse en composantes principales des spectres d'émissivité	53
2.2.3.1	Exemple simple d'analyse en composantes principales	53
2.2.3.2	Méthodologie	53
2.2.3.3	Résultats obtenus sur les spectres d'émissivité	54
2.2.3.4	Compression des émissivités infrarouges	55
2.2.4	Interpolation spectrale de l'émissivité infrarouge	56
2.2.5	La première ébauche	59
<b>2.3</b>	<b>Un nouvel algorithme de restitution de surface</b>	<b>60</b>
2.3.1	Inversion bayésienne du transfert radiatif	61
2.3.2	Sélection des paramètres de la restitution	65
2.3.2.1	La matrice de covariance d'erreur de $F$	65

2.3.2.2	La matrice de covariance d'erreur de la première ébauche	66
2.3.2.3	Les canaux sélectionnés	67
2.3.2.4	Le nombre de composantes de l'ACP	68
2.3.2.5	Conclusion	70
2.3.3	Influence de la première ébauche	70
2.3.4	Résultats et évaluation	71
2.3.4.1	Conditions expérimentales	71
2.3.4.2	Analyse spectrale	72
2.3.4.3	Évaluation de la température de surface	76
<b>2.4</b>	<b>Conclusion</b>	<b>82</b>

Les restitutions atmosphériques de sondages satellites restent limitées par la connaissance précise de l'interaction entre la surface et le rayonnement, spécialement dans les basses couches de l'atmosphère. Ces couches sont au coeur des problématiques de l'Homme, qui y est relativement confiné par ses dimensions réduites. Caractériser cette interaction permet de dissocier le rayonnement de l'atmosphère, de celui émis par la surface et donc d'améliorer les restitutions.

Deux variables permettent de déterminer le rayonnement de la surface : l'émissivité et la température. Notre connaissance de ces deux variables est, aujourd'hui, insatisfaisante, compte tenu de la précision recherchée dans les restitutions atmosphériques. Si leur variabilité reste faible au-dessus des océans, elle est plus complexe au-dessus des terres émergées. La variabilité de la température est nettement plus importante sur les continents. Elle est soumise au cycle diurne, particulièrement sur les zones arides et semi-arides. De plus, l'émissivité de la surface est plus variable au-dessus des continents, car l'hétérogénéité des surfaces continentales, comme la nature du sol, la quantité de végétation, l'humidité du sol ou la neige, influencent l'émissivité.

Diverses méthodes ont été utilisées dans les centres opérationnels mais aucune ne donne entière satisfaction. Certaines méthodes considèrent des modèles (Weng et al. 2001) qui utilisent les caractéristiques du sol, comme sa composition ou son humidité, pour estimer les émissivités de la surface. Cependant, une connaissance préalable de l'état de la surface est nécessaire et celle-ci n'est pas systématiquement disponible. De plus, ces modèles ne sont pas idéaux et ne représentent pas l'émissivité exacte des matériaux. Beaucoup de travaux ont été effectués, ces dernières années, pour mieux caractériser les émissivités, que ce soit dans l'infrarouge (Yu et al. 2008) ou dans le micro-onde (Aires et al. 2011b).

Il existe deux grandes familles d'algorithmes d'inversion des caractéristiques de surface. D'un côté, les algorithmes dit statistiques, qui sont basés sur une approche empirique. Une base d'apprentissage est constituée, avec des mesures satellites et les émissivités correspondantes. Cette base sert alors à calibrer un modèle statistique qui permet d'estimer l'émissivité. Une telle méthode est utilisée à la fois dans l'infrarouge (Zhou et al. 2011) et dans le micro-onde (Aires et al. 2001). D'un autre côté, d'autres algorithmes utilisent des

---

modèles de transfert radiatif pour inverser mathématiquement l'équation (Pequignot 2006; Prigent et al. 2006)<sup>1</sup>.

Les méthodes statistiques sont rapides, une fois calibrées. Cependant, certains phénomènes locaux peuvent influencer les restitutions si ces derniers ne sont pas bien considérés dans la phase d'apprentissage. Les aérosols, notamment au-dessus du Sahara perturbent fortement de tels modèles dans l'infrarouge. Les deux approches peuvent être complémentaires. Par exemple, Aires et al. (2001) utilisent une première ébauche et une climatologie restituées physiquement pour ensuite calculer une restitution statistique. Ils utilisent ici des réseaux de neurones pour les restitutions, mais une telle approche fonctionnerait avec des modèles d'assimilation au sein des NWP.

La sensibilité du rayonnement à la température peut être comparée quantitativement à sa sensibilité en émissivité (Hulley and Hook 2009). Une modification de 1,5 K de la température de surface équivaut à une variation de  $2 \times 10^{-2}$  en émissivité dans l'infrarouge. Les recommandations pour les instruments demandent une précision de l'ordre de 0,5 K en température de surface. Il faut alors connaître l'émissivité avec une précision d'au moins  $10^{-2}$ . Cet objectif n'est pas encore atteint aujourd'hui. Le fait que la température et l'émissivité de la surface soient tant liées dans le rayonnement de la surface rend difficile la restitution de l'un sans une bonne connaissance de l'autre (Vogel et al. 2011). Une erreur en température peut se compenser par une erreur sur l'émissivité de surface et inversement. C'est pourquoi, souvent, ces deux variables sont restituées conjointement, que ce soit dans l'infrarouge (Wan and Li 1997) ou dans le micro-onde (Aires et al. 2001).

Pour illustrer cela, nous avons mesuré la sensibilité moyenne des différents canaux de IASI à des modifications de l'émissivité ou de la température de surface. Pour ce faire, nous avons effectué des calculs de transfert radiatif, en modifiant l'un ou l'autre des paramètres. La sensibilité moyenne retrouvée sur 100.000 situations atmosphériques est présentée sur la Figure 2.1. La courbe bleue correspond à la modification moyenne des températures de brillance de IASI pour une augmentation de l'émissivité de 0,1 sur tout le spectre. La courbe verte présente les mêmes statistiques pour une augmentation de 5 K de la température de surface. Nous ne représentons pas, ici, la bande 3 de IASI qui a un bruit instrumental trop important pour pouvoir être exploitable dans de bonnes conditions. Une moyenne glissante sur 5 canaux a été utilisée pour lisser la courbe afin de gagner en visibilité. On voit apparaître de façon très claire les régions fenêtres du spectre infrarouge, entre  $770$  et  $1250 \text{ cm}^{-1}$  et au delà de  $2000 \text{ cm}^{-1}$ . Dans ces régions, l'influence des paramètres de surface est très importante et surtout très peu différenciable par paramètre. On voit bien qu'en couplant les deux courbes, il serait compliqué de dissocier les deux effets. C'est pourquoi nous choisissons de travailler sur une restitution simultanée de la température et de l'émissivité de surface.

---

1. Il existe également des méthodes variationnelles qui permettent, par itération successives, de déterminer les caractéristiques de surface. De telles méthodes sont notamment utilisées au sein des NWP.

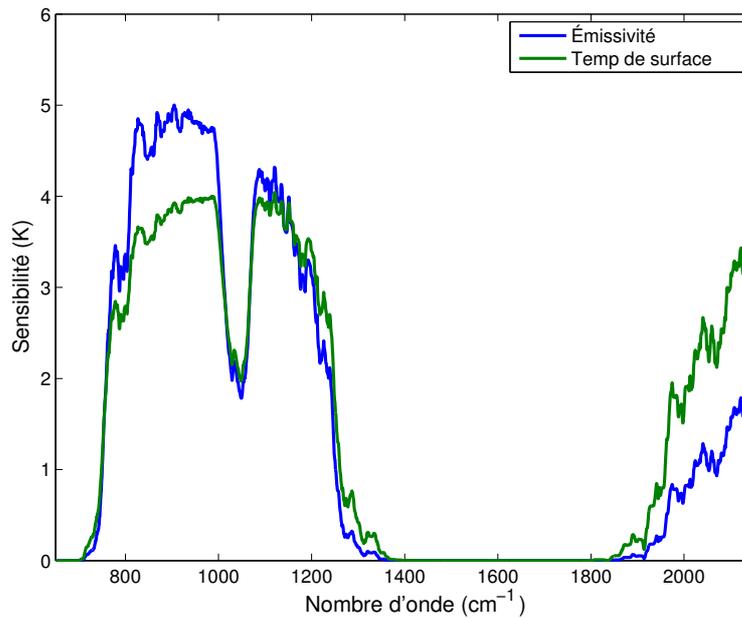


FIGURE 2.1 – Sensibilité moyenne des canaux IASI à une perturbation de 0,1 en émissivité (en bleu) et de 5 K en température de surface (en vert).

Dans le cadre de cette étude, nous disposons déjà d'un interpolateur d'émissivités micro-ondes en fréquence, en angle de visée et en polarisation (Aires et al. 2011b). Son principe est simple (Prigent et al. 2008) : après une paramétrisation de la dépendance en fréquence des émissivités micro-ondes, il peut calculer l'émissivité à une fréquence et à un angle de visée désirés, à partir d'un atlas d'émissivités dérivé d'observations satellites à des fréquences et des angles de visée donnés. Nous tâcherons dans cette étude de développer un interpolateur avec un principe similaire, mais fonctionnant dans l'infrarouge. Cet interpolateur servira de première ébauche à une restitution conjointe de la température et de l'émissivité de la surface. Nous nous attacherons ensuite à valider et à rendre opérationnel l'algorithme mis en place. Ce chapitre a fait l'objet d'une publication : Paul et al. (2012).

## 2.1 L'émissivité

Comme présenté dans la Section 1.1.1 (page 7), l'émissivité correspond au rapport entre l'énergie rayonnée par le corps considéré et l'énergie qu'il rayonnerait s'il était un corps noir. L'émissivité caractérise donc la capacité à rayonner d'un corps, à une longueur d'onde donnée. Si l'émissivité d'un corps vaut 1, alors il s'agit d'un corps noir (qui réémet toute l'énergie incidente). Inversement, plus l'émissivité sera faible, plus le rayonnement émis sera faible.

Dans le cadre de la télédétection par satellite, l'émissivité considérée est une émissivité correspondant à la résolution du sondeur, tant d'un point de vue spectral que spatial.

Spectralement, il faut prendre en considération la gamme de longueur d'onde sur laquelle chaque canal du sondeur considéré mesure. Cette plus ou moins grande largeur de bande peut impliquer des différences entre l'émissivité mesurée en laboratoire et celle perçue par le satellite. La résolution spatiale du satellite doit être prise en compte également. L'émissivité du point de vue du sondeur correspond à la moyenne de l'émissivité sur tout le champ de visée. Là encore, les mesures en laboratoire peuvent ne pas coïncider compte tenu de l'hétérogénéité de la surface à mesurer.

### 2.1.1 Émissivité de surface dans le micro-onde

Dans le domaine micro-onde, l'émissivité des continents est globalement supérieure à celle des océans. La modélisation de l'émissivité océanique est, aujourd'hui, arrivée à un niveau de maturité qui offre des résultats satisfaisants. Le modèle couramment utilisé, FASTEM ([English and Hewison 1998](#)), permet d'estimer l'émissivité des océans à partir de l'information du vent en surface, de la salinité de l'eau et de la température ([Liu et al. 2011](#)). L'émissivité océanique dans le micro-onde reste globalement inférieure à 0,6. Au-dessus des continents, l'émissivité est plus proche de l'unité.

[Prigent et al. \(2006\)](#) a montré que l'émissivité micro-onde dépendait surtout de la couverture végétale des sols. En effet, l'émissivité du sol dépend directement des propriétés diélectriques et de la rugosité du sol. Le fort contraste entre la constante diélectrique des végétaux et celle du sol engendre des différences d'émissivité. La dépendance de l'émissivité à la rugosité du sol est plus complexe à analyser car elle est fortement dépendante de l'angle de visée. On constate cependant que plus la rugosité du sol est importante, plus la diffusion est importante et plus l'émissivité du sol est élevée. Il est important de suivre également l'évolution de la différence d'émissivité suivant la polarisation du signal (horizontale ou verticale). Ici, une rugosité plus importante entraîne une diminution de la différence d'émissivité entre les deux polarisations.

Si la présence d'eau dans le sol aura tendance à faire décroître son émissivité et à augmenter la différence entre les deux polarisations, il est complexe de distinguer qui, de l'humidité du sol ou du taux de végétation du sol, en est le principal responsable. La corrélation entre le taux de végétation du sol et son taux d'humidité mélange leurs contributions à l'émissivité. Toutefois, la variation de la différence d'émissivité entre les deux polarisations est souvent utilisée pour restituer l'humidité du sol ([Lakshmi et al. 1997](#); [Vinnikov et al. 1999](#); [Ridder 2003](#)).

De nombreuses études ont prouvé qu'une bonne connaissance de l'émissivité du sol dans l'infrarouge permet d'avoir des informations sur sa composition ([Amer et al. 2010](#)). De plus en plus d'études mettent également en avant la corrélation entre le type de roche qui compose les sols arides et l'émissivité micro-onde de la surface. Le micro-onde a longtemps été utilisé pour ses qualités de pénétration dans le sol afin d'explorer sous la surface des zones désertiques, notamment à 900 MHz ([Grandjean et al. 2006](#)). Il a aussi été prouvé que l'émissivité

d'un sol aride dans le micro-onde avait des applications géologiques intéressantes, grâce à la dépendance de l'émissivité au type de roche sous-jacent (Prigent et al. 1999). Jiménez et al. (2010) a même établi un parallèle entre une carte géologique de l'Afrique du Nord et une carte d'émissivités micro-ondes. Il reste cependant à mesurer en laboratoire les propriétés diélectriques des matériaux de 1 à 600 GHz. Ces mesures sont déjà largement répandues en infrarouge, elles commencent seulement en micro-onde et manquent cruellement pour mieux interpréter l'émissivité micro-onde.

## TELSEM

Prigent et al. (2008) ont mis en place une paramétrisation de l'émissivité de 19 à 100 GHz en angle, fréquence et polarisation. Cette paramétrisation a été effectuée grâce à des restitutions d'émissivité, à partir de mesures de SSM/I (Special Sensor Microwave/Imager). Elle a permis de mettre au point un outil plus général appelé TELSEM (Aires et al. 2011b), permettant d'avoir accès à l'émissivité du sol sur une gamme spectrale plus large. Il s'agit d'une généralisation de la paramétrisation de l'émissivité.

Les données sont disponibles sous la forme de produits mensuels à la fréquence, à la polarisation et à l'angle de visée voulus sur une grille "equal-area" de  $0,25^\circ \times 0,25^\circ$  à l'équateur (voir Section 2.1.3, page 49). Cet outil nous servira, par la suite, pour définir les émissivités micro-ondes correspondant à AMSU-A et MHS. La précision de TELSEM, estimée à 0,02 sur les surfaces non-enneigées, nous a incités à développer un algorithme de restitution de l'émissivité de surface uniquement dans l'infrarouge. Cependant, la méthode qui est utilisée pour la restitution des émissivités infrarouge pourrait être appliquée aux données dans le micro-onde.

### 2.1.2 Émissivité de surface dans l'infrarouge

A l'inverse du domaine micro-onde, l'émissivité dans l'infrarouge au-dessus des océans est généralement plus élevée qu'au-dessus des continents. La faible variabilité de l'émissivité infrarouge de l'eau (quasi-constante à 0,98) implique que les océans sont généralement considérés comme uniformes en émissivité. Au-dessus de surfaces continentales, certains types de roche se distinguent par des signatures spectrales très remarquables. Ces signatures spectrales des roches terrestres dans l'infrarouge ont été longuement étudiées en laboratoire (Salisbury and D'Aria 1992, 1994).

Les structures les plus notables sont celles des silicates (il s'agit de la majorité des roches magmatiques et des quartz, ils composent 97% de la croûte terrestre). Ils présentent une grande bande d'absorption entre 1100 et 1150  $\text{cm}^{-1}$ . Cette bande d'absorption est due à leur teneur en quartz. La silice ( $\text{SiO}_2$ ) présente plusieurs bandes de vibration en élongation (variation de la longueur d'une liaison inter-atomique). Ces bandes ont été mesurées en laboratoire à 770, 1110 et 2500  $\text{cm}^{-1}$ . Ces fines bandes d'absorption sont particulièrement marquées spectralement, car l'effet Reststrahlen augmente leur largeur spectrale (Griffiths

1983). Cet effet est dû à un phénomène de réflexion interne au matériau, causé par les changements internes d'indice de réfraction à des longueurs d'onde proches des bandes d'absorption. Ainsi, ces trois structures spectrales sont bien marquées sur l'émissivité des sols contenant des silicates, même à l'échelle globale (Zhou et al. 2003).

Un autre matériau présentant une signature spectrale identifiable est le carbonate (composant de nombreuses roches sédimentaires). Les bandes d'absorption, dues aux vibrations de la liaison C-O dans l'anion carbonate, sont mesurées autour de 710, 910 et 1530  $\text{cm}^{-1}$ . Ici encore, l'effet Reststrahlen augmente leur largeur spectrale et les rend détectables à l'échelle globale (Chedin et al. 2004).

Les structures spectrales liées à la végétation sont très peu marquées comparées aux signatures des différentes roches précitées. L'émissivité d'une zone végétalisée reste supérieure à 0,9 tandis qu'au niveau d'une bande d'absorption des silicates, elle peut descendre jusqu'à 0,6. La variabilité spatiale de l'émissivité infrarouge est donc fortement marquée par les silicates, carbonates, ou, à plus faible échelle, d'autres roches. Au niveau de la neige, le comportement de l'émissivité est plus compliqué à analyser, compte tenu de la forte variabilité de la densité de la neige, de la forme des cristaux ou de la présence de glace en surface. Cependant, les variations spatiales restent moins marquées que les fortes signatures des silicates.

L'importante variabilité spatiale de l'émissivité infrarouge nous a incités à mettre au point un algorithme à même de restituer l'émissivité de surface dans l'infrarouge, à la résolution spectrale et spatiale de IASI. Il nous permettra de restituer des profils atmosphériques au-dessus des continents. En effet, les variations de la quantité de végétation du sol entraînent d'importantes modifications de son émissivité, notamment dans les bandes d'absorption des silicates.

### 2.1.2.1 Bases de données d'émissivités mesurées en laboratoire

#### Base de données de l'UCSB

Cette base de données est constituée de mesures en laboratoire de l'émissivité de différents matériaux ou d'échantillons de sols. Elle a été rassemblée par le Dr. Zhengming Wan à l'Institute for Computational Earth System Science à l'Université de Californie à Santa Barbara (UCSB) (<http://www.icesb.ucsb.edu/modis/EMIS/html/em.html>). L'émissivité d'un matériau plat est déterminée par la mesure de sa réflectance en utilisant un spectromètre TIR (Transformed Infrared) associé à une sphère d'intégration. La mesure de leur réflectance est ensuite convertie en émissivité hémisphérique directionnelle en utilisant la loi de Kirchhoff :  $\varepsilon = 1 - R$  (voir Section 2.1, page 40). Cette base de données apporte des informations très intéressantes sur l'émissivité des sols. Cependant, à l'échelle d'un satellite avec un champ de vue de l'ordre de la dizaine de kilomètres (voir Section 1.3, page 25), les surfaces considérées sont généralement plus complexes que les matériaux bruts décrits dans la base. Elles peuvent avoir différentes rugosités de surface ou mélanger plusieurs maté-

riaux, ce qui influencera leur émissivité. Celle-ci sera alors différente des mesures effectuées en laboratoire.

### Base de données ASTER

Cette autre base de données contient des mesures en laboratoire de l'émissivité rassemblées par le Jet Propulsion Laboratory à Pasadena, l'Université John Hopkins à Baltimore et le "United States Geological Survey" basé à Reston (<http://speclib.jpl.nasa.gov>). Différents types de spectres sont inclus, des roches sous différents états (solides, émietées ou même en poudre) jusqu'aux différents types de végétaux (Salisbury et al. 1994; Baldrige et al. 2009).

Cette base de données complète la précédente. Les spectres de silicates de cette base sont particulièrement importants car la base de données MODIS UCSB n'en contient que très peu.

#### 2.1.2.2 L'émissivité restituée à partir de MODIS

Nous utiliserons dans la suite de cette étude les produits MODIS MYD11 (Wan and Li 1997; Wan 2008). MODIS (Moderate Resolution Imaging Spectroradiometer) est une série d'instruments imageurs par satellite. Ils observent le rayonnement de la Terre dans 36 bandes spectrales entre 694 et 25.000  $\text{cm}^{-1}$ . Leur résolution au sol varie de 250 m à 1 km, en fonction de la fréquence. Ils ont été conçus pour effectuer des mesures à grande échelle de la couverture nuageuse, du bilan radiatif, des courants océaniques ou de la basse atmosphère, ainsi que de leurs variations climatiques.

Les produits MYD11 consistent en des moyennes mensuelles des émissivités restituées à partir de MODIS à 833,3; 909,1; 1162,8; 2500; 2564 et 2631,6  $\text{cm}^{-1}$ . L'algorithme de restitution est basé sur la différence des mesures de jour et de nuit (uniquement en ciel clair).

Considérons un point de la surface de la Terre, survolé une fois de nuit et une fois de jour. Il y a donc douze équations correspondant aux douze équations de transfert radiatif pour chacune des six longueurs d'onde de jour et de nuit. On se place bien sûr dans une situation en ciel clair. Le rayonnement infrarouge ne traversant pas les nuages, la surface n'est pas visible du point de vue d'un capteur infrarouge lorsqu'il y a des nuages. L'atmosphère est considérée comme transparente à ces longueurs d'onde (le coefficient de transmission atmosphérique  $\tau$  vaut 1), il ne reste donc que le terme de surface des équations de transfert radiatif (voir Section 1.1.2, page 12). En indexant par 1, 2, 3, 4, 5 et 6 les différentes longueurs

d'onde considérées, obtient alors le système d'équation suivant :

$$\begin{aligned}
 I_1^{jour} &= \varepsilon_1^{jour} \times B(T_{surf}^{jour}) & I_1^{nuit} &= \varepsilon_1^{nuit} \times B(T_{surf}^{nuit}) \\
 I_2^{jour} &= \varepsilon_2^{jour} \times B(T_{surf}^{jour}) & I_2^{nuit} &= \varepsilon_2^{nuit} \times B(T_{surf}^{nuit}) \\
 I_3^{jour} &= \varepsilon_3^{jour} \times B(T_{surf}^{jour}) & I_3^{nuit} &= \varepsilon_3^{nuit} \times B(T_{surf}^{nuit}) \\
 I_4^{jour} &= \varepsilon_4^{jour} \times B(T_{surf}^{jour}) & I_4^{nuit} &= \varepsilon_4^{nuit} \times B(T_{surf}^{nuit}) \\
 I_5^{jour} &= \varepsilon_5^{jour} \times B(T_{surf}^{jour}) & I_5^{nuit} &= \varepsilon_5^{nuit} \times B(T_{surf}^{nuit}) \\
 I_6^{jour} &= \varepsilon_6^{jour} \times B(T_{surf}^{jour}) & I_6^{nuit} &= \varepsilon_6^{nuit} \times B(T_{surf}^{nuit})
 \end{aligned}$$

La température de surface  $T_{surf}$  est la même pour les différentes longueurs d'onde mais est différente suivant le moment de la mesure. En considérant que la variation journalière de l'émissivité est négligeable (les émissivités de jour et de nuit sont égales), il y a huit inconnues et non 14 : les 6 émissivités et la température de surface de jour et de nuit<sup>2</sup>. Le système est composé de 12 équations pour 8 inconnues, il est inversible mathématiquement car il a plus d'équations que d'inconnues.

Cet algorithme fonctionne uniquement en ciel clair puisque les nuages sont opaques au rayonnement infrarouge. C'est pourquoi dans certaines zones du globe, notamment les tropiques qui sont constamment sous les nuages, la moyenne mensuelle des restitutions peut être incomplète.

La variabilité annuelle de ce produit mensuel est présentée sur la Figure 2.2. Sur cette figure, la variabilité annuelle n'est pas calculée partout avec le même nombre d'occurrences. En effet, certains pixels restent sous les nuages durant tout un mois et n'ont donc pas d'émissivité MODIS associée. De plus, les points proches des pôles n'ont pas l'alternance jour/nuit nécessaire à l'algorithme de restitution<sup>3</sup>. Ainsi, la carte de variabilité n'est pas complète.

L'émissivité des surfaces présentant une saisonnalité de la végétation est plus variable que celle des déserts, mais la variation temporelle reste faible. Les plus fortes variations sont observées dans les zones de transition semi-arides. On voit par exemple que la région du sub-Sahel présente une importante variabilité de son émissivité. Elle est très aride pendant une partie de l'année et recouverte de végétation à la saison humide.

2. Cette approximation néglige les variations journalières de l'émissivité, qui peuvent être importantes.

3. L'algorithme nécessite cette alternance jour/nuit et non pas simplement deux survols de la zone afin d'avoir un vrai contraste entre les deux observations.

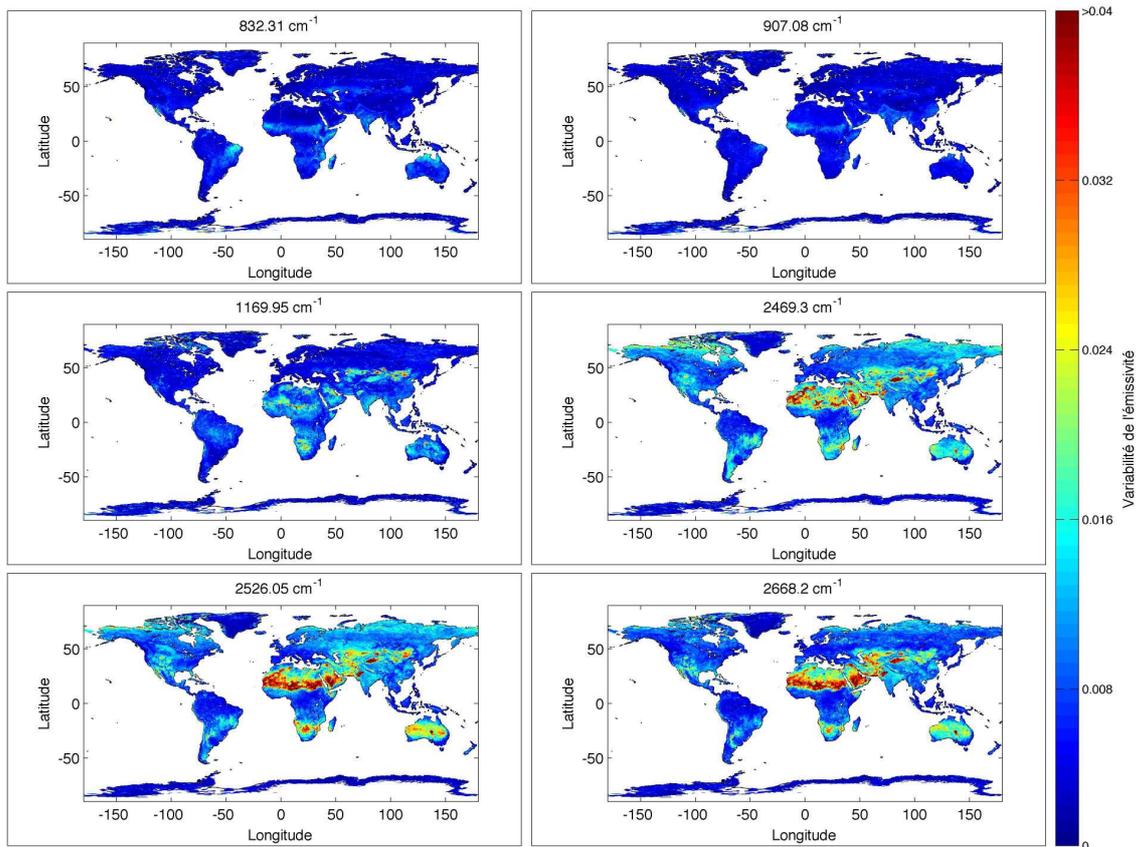


FIGURE 2.2 – Carte représentant la variabilité annuelle des émissivités MODIS. Les données utilisées sont les produits MYD11 de l’année 2007.

Afin de pouvoir comparer les résultats obtenus avec ceux de [Seemann et al. \(2008\)](#) (voir Section 2.1.2.3, page 46), la version 4.1 de l’année 2007 est utilisée.

### 2.1.2.3 La base de données UWIRemis

Les mesures MODIS fournissent six émissivités de surface (voir Section 2.1.2.2, page 44) de façon globale avec une résolution de  $0.05^\circ \times 0.05^\circ$ . De ces six émissivités, [Seemann et al. \(2008\)](#) déduisent dix émissivités à dix longueurs d’onde différentes : 699,3, 826,5 ; 925,9 ; 1075,3 ; 1204,8 ; 1315,8 ; 1724,1 ; 2000 ; 2325,6 et 2777,8  $\text{cm}^{-1}$ . Une forme spectrale caractéristique est définie grâce aux émissivités de laboratoire (voir Section 2.1.2.1, page 43) : une lente décroissance de l’émissivité à partir de 714,3  $\text{cm}^{-1}$  jusqu’au minimum de l’émissivité dans la bande d’absorption du quartz à 1162,8  $\text{cm}^{-1}$ , suivi d’une rapide remontée au début de la bande d’absorption de la vapeur d’eau, une lente diminution et enfin une partie fortement décroissante entre 2000 et 2500  $\text{cm}^{-1}$ .

La Figure 2.3 présente des exemples de spectres d’émissivité à diverses résolutions. Les étoiles correspondent à l’émissivité restituée à partir de MODIS (le produit MYD11). Les lignes en pointillés représentent les 10 points caractéristiques considérés par Seemann, reliés

les uns aux autres afin de mettre en avant l'évolution graphique du spectre précédemment décrite. Les traits continus représentent l'émissivité à la résolution IASI. Les spectres rouges correspondent à un sol couvert de végétation, avec une émissivité forte quasi constante. Les spectres verts correspondent à un sol partiellement aride et les spectres bleus correspondent à un sol aride.

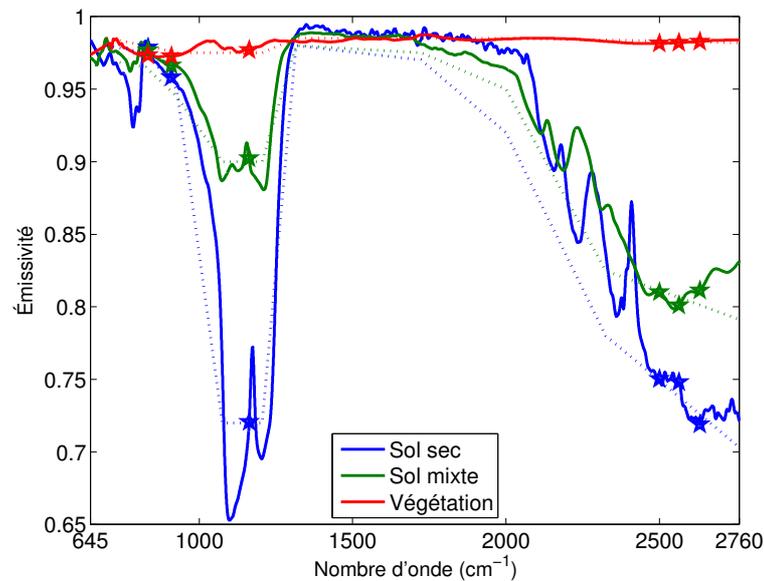


FIGURE 2.3 – Spectres d'émissivité d'un sol sec (bleu), d'un sol mixte (vert) et d'un sol avec de la végétation (rouge). Chaque spectre complet (issu de mesures en laboratoire) est représenté en trait continu, les étoiles représentent le spectre à la résolution de l'instrument MODIS, et la courbe en pointillés est la représentation linéaire de Seemann.

On peut remarquer les signatures des silicates à  $1100$  et  $2500\text{ cm}^{-1}$ . Ces structures spectrales sont encore plus remarquables sur les spectres bleus correspondant à un sol aride. La résolution spectrale équivalente à IASI est obtenue à partir de ces dix émissivités grâce à une régression linéaire (Borbis et al. 2007; Borbis and Ruston 2010). Une analyse en composantes principales (voir Annexe B.1, page 211) est calculée sur les spectres d'émissivité haute résolution mesurés en laboratoire (voir Section 2.1.2.1, page 43). Une régression linéaire est utilisée pour transformer les dix émissivités extraites à partir de MODIS en composantes de l'analyse en composantes principales sur les émissivités à haute résolution spectrale. Après divers tests, seuls six composantes sont utilisées dans cette régression. Il est possible d'obtenir des structures spectrales même entre les 10 points utilisés par Seemann grâce à l'utilisation de l'analyse en composantes principales (ACP) car les composantes restituées, même linéairement à partir des 10 points, contiennent des structures spectrales plus fines.

En respectant la notation de Seemann et al. (2008) et Borbis and Ruston (2010), cette base de données sera appelée "la base UWiremis".

#### 2.1.2.4 L'émissivité restituée à partir de IASI par D. Zhou à la NASA

Zhou et al. (2011) restituent des spectres d'émissivité directement à partir des mesures IASI. Dans un premier temps, une base d'apprentissage est construite. Cette base est constituée de profils atmosphériques auxquels sont associés une température de surface et une émissivité de surface. Les émissivités de surface considérées sont issues des mesures en laboratoire présentées en Section 2.1.2.1 (page 43) et prennent en compte le type de surface associé au profil atmosphérique afin d'avoir une émissivité cohérente. Les mesures IASI associées à ces différents profils sont un mélange de données réelles et de données simulées en utilisant CRTM (Community Radiative Transfer Model, un outil similaire à RTTOV présenté en Section 1.4, page 33). Cette base d'apprentissage permet de calculer des coefficients de régression séparés en plusieurs catégories suivant la couverture nuageuse.

Chaque profil à restituer est d'abord calculé en utilisant la régression de chaque catégorie (claire, nuageuse ou mixte) et un critère de convergence sert alors à déterminer si la situation est claire ou nuageuse. Ensuite, la restitution de la température et de l'émissivité obtenue passe par un algorithme 1D-var afin d'être améliorée. L'avantage de cet algorithme est qu'il est relativement rapide par rapport à des approches plus physiques<sup>4</sup>.

Ce produit est distribué sous la forme de moyennes mensuelles pour des restitutions en ciel clair. Nous utiliserons les données de 2008. Ce jeu de données sera appelé par la suite les "données NASA".

#### 2.1.2.5 L'émissivité restituée à partir de IASI par le groupe ARA au LMD

Cette base de données a été construite par le groupe ARA, Analyse du Rayonnement Atmosphérique, du LMD, Laboratoire de Météorologie Dynamique (Pequignot 2006; Pequignot et al. 2008; Capelle et al. 2012). Un réseau de neurones est utilisé afin de restituer la température de surface à partir de mesures de AIRS (un instrument proche de IASI, voir Section 1.3.1, page 29) à des longueurs d'onde données. L'émissivité est déterminée pour les canaux "fenêtres" en utilisant directement l'équation de transfert radiatif (voir Section 1.1.2, page 12). Afin de pouvoir calculer tous les termes de l'équation de transfert radiatif, une connaissance préalable du profil atmosphérique correspondant est nécessaire. Ici, une base de référence est utilisée. Cette base est constituée de profils atmosphériques associés aux spectres de température de brillance de IASI correspondants. Cette base s'appelle TIGR (Thermodynamic Initial Guess Retrieval). Une recherche du plus proche voisin est effectuée à partir des températures de brillance de IASI et le profil atmosphérique correspondant est utilisé dans les calculs.

Afin d'obtenir la résolution spectrale de IASI (ou de AIRS) à partir des émissivités aux niveaux des canaux fenêtres sélectionnés, un algorithme de reconnaissance de forme est utilisé, basé sur les spectres d'émissivité mesurés en laboratoire (voir Section 2.1.2.1,

---

4. Un algorithme 1D-var direct avec une première ébauche trop éloignée de la vraie solution est plus lent à converger.

page 43). Cette base de données est disponible uniquement dans les tropiques et sous-tropiques ( $\pm 30^\circ$  de latitude). Elle sera appelée “la base ARA” dans la suite.

### 2.1.3 Mise en coïncidence des bases

Toutes les bases de données présentées ici sont projetées sur une grille “equal-area” à  $0,25^\circ \times 0,25^\circ$  afin de pouvoir comparer les différentes émissivités sans avoir de différences de résolution spatiale. Une grille “equal-area” correspond à un maillage de la terre où chaque maille a la même surface au sol. Ainsi, sur une projection Mercator classique (celle utilisée sur la carte 2.2), les mailles seront plus petites au niveau de l'équateur qu'au niveau des pôles. Une représentation sur un globe en trois dimensions montrerait un maillage régulier.

Nous utilisons ici une grille à  $0,25^\circ \times 0,25^\circ$ , c'est-à-dire qu'une maille située au niveau de l'équateur fait  $0,25^\circ \times 0,25^\circ$ . Le découpage en latitude est régulier tout les  $0,25^\circ$ , mais le découpage en longitude varie suivant la latitude pour avoir toujours la même surface de maille au sol. Cette grille représente, au sol, des carrés déformés d'environ 30 km de côté.

## 2.2 Construction d'une base de données d'émissivités infrarouges

L'objectif est de mettre en place un schéma de restitution de la température de surface et de l'émissivité à la résolution de IASI. De nombreux schémas de restitution des caractéristiques de surface sont déjà utilisés, certains d'entre eux ont été présentés précédemment (voir Sections 2.1.2.3, page 46, 2.1.2.4, page 48, et 2.1.2.5, page 48). Notre objectif est de tirer profit de tout ce travail déjà mis en place afin d'établir un nouvel algorithme encore plus performant, qui utilise les atouts de chacune de ces méthodes. Il s'agit alors de combiner des méthodes statistiques (l'interpolateur) et analytique (la restitution en elle-même) pour mettre en place un algorithme de restitution global et efficace.

Dans un premier temps, nous avons créé un interpolateur d'émissivités infrarouges pour obtenir un spectre d'émissivité à la résolution IASI à partir des produits MODIS. Cette émissivité servira, dans un deuxième temps, de première ébauche (de bonne qualité, l'interpolateur a une erreur inférieure à  $4 \times 10^{-3}$ ) pour un schéma de restitution analytique basé sur une inversion mathématique du transfert radiatif (présenté à la section suivante).

Notre interpolateur d'émissivité a été développé en parallèle de celui de [Seemann et al. \(2008\)](#). Il lui est similaire. Nous avons cependant décidé d'utiliser un réseau de neurones à la place de la régression linéaire, pour passer des émissivités MODIS aux composantes de l'ACP sur les émissivités. Afin de pouvoir calibrer le réseau de neurones, il est nécessaire de disposer d'une base d'apprentissage. Cette base de données doit mettre en parallèle les sorties du réseau, c'est-à-dire les composantes de l'ACP sur les spectres à haute-résolution (celle de IASI) et les entrées du réseau, c'est-à-dire les valeurs d'émissivité issues des restitutions de MODIS. Cette base servira également à calculer l'ACP sur les émissivités à haute-résolution

spectrale. Il s'agira ensuite d'effectuer l'apprentissage du réseau de neurones afin d'obtenir l'interpolateur d'émissivité.

### 2.2.1 Les réseaux de neurones

Un réseau de neurones est un outil de traitement de l'information qui est inspiré par la manière dont les systèmes nerveux biologiques, comme le cerveau, exploitent les informations (Lettvin et al. 1959). Il est composé d'un grand nombre d'éléments interconnectés (neurones) qui travaillent ensembles afin de résoudre un problème. Pour la télédétection, le réseau le plus couramment utilisé est le perceptron multicouche entièrement connecté (Hornik et al. 1989). Il peut se décomposer en trois parties. La première partie constitue la couche d'entrée, suivie des diverses couches cachées et enfin la couche de sortie. Chaque couche est composée de neurones connectés aux neurones de la couche précédente et de la suivante, comme présenté sur le schéma 2.4. Chacun de ces liens synaptiques est associé à un poids synaptique.

Chaque neurone effectue deux opérations :

- le neurone  $i$  calcule la somme pondérée  $h_i$  à partir de ses  $p$  entrées  $x_j$  :  $h_i = \sum_{j=1}^p w_{j,i}x_j$ , où  $w_{j,i}$  correspond au poids synaptique de l'entrée  $j$  pour le neurone  $i$ .
- il applique ensuite à  $h_i$  une sigmoïde d'activation  $\sigma$  et ajoute au total un biais éventuel  $b_i$  :  $y_i = \sigma(h_i) + b_i$ .

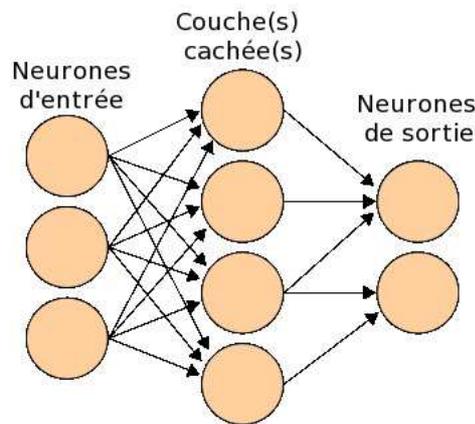


FIGURE 2.4 – Schéma d'un réseau de neurones.

D'un réseau à l'autre, la sigmoïde peut être différente (arctangente, fonction seuil...). Les fonctions d'activations utilisées sont, dans notre cas, des arctangentes. La non-linéarité de ces sigmoïdes rend le réseau capable de résoudre des problèmes non-linéaires. Son apprentissage (*i.e.*, le calcul des différents poids synaptiques) est effectué grâce à une fonction de poids qui somme les erreurs quadratiques entre les sorties du réseau et les sorties cibles définies dans la base d'apprentissage. Une descente de gradient permet de diminuer l'erreur sur les sorties du réseau en modifiant les différents poids synaptiques. La rétropropagation du gradient et

## 2.2. CONSTRUCTION D'UNE BASE DE DONNÉES D'ÉMISSIVITÉS INFRAROUGES

l'algorithme de descente du gradient (*i.e.*, la méthode Levenberg-Marquardt) est détaillée en Annexe B.2.2 (page 215) (Rumelhart et al. 1986).

Par la suite, les différents réseaux de neurones qui seront utilisés (même pour les restitutions) seront des perceptrons multicouches de même type. Il sera fait référence à cette brève description d'un réseau de neurones pour que le lecteur puisse comprendre mais une présentation plus approfondie est fournie en Annexe B.2.2 (page 215).

Dans le cadre de l'interpolateur d'émissivités infrarouges, nous utiliserons un perceptron multicouche constitué d'une seule couche cachée, comprenant 20 neurones. Cet architecture a été choisie ici car elle offrait les meilleurs résultats.

### 2.2.2 Création d'une base de première ébauche d'émissivités hyperspectrales à partir des observations MODIS

#### 2.2.2.1 Une base de première ébauche indépendante de IASI

Comme expliqué en préambule, il s'agit de créer une base d'apprentissage pour l'interpolateur d'émissivité. Ce dernier a, comme entrées, les émissivités restituées à partir de MODIS et comme sorties, des composantes d'une ACP sur les émissivités à haute-résolution spectrale. Il est important que l'interpolateur soit indépendant des mesures utilisées pour la restitution (*i.e.*, les données réelles IASI). Il nous permet d'obtenir une première ébauche qui nous servira à effectuer la restitution. Une corrélation entre les deux peut entraîner une très forte augmentation des erreurs engendrées, les erreurs de l'un pouvant résonner avec celles de l'autre (Rodgers 2000). L'objectif est donc ici de créer une base d'émissivité de surface à la résolution de IASI qui représente au mieux la variabilité naturelle de l'émissivité. Nous utiliserons uniquement les spectres d'émissivité mesurés en laboratoire et non ceux déjà restitués par d'autres méthodes afin de rendre cette méthode indépendante des autres schémas de restitution et de pouvoir ensuite comparer les différents résultats de façon décorrélée, ce qui est nettement plus intéressant. Par la suite il suffira d'extraire, de ces spectres à haute-résolution, les bandes spectrales correspondant aux longueurs d'onde de MODIS, pour obtenir des spectres à la résolution voulue.

#### 2.2.2.2 Une classification de surface

La classification de surface que nous utiliserons est la classification de l'IGBP-DIS (International Geosphere Biosphere Program Data and Information System). Grâce à l'expertise cartographique de l'U.S. Geological Survey et de l'European Commission's Joint Research Center, l'IGBP-DIS a généré une base de données de mesures du radiomètre AVHRR (Advanced Very High Resolution Radiometer, radiomètre mesurant dans 6 bandes dans le rouge, le proche infrarouge et l'infrarouge). Plus de 4,4 To de données, provenant de 23 stations de réception, ont été recueillies, rassemblées et traitées. DISCover a été créé à partir de ces données. Il consiste en une classification des surfaces continentales en 17 catégories. Ces

catégories ont été conçues pour représenter les éléments généraux de structure des couverts végétaux et des terres. Le produit final, DISCover, a été basé sur des différences mensuelles normalisées d'indice de végétation composite entre 1992 et 1993. Les catégories considérées sont les suivantes :

- |  |  |
|--|--|
| 1. Eau                                       | 10. Savane                                     |
| 2. Forêt résineuse                           | 11. Prairie                                    |
| 3. Forêt de feuillus à feuilles persistantes | 12. Marécage permanent                         |
| 4. Forêt de conifères à feuilles caduques    | 13. Terre cultivée                             |
| 5. Forêt de feuillus à feuilles caduques     | 14. Zone urbaine ou bâtie                      |
| 6. Forêt mixte                               | 15. Mosaïque de terres cultivées ou naturelles |
| 7. Broussaille dense                         | 16. Neige et glace                             |
| 8. Broussaille clairsemée                    | 17. Désertique ou quasi désertique             |
| 9. Savane boisée                             |  |

Afin de prendre en compte les variations saisonnières des surfaces (car le produit DISCover est annuel), nous utilisons une climatologie mensuelle de neige et de glace. La classification mensuelle des surfaces présente donc une variabilité mensuelle uniquement dans les zones susceptibles d'être recouvertes par la neige.

### 2.2.2.3 Création de la base à partir des observations MODIS

Chaque pixel de chaque mois de la grille "equal-area" à  $0,25^\circ \times 0,25^\circ$  est associé à un type de surface particulier et à un spectre MODIS d'émissivité (6 longueurs d'onde différentes). On cherche à associer chacun de ces pixels à un spectre d'émissivité complet (à la résolution IASI), mesuré en laboratoire. Pour ce faire, on associe à chaque classe de l'IGBP-DIS les spectres mesurés en laboratoire pouvant correspondre à ce type de sol. Ce regroupement reste assez large pour permettre à l'algorithme de compenser certaines erreurs de classification. Pixel par pixel, on cherche alors les spectres de laboratoire, parmi ceux correspondant au type de surface, qui s'approchent le plus des 6 émissivités MODIS. Pour cela, on minimise l'écart quadratique moyen sur les 6 bandes spectrales. Les spectres ainsi obtenus sont moyennés, de façon pondérée par leur distance aux émissivités MODIS. Prendre en compte les 5 plus proches voisins des 6 émissivités MODIS permet de mieux représenter la variabilité naturelle des émissivités.

Du fait de la non-disponibilité de restitution d'émissivité MODIS au-dessus des pôles (due à l'absence d'alternance jour/nuit), il n'y a aucune données à ce niveau. Cependant, la prise en compte de bases mensuelles liées à une climatologie de la glace et la neige permet d'inclure de nombreux spectres de glace et de neige, compensant l'absence des pôles dans la base de données. Seuls quatre mois sont utilisés afin de réduire la taille totale de la base de données. Nous ne conservons que les mois de janvier, d'avril, de juillet et d'octobre, ce qui nous laisse tout de même plus de 300.000 spectres dans la base de données complète.

### 2.2.3 Analyse en composantes principales des spectres d'émissivité

Nous cherchons ici à réduire la dimension des émissivités hyperspectrales. L'objectif est de caractériser leur variabilité à partir d'un nombre minimum de variables.

#### 2.2.3.1 Exemple simple d'analyse en composantes principales

Prenons un exemple simple : considérons une base de données, où chaque échantillon est constitué de deux variables  $(X_1, Y_1)$ . Cette base de données est représentée sur la Figure 2.5. L'analyse en composantes principales conduirait au changement de coordonnées de  $(X_1, Y_1)$  en  $(X_2, Y_2)$ . On peut remarquer qu'avec cette nouvelle projection, la variable  $X_2$  seule est déjà une bonne caractérisation de la donnée. Les variations en  $Y_2$  ne semblent être que du bruit. On peut donc caractériser cette base uniquement par la variable  $X_2$ .

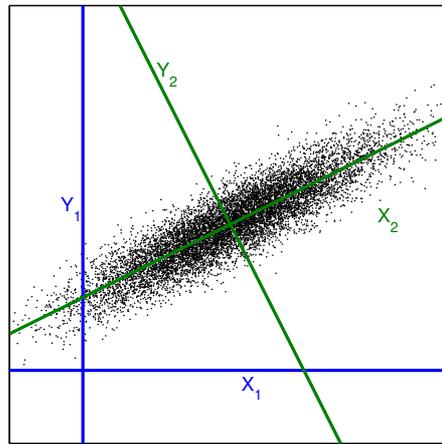


FIGURE 2.5 – Exemple de représentation graphique d'une base de données de dimension 2. Elle peut être projetée sur les axes  $(X_1, Y_1)$  ou  $(X_2, Y_2)$ .

#### 2.2.3.2 Méthodologie

L'Analyse en Composantes Principales (ci-après désignée ACP) permet de compresser des données de grande dimension en limitant la perte d'informations (Jolliffe 2002). L'avantage de cette méthode est qu'elle permet également de réduire l'importance des faibles variations des variables par rapport aux autres (*i.e.*, augmenter le rapport signal sur bruit). L'ACP consiste en une projection des données sur une nouvelle base orthonormée de façon à diagonaliser la matrice de covariance des variables. Il suffit alors de prendre un nombre restreint de composantes de la nouvelle base pour compresser les données. Les vecteurs qui forment la nouvelle base orthonormée sont appelés vecteurs propres, il s'agit de combinaisons linéaires des variables originales. Plus de précisions sont fournies à l'Annexe B.1 (page 211).

Nous effectuons donc une ACP sur les émissivités hyperspectrales infrarouges. Le principe est de calculer la matrice de covariance  $V = cov(\varepsilon)$ . Un algorithme de Cholesky permet

de rapidement diagonaliser cette matrice (Cholesky 1910). On obtient alors les différents vecteurs propres et les valeurs propres qui leur sont associées. Chacun de ces vecteurs propres (souvent appelés EOF, pour Empirical Orthogonal Function) représente une variation par rapport à l'émissivité moyenne, avec une structure spatiale particulière. Chacun décrit une partie de la variabilité spectrale des émissivités. Plus la valeur propre associée est grande plus cette variabilité est importante dans la base complète. Les vecteurs propres sont donc triés par valeur propre décroissante. Seuls les vecteurs propres associés aux plus grandes valeurs propres seront pris en compte. On a alors pour chaque spectre :

$$(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_{8461}) = C \cdot VP + \bar{\varepsilon} \quad (2.1)$$

La matrice  $C$  de taille  $(1 \times 8461)$  représente les nouvelles variables dont seulement une partie sera prise en compte. La matrice  $VP$  de taille  $(8461 \times 8461)$  représente la matrice de changement de base, elle contient donc les vecteurs propres projetés sur l'espace des longueurs d'ondes et triés par valeurs propres décroissantes ;  $\bar{\varepsilon}$  représente la moyenne de l'émissivité. Chaque composante est une anomalie par rapport à l'émissivité moyenne.

Au lieu de prendre en compte les 8461 vecteurs propres pour représenter le spectre d'émissivité, un nombre restreint  $N$  de composantes est sélectionné. Dans ce cas, la matrice  $C$  de l'équation (2.1) a pour dimensions  $(1 \times N)$  et la matrice  $VP$  de l'équation (2.1) est remplacée par  $VP_N$ , une matrice  $(N \times 8461)$ . On peut ensuite facilement revenir au spectre d'émissivité car la matrice  $VP$  est inversible. Par définition,  $VP$  est la matrice des vecteurs propres donc  $VP^{-1} = VP^T$ .

### 2.2.3.3 Résultats obtenus sur les spectres d'émissivité

Une étude des variations spectrales et spatiales des vecteurs propres aide à mieux comprendre l'ACP. La Figure 2.6 présente deux composantes principales, sous forme de spectre et de répartition géographique. La partie haute de la figure présente le spectre de la première composante sur la gauche et la carte de son impact sur la droite.

Cette étude est effectuée sur la base créée précédemment à partir des spectres de laboratoire (voir Section 2.2.2, page 51). On remarque sur le spectre la signature de silicates à  $1100 \text{ cm}^{-1}$ . Le signe des composantes sur la carte n'a pas vraiment d'intérêt, ce qu'il faut regarder est l'écart par rapport à zéro. On retrouve la signature spatiale des silicates au niveau des déserts de sable (dans la péninsule arabique et le Sahara.). Cette signature est la plus marquée dans la variabilité spectrale des émissivités infrarouges.

Si l'ACP permet de bien compresser des données, elle a l'inconvénient de mélanger les différents signaux pour maximiser la variabilité expliquée par chacune des composantes. Plus la composante considérée est d'ordre élevé, plus c'est un mélange de signatures physiques et d'anomalies par rapport aux premières composantes. Ainsi, il est difficile d'interpréter les composantes d'ordre élevé. La partie basse de la figure présente la même étude sur la cin-

## 2.2. CONSTRUCTION D'UNE BASE DE DONNÉES D'ÉMISSIVITÉS INFRAROUGES

quième composante. Il est complexe de donner une signification physique à cette composante, tant spectralement que spatialement. Elle présente de la variabilité tant sur les zones arides qu'humides. Cependant, sa plus ample diversité sur les zones couvertes de végétation tend à indiquer que cette composante comporte la signature de l'eau et de la végétation. Le mélange de signatures de l'ACP rend complexe l'analyse des diverses composantes (Aires et al. 2002b). Seules les fortes variabilités (ici les silicates) peuvent être facilement identifiées.

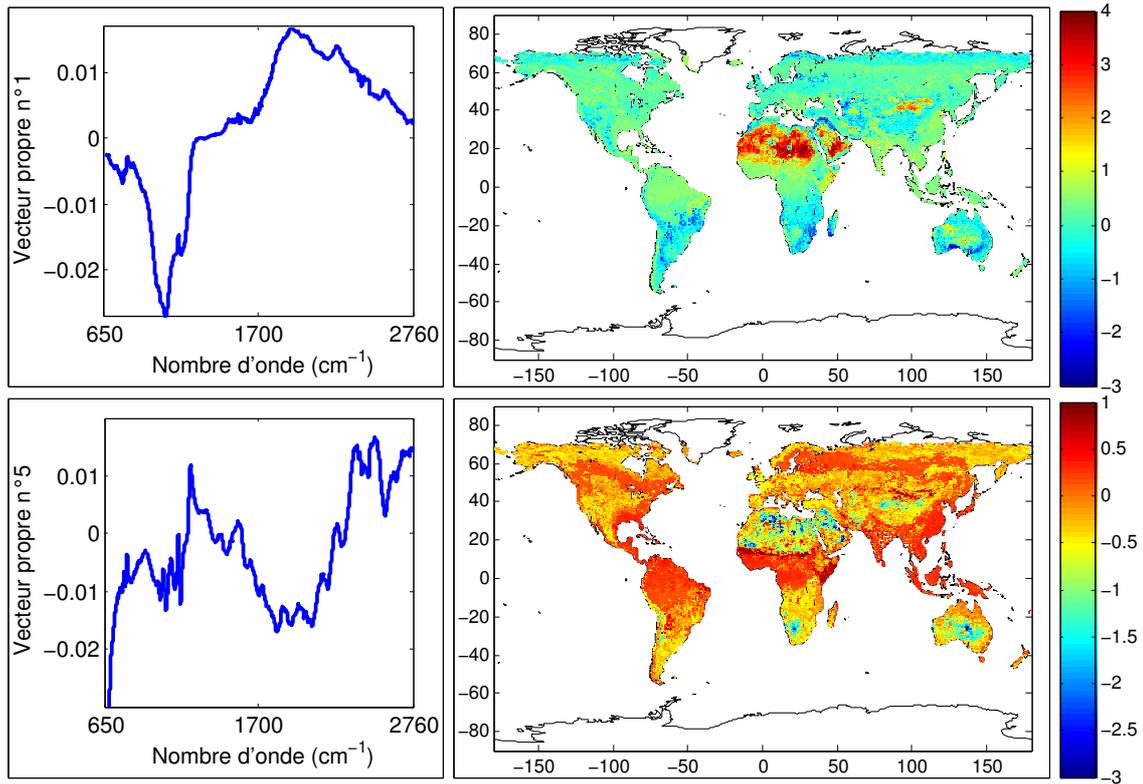


FIGURE 2.6 – Étude de deux composantes principales. En haut, la première composante, son spectre (à gauche) et sa répartition géographique (à droite). En bas, la même figure mais pour la cinquième composante.

### 2.2.3.4 Compression des émissivités infrarouges

Il reste désormais à étudier la capacité de compression de l'ACP. Pour cela, des calculs d'erreur suite à une compression suivie d'une décompression sont effectués. La base de données d'émissivités hyperspectrales est compressée sur un nombre  $N$  de composantes puis décompressée. On mesure ensuite la racine de l'erreur quadratique moyenne (aussi appelée RMS de l'erreur) entre les spectres originaux et les nouveaux. La Figure 2.7 représente spectralement les résultats obtenus pour différents nombres de composantes.

On remarque sur cette figure que l'erreur de compression est déjà faible en n'utilisant que 5 composantes, de l'ordre de  $2 \times 10^{-2}$ . Cette faible erreur de compression, avec un nombre restreint de composantes, s'explique par l'importante corrélation qui existe entre

les différentes émissivités des canaux voisins de IASI, ce qui réduit le nombre de degrés de liberté dans la variabilité de l'émissivité. Cette forte corrélation spectrale permet d'expliquer la variabilité de l'émissivité avec nettement moins de composantes que les 8461 canaux de IASI. À partir de 10 composantes utilisées, l'erreur descend en dessous de  $10^{-3}$ . Plus le nombre de composantes augmente, plus l'erreur diminue. L'émissivité que l'on souhaite obtenir à partir des composantes de l'ACP servira de première ébauche à la restitution.

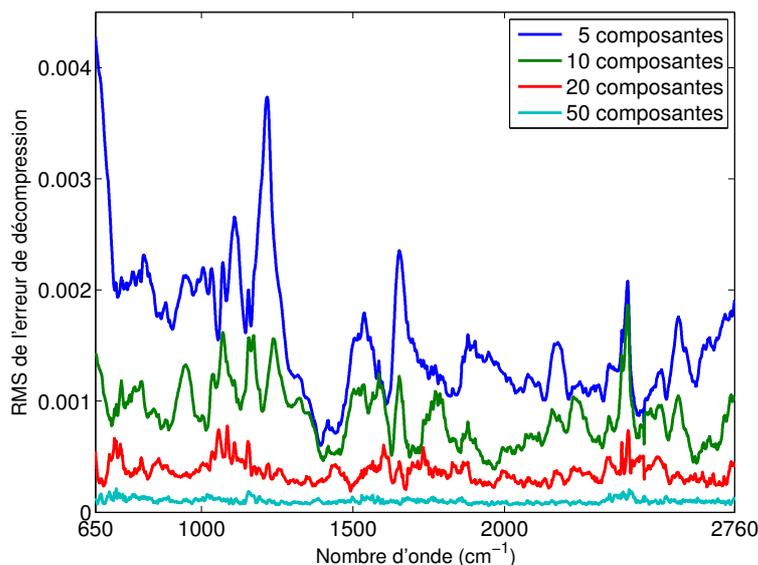


FIGURE 2.7 – RMS de l'erreur de compression puis décompression de l'émissivité sur différents nombres de composantes : 5 (en bleu foncé), 10 (en vert), 20 (en rouge) et 50 (en bleu clair).

Nous choisissons de ne prendre en compte que 10 composantes de l'émissivité. Ce nombre est un bon compromis dans la mesure où la précision que l'on recherche dans la restitution de l'émissivité est de l'ordre de  $10^{-2}$ , et limiter le plus possible le nombre de composantes augmente la stabilité et la rapidité de l'algorithme de restitution qui sera mis en place.

De plus, l'interpolateur d'émissivités infrarouges n'a, comme entrées, que 6 informations, correspondant aux 6 émissivités aux fréquences de MODIS. Prendre en compte un nombre trop élevé de composantes de l'ACP n'aurait alors pas beaucoup de sens dans la mesure où on cherche à les reconstituer avec seulement 6 informations. Restreindre les degrés de liberté des algorithmes améliore leur stabilité.

#### 2.2.4 Interpolation spectrale de l'émissivité infrarouge

Nous avons désormais en main tous les outils nécessaires à la création de l'interpolateur d'émissivité. Un réseau de neurones est mis en place, comme décrit dans la Section 2.2.1 (page 50). La base de données construite précédemment (voir Section 2.2.2, page 51) servira

## 2.2. CONSTRUCTION D'UNE BASE DE DONNÉES D'ÉMISSIVITÉS INFRAROUGES

de base d'apprentissage pour le réseau de neurones. Ce dernier a, comme entrées, les 6 émissivités extraites de la base hyperspectrale, correspondant aux longueurs d'onde MODIS, et comme sorties les 10 premières composantes de l'ACP sur les émissivités hyperspectrales. À partir de ces 10 composantes de l'ACP nous pouvons obtenir un spectre d'émissivité à la résolution de IASI. L'interpolateur neuronal nous permettra de passer d'émissivités à la résolution de MODIS à celle de IASI (de 6 à 8461 canaux).

La base de données constituée de 300.000 spectres est séparée en trois parties distinctes :

- **Une base d'apprentissage** à proprement parler qui est présentée au réseau et qui modifie les poids des différents neurones (240.000 situations) ;
- **Une base de validation** qui permet de vérifier au cours de l'apprentissage que le réseau ne fait pas de surapprentissage. Il s'agit, à chaque itération de l'apprentissage, de vérifier que l'erreur de restitution sur cette base est décroissante. Le fait que cette base de données ne soit pas présentée au réseau permet de vérifier le pouvoir de généralisation du réseau (30.000 situations) ;
- **Une base de test** qui permet de calculer l'erreur du réseau une fois l'apprentissage effectué. Séparer cette base de test de la base de validation empêche le réseau "d'apprendre" la base de validation. On a ainsi une meilleure estimation de l'erreur commise par le réseau de neurones (30.000 situations) en condition opérationnelle.

La Figure 2.8 présente les statistiques de la RMS de l'erreur de l'interpolateur sur la base de test. On remarque sur cette figure que l'erreur d'interpolation est supérieure à l'erreur de compression avec 10 composantes. Il serait inutile de chercher à représenter plus de composantes à partir de seulement 6 émissivités. Avec environ  $5 \times 10^{-3}$  contre  $10^{-3}$ , notre choix de ne prendre en compte que 10 composantes dans l'ACP est ainsi validé.

Sur l'axe des abscisses, on a représenté en vert les bandes spectrales correspondant aux 6 émissivités MODIS. Ces dernières sont représentées en fonction de leur largeur spectrale. MODIS est un instrument moins bien résolu spectralement que IASI. Un canal de MODIS correspond à la moyenne de plusieurs canaux IASI. Dans la réalité, cette moyenne est pondérée par la fonction spectrale de l'instrument (*i.e.*, certaines zones spectrales influent plus que d'autres), mais celle-ci n'est pas prise en compte car elle est négligeable pour notre utilisation. La largeur spectrale est utilisée au moment d'extraire les émissivités aux fréquences de MODIS de la base d'émissivité hyperspectrale.

L'erreur d'interpolation est plus faible au niveau des bandes spectrales de MODIS, représentée en vert. Entre 1300 et 2300  $\text{cm}^{-1}$ , il n'y a aucune information issue de MODIS, mais l'erreur d'interpolation reste faible. L'utilisation de l'ACP et de la base d'apprentissage a permis à l'interpolateur de reconstituer des parties du spectre sur lesquelles il n'avait aucune information en entrée. L'erreur moyenne sur l'ensemble du spectre est d'environ  $3,5 \times 10^{-3}$  ce qui est satisfaisant pour l'interpolateur.

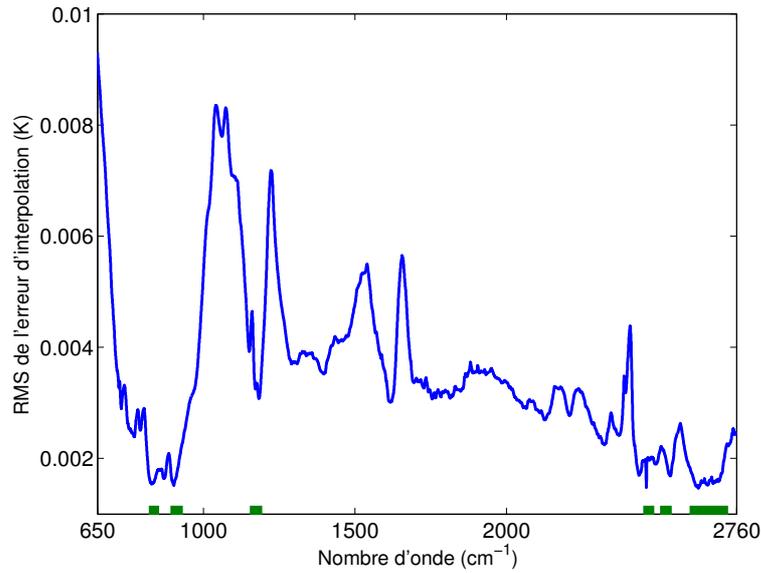


FIGURE 2.8 – RMS de l'erreur d'interpolation de l'émissivité en bleu, à partir des 6 émissivités MODIS, indiquées en vert.

Les erreurs d'interpolation sur la base d'apprentissage, de validation et de test sont très proches, ce qui vient valider plus encore notre schéma d'interpolateur. La Figure 2.9 représente la moyenne spectrale de la racine de l'erreur quadratique d'interpolation pour chaque point, pour le mois de juillet.

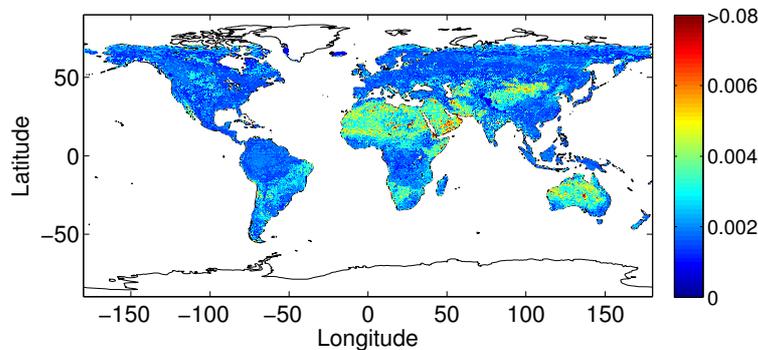


FIGURE 2.9 – Moyenne spectrale de la RMS de l'erreur d'interpolation de l'émissivité pour le mois de juillet.

Sans surprises, les régions où l'erreur d'interpolation est la plus importante sont celles où les variations spectrales sont élevées (*i.e.*, les déserts avec une forte signature des silicates). Les erreurs restent cependant faibles comparées à la valeur de l'émissivité. Le rapport reste proche de 1%.

Pour terminer cette étude sur l'interpolateur, un exemple d'interpolation d'un spectre d'émissivité au-dessus du Sahara est présenté sur la Figure 2.10. On représente sur cette

figure le spectre original issu de la base hyperspectrale de la Section 2.2.2 (page 51) correspondant à un spectre au-dessus du Sahara (en bleu). On représente également l'entrée de l'interpolateur, c'est-à-dire les émissivités correspondant aux fréquences de MODIS extraites de ce spectre (en vert). On représente enfin le spectre interpolé (en rouge). On retrouve sur cette figure les résultats précédents : à savoir la coïncidence spectrale au niveau des bandes de MODIS (ce qui valide le travail d'interpolation du réseau de neurones) et la capacité de l'interpolateur à retrouver les structures spectrales originales (grâce à la représentation spectrale de l'ACP), même dans les bandes spectrales où il n'a pas d'observations issues de MODIS.

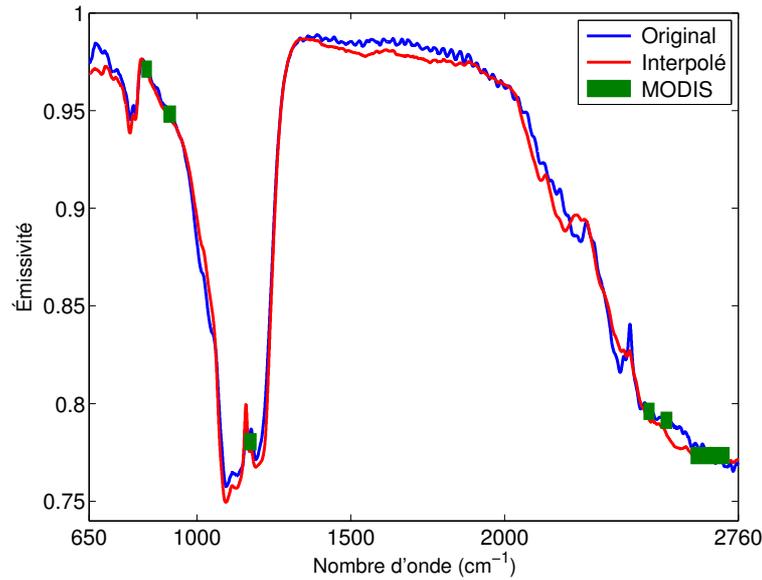


FIGURE 2.10 – Exemple d'interpolation d'un spectre au-dessus du Sahara. La courbe bleue correspond au spectre original. Les tirets verts correspondent aux émissivités MODIS qui en sont extraites pour interpoler le spectre. La courbe rouge correspond au spectre interpolé par notre schéma.

### 2.2.5 La première ébauche

Nous avons réussi à créer un outil capable de transformer six valeurs d'émissivités, restituées à partir de MODIS, en un spectre complet d'émissivité. À l'aide de cet interpolateur, nous allons construire une base mensuelle d'émissivités de surface hyperspectrales. Pour cela, nous utilisons les données MODIS réelles, présentées à la Section 2.1.2.2 (page 44). Ces bases de données mensuelles d'émissivités sont transformées en bases hyperspectrales grâce à l'interpolateur que l'on vient de construire.

L'algorithme de restitution des émissivités à partir de MODIS est basé sur un algorithme jour/nuit (voir Section 2.1.2.2, page 44). Il est donc nécessaire d'avoir une donnée de jour et une donnée de nuit, pour chaque point où l'on souhaite restituer l'émissivité. Les pôles (sud ou nord suivant la saison) sont alternativement ensoleillés ou dans la nuit pendant 6 mois. On

ne peut donc pas restituer d'émissivité à partir de MODIS au-dessus de ces points. De plus, l'algorithme nécessite une situation sans nuages. En utilisant les moyennes mensuelles de ces produits, on s'affranchit de nombreuses situations nuageuses. La projection de ces bases de données sur la grille "equal-area" à  $0,25^\circ \times 0,25^\circ$  permet également de moyenniser certaines restitutions et de couvrir une plus grande partie du globe. Cependant, certaines régions particulièrement nuageuses (notamment au niveau des tropiques et plus particulièrement l'Indonésie et l'Amazonie) peuvent rester sous un couvert nuageux tout un mois durant. Si, pendant un mois complet, une maille de la grille "equal-area" n'a pas une seule journée avec un survol par MODIS de jour puis de nuit (à la suite) sans nuages, alors il n'y aura pas de données mensuelles.

Afin de compléter notre base de données d'émissivités et de disposer d'une base globale couvrant tout le globe, il nous a fallu ajouter des données à la base interpolée à partir de MODIS. Pour cela, nous projetons les émissivités restituées par la NASA (voir Section 2.1.2.4, page 48) sur la grille "equal-area". Chaque fois qu'aucune restitution à partir de MODIS n'est disponible, nous utilisons l'émissivité restituée par la NASA. Ceci nous permet d'avoir une base mensuelle globale, au prix de discontinuités significatives à la frontière de transition. Les zones complétées grâce aux données de la NASA sont très facilement identifiables sur les cartes d'émissivités de première ébauche. Ces délimitations permettent de surveiller la sensibilité de l'algorithme que l'on construit aux erreurs de la première ébauche et ainsi de nous assurer de sa stabilité.

Cette base de données globale nous servira de première ébauche pour les restitutions de l'émissivité et de la température de la surface.

### 2.3 Un nouvel algorithme de restitution de la température de surface et de l'émissivité infrarouge

Nous disposons désormais d'une base de données globale mensuelle d'émissivités hyperspectrales. Ces émissivités sont représentées par 10 composantes de l'ACP qui codent le spectre complet. Pour certaines applications, une telle climatologie mensuelle d'émissivité est suffisante. Mais il est possible d'améliorer cette première ébauche grâce à un schéma d'inversion appliqué à des observations en temps réel. Notre objectif est de construire ce schéma de restitution d'émissivités et de température de surface en temps réel. Par rapport aux différentes méthodes décrites à la Section 1.2 (page 18), nous utilisons une méthode physique basée sur une inversion analytique (*i.e.*, mathématique) de l'équation de transfert radiatif.

Du fait du volume de données important généré par IASI, nous avons décidé d'effectuer la construction de l'algorithme de restitution sur seulement 4 semaines de données. Ces données représentent tout de même plusieurs centaines de Go (plusieurs millions de situations). Elles ont été sélectionnées pour représenter la variabilité annuelle : la première semaine des

mois de janvier, avril, juillet et octobre. Ceci est suffisant pour représenter la variabilité de l'émissivité. Nous prenons les données de l'année 2008 pour être cohérents avec les données de la NASA à notre disposition.

Nous utilisons une première information sur l'état de l'atmosphère et de la surface issue des re-analyses de l'ECMWF. Ces re-analyses sont, par souci de clarté, présentées plus en détail à la Section 5.1.1 (page 123). Il s'agit de re-analyses des données issues des modèles de prévision et des mesures satellites. Elles sont disponibles de façon globale, échantillonnées toutes les 6 heures, à une résolution spatiale de 1,125°. Elles sont composées de plus de 350 variables décrivant l'atmosphère (tableau fourni en Annexe C, page 219), depuis la température jusqu'aux propriétés nuageuses. À l'aide d'une interpolation bilinéaire en temps et en espace, chaque observation satellite est associée à une re-analyse de l'ECMWF. Le masque nuageux fourni par l'ECMWF nous sert à ne conserver que les situations en ciel clair. On considère comme claires les situations dont la couverture nuageuse est inférieure à 5%, il reste alors plus de 250.000 situations.

### 2.3.1 Inversion bayésienne du transfert radiatif

D'après la Section 1.1.2 (page 12), on peut écrire la radiance  $I_{obs}$  mesurée par l'instrument IASI, à une longueur d'onde  $\lambda$  au sommet de l'atmosphère, comme :

$$I_{obs}(\lambda) = \tau_{vrai}(\lambda) \cdot \varepsilon_{vrai}(\lambda) \cdot B_{\lambda}(T_{s_{vrai}}) + atm_{vrai} \uparrow(\lambda) \\ (+\tau_{vrai}(\lambda) \cdot (1 - \varepsilon_{vrai}(\lambda)) \cdot atm_{vrai} \downarrow(\lambda)) \quad (2.2)$$

où  $\tau_{vrai}(\lambda)$  est la transmission atmosphérique à la longueur d'onde  $\lambda$ .  $atm_{vrai} \uparrow$  est le flux atmosphérique montant.  $atm_{vrai} \downarrow$  est le flux atmosphérique descendant.  $\varepsilon_{vrai}$  est l'émissivité de la surface, et  $T_{s_{vrai}}$  sa température. Toutes les variables sont appelées *vrai* pour les différencier facilement de leurs estimations à venir.

Les re-analyses de l'ECMWF nous fournissent l'information de l'état de l'atmosphère et une première ébauche pour la température de surface ( $T_{s_{fg}}$  pour "first-guess", traduction anglaise de première ébauche). L'estimation de l'état de l'atmosphère dans ces analyses n'est pas parfaite, mais nous allons nous focaliser sur des canaux "fenêtres" qui seront, par définition, moins sensibles à ces erreurs. La première ébauche en émissivité ( $\varepsilon_{fg}$ ) nous est fournie par la climatologie mensuelle d'émissivité construite à la Section 2.2.5 (page 59) à l'aide de l'interpolateur. Cette première ébauche nous permet de calculer, à l'aide de RTTOV, une radiance  $I_{calc}$  qui serait mesurée par IASI à la longueur d'onde  $\lambda$  dans ces conditions. On retrouve alors la même formulation qu'à l'équation (2.2) mais avec *calc* et *fg* à la place de *vrai* :

$$I_{calc}(\lambda) = \tau_{calc}(\lambda) \cdot \varepsilon_{fg}(\lambda) \cdot B_{\lambda}(T_{s_{fg}}) + atm_{calc} \uparrow(\lambda) \\ (+\tau_{calc}(\lambda) \cdot (1 - \varepsilon_{fg}(\lambda)) \cdot atm_{calc} \downarrow(\lambda)) \quad (2.3)$$

avec les mêmes notations que précédemment, où  $fg$  indique les données de première ébauche et  $calc$  celles issues des calculs de transfert radiatif à partir de la première ébauche et des informations atmosphériques des re-analyses.

Nous cherchons à estimer la différence  $I_{obs} - I_{calc}$ , afin de déterminer les variables *vraies*.

Nous considérons ces équations à  $N$  longueurs d'onde,  $\lambda_1, \lambda_2, \dots, \lambda_N$ , où l'atmosphère est transparente. Ces longueurs d'onde seront sélectionnées par la suite, elles seront également appelées "canaux" car elles correspondent aux longueurs d'onde des canaux IASI sélectionnés. Du fait de la transparence de l'atmosphère, les contributions atmosphériques, c'est-à-dire les termes  $atm \uparrow$  et  $atm \downarrow$ , peuvent être négligées par rapport aux contributions de surface. De plus, les termes atmosphériques, déjà négligeables par rapport aux termes de surface, sont très bien estimés par le calcul de transfert radiatif car la description atmosphérique des re-analyses que nous utilisons est raisonnable. La différence entre les termes atmosphériques réels et calculés est donc négligeable. Nous considérons également, pour les mêmes raisons, que la transmission atmosphérique est très bien estimée par le calcul de transfert radiatif ( $\tau_{vrai} = \tau_{calc} = \tau$ ). On peut alors écrire la différence entre la radiance vraie et la radiance calculée d'après les équations (2.2) et (2.3) comme :

$$(I_{obs} - I_{calc})(\lambda) = \tau(\lambda) \cdot (\varepsilon_{vrai}(\lambda) \cdot B_\lambda(T_{s_{vrai}}) - \varepsilon_{fg}(\lambda) \cdot B_\lambda(T_{s_{fg}}))$$

C'est-à-dire que l'écart entre les radiances vraies et celles calculées à l'aide RTTOV n'est dû qu'à une différence entre les termes liés à l'émission de la surface. On peut réécrire cette équation comme :

$$\varepsilon_{vrai}(\lambda) \cdot B_\lambda(T_{s_{vrai}}) = \underbrace{\frac{(I_{obs} - I_{calc})(\lambda)}{\tau(\lambda)} + \varepsilon_{fg}(\lambda) \cdot B_\lambda(T_{s_{fg}})}_{=A_\lambda} \quad (2.4)$$

pour les  $N$  longueurs d'onde considérées.  $A_\lambda$  est différent à chaque longueur d'onde, on note par la suite  $A$  le vecteur des  $A_\lambda$  pour toutes les longueurs d'onde considérées.

En utilisant l'ACP décrite à la Section 2.2.3 (page 53), on peut écrire les  $N$  émissivités comme :

$$\varepsilon = (c_1, c_2, \dots, c_P) \cdot VP_{P,N} + \overline{\varepsilon_N}$$

où  $P$  est le nombre de composantes considérées dans la restitution.  $\overline{\varepsilon_N}$  représente les  $N$  émissivités moyennes aux  $N$  longueurs d'onde considérées.  $VP_{P,N}$  est une partie de la matrice des vecteurs propres définie par :

$$VP_{P,N} = \begin{pmatrix} vp_{1,\lambda_1} & \cdots & vp_{1,\lambda_N} \\ \vdots & \ddots & \vdots \\ vp_{P,\lambda_1} & \cdots & vp_{P,\lambda_N} \end{pmatrix}$$

En remplaçant l'émissivité dans les  $N$  équations (2.4) par son expression en fonction des

composantes, on obtient une seule équation matricielle que l'on peut écrire comme :

$$((c_1, c_2, \dots, c_P) \cdot VPP_{P,N} + \bar{\varepsilon}) \cdot B(T_{s_{vrai}}) = A$$

où  $B(T_{s_{vrai}})$  correspond au vecteur des  $N$  fonctions de Planck aux  $N$  longueurs d'onde considérées. On peut alors réécrire cette équation comme :

$$(c_1, c_2, \dots, c_P) \cdot VPP_{P,N} + \underbrace{\frac{-A}{B(T_{s_{vrai}})}}_{=F(T_{s_{vrai}})} = -\bar{\varepsilon} \quad (2.5)$$

Dans cette équation, la fonction  $F$  est définie pour chaque longueur d'onde  $\lambda_i$  considérée. Pour une longueur d'onde donnée  $\lambda_i$ ,  $F_i$  peut s'écrire :

$$F_i(T_{s_{vrai}}) = -A_i \cdot \frac{\lambda_i^5}{2 \cdot h \cdot c^2} \cdot \left( e^{\frac{h \cdot c}{\lambda_i \cdot k \cdot T_{s_{vrai}}}} - 1 \right) \quad (2.6)$$

Afin de linéariser le problème, on utilise un développement de Taylor de  $F$  au premier degré, autour de la première ébauche. Cette approximation est justifiée si la première ébauche est assez proche des variables réelles pour permettre la linéarisation. Dans ce cas, le développement de Taylor est une bonne estimation de la fonction  $F$  localement, car la première ébauche ainsi que les re-analyses sont de bonne qualité. On peut alors écrire chaque fonction  $F_i$  comme :

$$F_i(T_{s_{vrai}}) \approx F_i(T_{s_{calc}}) + F'_i(T_{s_{calc}}) \cdot \underbrace{(T_{s_{vrai}} - T_{s_{calc}})}_{\Delta T_s}$$

où  $F'_i(T_{s_{calc}})$  est la dérivée de  $F_i$  en  $T_{s_{calc}}$ , qui peut se calculer facilement d'après l'expression de  $F_i$  donnée à l'équation (2.6) :

$$F'_i(T_{s_{calc}}) = A_i \cdot \frac{\lambda_i^5}{2 \cdot h \cdot c^2} \cdot \frac{h \cdot c}{\lambda_i \cdot k \cdot T_{s_{calc}}^2} \cdot e^{\frac{h \cdot c}{\lambda_i \cdot k \cdot T_{s_{calc}}}}$$

On garde les notations  $F'_i$  pour rendre les équations suivantes plus facilement intelligibles et éviter de perdre le lecteur sous un flot de notations. On reprend donc l'équation (2.5) dans laquelle on introduit le développement de Taylor :

$$(c_1, c_2, \dots, c_P) \cdot VPP_{P,N} + F(T_{s_{calc}}) + F'(T_{s_{calc}}) \cdot \Delta T_s \approx -\bar{\varepsilon}$$

où  $F$  et  $F'$  correspondent respectivement aux vecteurs des  $F_i$  et  $F'_i$  pour les  $N$  longueurs

d'onde considérées. On a alors :

$$\left( \begin{array}{cccc} c_1 & \cdots & c_P & \Delta T_s \end{array} \right) \cdot \underbrace{\left( \begin{array}{ccc} vp_{1,\lambda_1} & \cdots & vp_{1,\lambda_N} \\ \vdots & \ddots & \vdots \\ vp_{P,\lambda_1} & \cdots & vp_{P,\lambda_N} \\ F'_1(T_{s_{calc}}) & \cdots & F'_N(T_{s_{calc}}) \end{array} \right)}_{=M} = -\bar{\varepsilon} - F(T_{s_{calc}})$$

On peut donc écrire cette équation matriciellement de façon très simple :

$$\left( \begin{array}{cccc} c_1 & \cdots & c_P & \Delta T_s \end{array} \right) \cdot M = -\bar{\varepsilon} - F(T_{s_{calc}}) \quad (2.7)$$

où les  $p + 1$  inconnues à déterminer sont  $c_1, c_2, \dots, c_P$  et  $\Delta T_s$ . Il ne reste donc plus qu'à inverser mathématiquement cette équation pour obtenir la température et l'émissivité de surface.

Dans une première version de l'algorithme, nous nous contentions d'utiliser une pseudo-inverse ( $M^1 = (M^T \cdot M)^{-1} \cdot M^T$ ) (Aires 1999). Cette méthode permet d'approximer l'inverse d'une matrice directement. On inverse alors la matrice  $M$  de l'équation (2.7) pour obtenir :

$$\left( \begin{array}{cccc} c_1 & \cdots & c_P & \Delta T_s \end{array} \right) \approx (-\bar{\varepsilon} - F(T_{s_{calc}})) \cdot (M^T \cdot M)^{-1} \cdot M^T$$

L'inconvénient de cette méthode est son manque de stabilité. Aucune contrainte ne s'exerce sur la solution et si certains paramètres sont mal ajustés, le résultat peut devenir aberrant. En particulier, une mauvaise estimation de la première ébauche peut mener à des erreurs supérieures à celles de la première ébauche en elle-même, l'algorithme peut diverger. De plus, l'algorithme que l'on a mis en place n'est valable que dans une zone relativement proche de la première ébauche, du fait de l'utilisation du développement de Taylor. Nous avons donc décidé de modifier l'algorithme pour mieux contraindre la solution et ainsi gagner en stabilité.

Pour ce faire, nous utilisons un estimateur bayésien (Aires 1999). Il s'agit d'un estimateur de maximum *a posteriori* qui utilise la relation de Bayes. Afin d'épargner au lecteur une nouvelle série d'équations, nous utilisons ici directement la solution obtenue. L'Annexe B.2.1 (page 213) fournit plus de précisions sur cet estimateur<sup>5</sup>.

Nous avons la solution :

$$\left( \begin{array}{cccc} c_1 & \cdots & c_P & \Delta T \end{array} \right) \approx f_g + \left( M^T \cdot S_F^{-1} \cdot M + S_{f_g}^{-1} \right)^{-1} \cdot M^T \cdot S_F^{-1} \cdot (-\bar{\varepsilon} - F(T_{s_{fg}}) - M \cdot f_g) \quad (2.8)$$

où  $f_g$  correspond à la première ébauche, c'est-à-dire à l'émissivité de première ébauche et à l'écart en température de surface de première ébauche  $\Delta T_s = 0$ .  $S_F$  correspond à la matrice

5. Le fait d'avoir  $f \cdot M = F$  et non  $M \cdot f = F$  comme dans l'annexe n'est pas un problème, il suffit de considérer la transposée de l'équation pour revenir au même résultat :  $M^T \cdot f^T = F^T$

de covariance d'erreur de  $F$  (car  $\bar{\varepsilon}$  n'est pas une source d'erreur, il s'agit d'une simple valeur moyenne).  $S_{fg}$  correspond à la matrice de covariance d'erreur de la première ébauche.

Pour pouvoir utiliser cet estimateur bayésien on a donc considéré que  $F$  suivait une loi gaussienne de matrice de covariance  $S_F$ , et que l'erreur de la première ébauche était également gaussienne avec une matrice de covariance d'erreur  $S_{fg}$ . Il est courant, voire systématique, de considérer que ce type d'erreurs suit des lois gaussiennes lorsque des méthodes analytiques sont utilisées. Cette approximation en elle-même ne génère pas beaucoup d'erreur. Il faut alors estimer correctement les matrices de covariance d'erreur, notamment celle de la première ébauche, qui est la plus compliquée à estimer.

Cette inversion analytique de l'équation de transfert radiatif nous permet de restituer directement l'émissivité (sous forme de composantes) et la température de surface, en connaissant les radiances mesurées par le satellite, une bonne information de l'état de l'atmosphère et une première ébauche de l'état des caractéristiques de la surface.

L'avantage de l'utilisation de l'estimateur bayésien est que l'on a accès directement à la matrice de covariance d'erreur de l'inversion. D'après l'Annexe B.2.1 (page 213), la matrice de covariance d'erreur  $Q$  de la solution est donnée par :

$$Q = \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \quad (2.9)$$

Cette estimation de l'erreur d'inversion donne directement accès à la qualité théorique de la restitution. Elle sera utilisée par la suite.

### 2.3.2 Sélection des paramètres de la restitution

Il faut désormais déterminer les différents paramètres qui entrent en jeu dans l'algorithme d'inversion : les deux matrices de covariance  $S_F$  et  $S_{fg}$  qui sont utilisées dans l'estimateur bayésien, mais également le nombre de composantes  $P$  prises en compte et les  $N$  longueurs d'onde sélectionnées pour la restitution.

#### 2.3.2.1 La matrice de covariance d'erreur de $F$

Il faut construire les matrices de covariance d'erreur de la fonction  $F$ ,  $S_F$ , et de la première ébauche,  $S_{fg}$ . D'après l'expression de la fonction  $F$  de l'équation (2.6), elle est dépendante de la matrice  $A$ , elle-même dépendante, d'après l'équation (2.4), de la radiance réelle  $I_{obs}$  et de la radiance calculée  $I_{calc}$ . Nous avons donc décidé de prendre en compte les erreurs sur les radiances observées et calculées pour déterminer la matrice de covariance de  $F$ . L'erreur instrumentale de IASI est donnée par le constructeur. Elle est présentée plus en détails à la Section 5.1.2 (page 124). Cette erreur instrumentale nous permet de connaître l'erreur sur la radiance observée. La matrice de covariance d'erreur est ici diagonale, l'erreur sur chaque canal étant considérée comme indépendante. Pour déterminer l'erreur sur la radiance calculée, nous utilisons les résultats fournis par [Matricardi \(2009\)](#), qui décrivent

les erreurs du code de transfert radiatif RTTOV que nous utilisons. Ici aussi, les erreurs des différents canaux sont indépendantes entre elles. On somme les deux matrices de covariance d’erreur pour obtenir une matrice diagonale de covariance d’erreur sur la différence entre les radiances (ce qui est couramment fait en assimilation variationnelle dans les NWP). Il faut ensuite appliquer la formule de  $F$  donnée à l’équation (2.6) pour obtenir la matrice de covariance  $S_F$ .

Nous avons considéré que les erreurs instrumentales étaient indépendantes (*i.e.*, matrice de covariance d’erreur diagonale). Cette approximation, qui est communément effectuée, surestime néanmoins le bruit en ne prenant pas en compte les termes non-diagonaux qui déterminent des corrélations entre les erreurs sur chaque canal. De plus, nous n’avons pris en compte que l’erreur due au code de transfert radiatif pour la détermination de la matrice de covariance d’erreur de  $I_{calc}$ . Il faudrait, pour être plus précis, prendre en compte les erreurs d’estimation des différents paramètres atmosphériques et leur impact sur la radiance calculée. Le fait de ne prendre en compte que des canaux “fenêtres” permet de minimiser l’impact de cette approximation. Cependant, une mauvaise estimation de l’état de l’atmosphère entraînerait des corrélations d’erreur entre les différents canaux et minimiserait l’erreur totale. Nous avons préféré surestimer l’erreur sur les radiances plutôt que de chercher à caractériser plus précisément les différentes sources d’erreur. En essayant de trop déterminer les différentes erreurs, nous risquons de les sous-estimer et de fausser la restitution.

### 2.3.2.2 La matrice de covariance d’erreur de la première ébauche

Pour déterminer la matrice  $S_{fg}$ , nous avons calculé la matrice de covariance des émissivités de la première ébauche pour chaque point de la grille “equal-area”. Nous disposons d’une base mensuelle d’émissivité sur une grille “equal-area” à  $0,25^\circ \times 0,25^\circ$  (voir Section 2.2.5, page 59). Nous avons alors 12 valeurs d’émissivité (pour les douze mois de l’année) pour chaque point de la carte.

Afin d’augmenter la variabilité des émissivités, pour un pixel donné, on calcule la variance des valeurs d’émissivité du pixel et de ses 4 plus proches voisins (spatialement) sur les douze mois de l’année. On dispose ainsi de 60 valeurs d’émissivité pour chaque pixel, ce qui permet une meilleure estimation de la variance. Cette façon de procéder permet de construire une matrice de covariance d’erreur des émissivités de la première ébauche non-diagonale et semblable à la dispersion saisonnière. La projection de cette matrice de covariance d’erreur sur les composantes de l’émissivité est simple, du fait de la linéarité de l’ACP.

Nous avons donc déterminé le bloc supérieur droit de taille  $P \times P$  de la matrice  $S_{fg}$  (de taille  $P + 1 \times P + 1$ ). La dernière ligne (et la dernière colonne) correspond aux covariances d’erreur sur la température de surface. Nous avons décidé ici de considérer que l’écart-type de l’erreur d’estimation de la température de surface est d’environ 5 K, sans corrélation avec les erreurs sur les émissivités (Tsuang et al. 2008). La matrice de covariance d’erreur de la première ébauche peut donc s’écrire en fonction de la matrice de covariance d’erreur

$S_{compo}$  (matrice  $P \times P$ ) sur les émissivités et de l'écart-type de l'erreur sur la température de surface  $\sigma_{T_s}$  comme :

$$S_{fg} = \begin{pmatrix} & & & 0 \\ & S_{compo} & & \vdots \\ & & & 0 \\ 0 & \cdots & 0 & \sigma_{T_s}^2 \end{pmatrix}$$

Nous avons ainsi déterminé les matrices de covariance d'erreur utiles à l'estimateur bayésien. Elles sont dépendantes de la situation considérée. La matrice de covariance d'erreur sur la première ébauche dépend du point considéré. Ceci permet de mieux contraindre l'émissivité sur les zones à faible variabilité et inversement.

### 2.3.2.3 Les canaux sélectionnés

Nous nous concentrons maintenant sur le paramètre  $N$ . La sélection des longueurs d'onde  $\lambda_1, \lambda_2, \dots, \lambda_N$  prises en compte dans la restitution est cruciale. L'algorithme de restitution part du principe que la contribution atmosphérique calculée à partir de la première ébauche est égale à la contribution réelle. Si on ne sélectionne que des canaux "fenêtres", la contribution atmosphérique sera faible par rapport à celle de la surface et l'affirmation précédente n'en sera que plus proche de la réalité. IASI comporte trois régions considérées comme "fenêtres" (voir Section 1.3.1, page 29) : de 770 à 980  $\text{cm}^{-1}$ , de 1080 à 1150  $\text{cm}^{-1}$  et de 2420 à 2700  $\text{cm}^{-1}$ . Nous ne considérons que les deux premières, la troisième étant trop contaminée par le bruit instrumental de IASI. Le rapport signal sur bruit de la troisième bande est très faible et elle souffre également de contamination solaire.

Parmi ces deux régions spectrales, on calcule la transmission atmosphérique en fonction de la longueur d'onde. La moyenne de celle-ci sur les 4 semaines de données considérées (premières semaines de janvier, avril, juillet et octobre 2008) est représentée en bleu sur la Figure 2.11. L'échelle correspondante est celle de droite.

On peut constater que parmi ces zones spectrales supposées "fenêtres", il y a de nombreux canaux pour lesquels la transmission atmosphérique est faible. Ceci est dû à des bandes d'absorption très fines de divers composants atmosphériques. Si nous nous contentons de sélectionner les canaux IASI pour lesquels la transmission atmosphérique est importante, on risque d'obtenir des canaux très proches de certaines bandes d'absorption et risquer d'être influencé par la composition atmosphérique.

Nous avons calculé le gradient de la transmission atmosphérique en fonction du nombre d'onde. La moyenne de ce gradient sur les 4 semaines de données est représentée en vert sur la Figure 2.11. L'échelle correspondante est celle de gauche. Sélectionner les longueurs d'onde en fonction du gradient de la transmission atmosphérique, plutôt que de sa valeur en elle-même, nous permet d'éviter les raies d'absorption des composants mineurs de l'atmosphère.

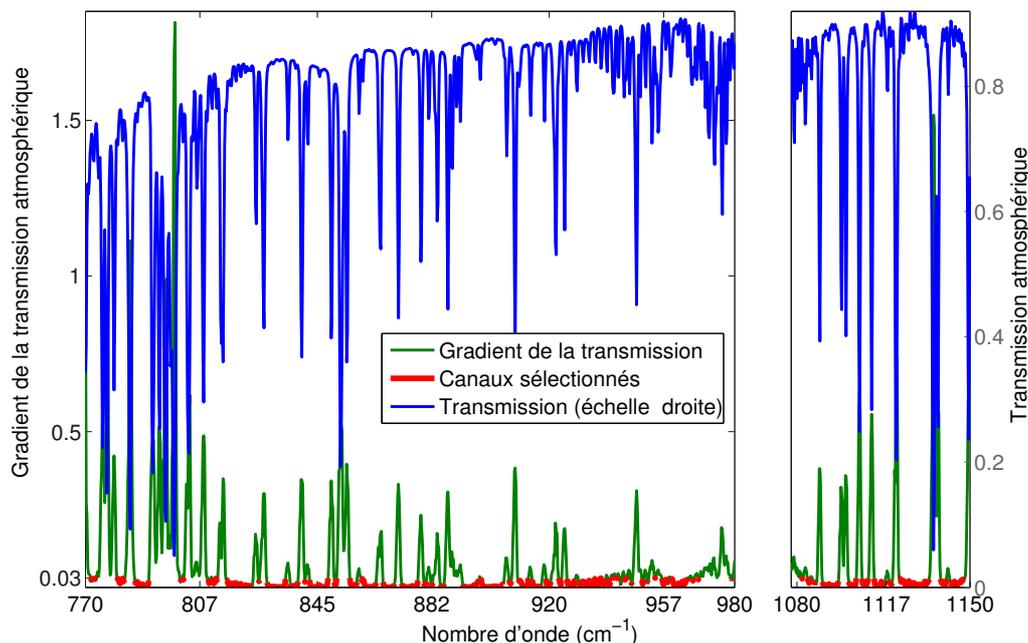


FIGURE 2.11 – Sélection des canaux “fenêtres” pour la restitution : la courbe bleue correspond à la transmission atmosphérique moyenne (échelle de droite), la courbe verte correspond au gradient moyen en nombre d’onde de la transmission atmosphérique (échelle de gauche) et les points rouges correspondent aux canaux sélectionnés.

À la suite de divers tests sur le nombre de canaux sélectionnés et leurs impacts sur l’algorithme de restitution, nous avons décidé de sélectionner tous les canaux dont le gradient moyen est inférieur à 0,03. Ces canaux sont indiqués par des points rouges sur la Figure 2.11. Cette sélection isole 512 longueurs d’onde différentes correspondant à 512 canaux IASI qui seront utilisés dans la restitution.

#### 2.3.2.4 Le nombre de composantes de l’ACP

Afin de choisir le nombre  $P$  de composantes de l’ACP à prendre en compte dans l’algorithme, de nombreux tests ont été effectués. Augmenter le nombre de composantes sélectionnées augmente les degrés de liberté de l’inversion et diminue ainsi sa stabilité. Inversement, si le système est trop peu contraint, la représentation des spectres d’émissivité est insuffisante, certaines structures spectrales ne sont pas atteignables, les spectres seraient alors trop lisses. L’utilisation de l’estimateur bayésien pour mieux contraindre le problème ne nous affranchit pas des précautions nécessaires aux algorithmes d’inversion.

Des tests *a posteriori* ont été mis en place afin de choisir le nombre de composantes optimales. On considère que les restitutions où le changement de la température de surface est supérieur à 20K sont erronées. Le nombre de telles situations est un bon indicateur de la stabilité de la méthode. La Figure 2.12 illustre cette méthode de sélection. Dans

la partie haute de la figure, nous avons représenté la racine de la différence quadratique moyenne entre la température de brillance observée et celle calculée à partir des variables de surface restituées, en fonction du nombre de composantes prises en compte. L'erreur diminue jusqu'à  $P \approx 15$ , puis augmente. La première partie décroissante est liée à la meilleure représentation du spectre d'émissivité par un nombre plus important de composantes. Au delà de 15 composantes, l'algorithme n'arrive plus à prendre en compte tant d'informations et l'erreur augmente alors à nouveau.

La partie basse de la Figure 2.12 représente le pourcentage de situations stables en fonction du nombre de composantes sélectionnées. On remarque qu'à partir de 10 composantes la quantité de situations stables est fortement décroissante. Malgré une légère croissance aux alentours de 20 composantes, le nombre de situations stables diminue jusqu'à ce qu'il n'y ait plus de situations stables pour 40 composantes. La forte décroissance de la stabilité de l'algorithme à partir de 10 composantes est liée à un trop plein de variables à restituer.

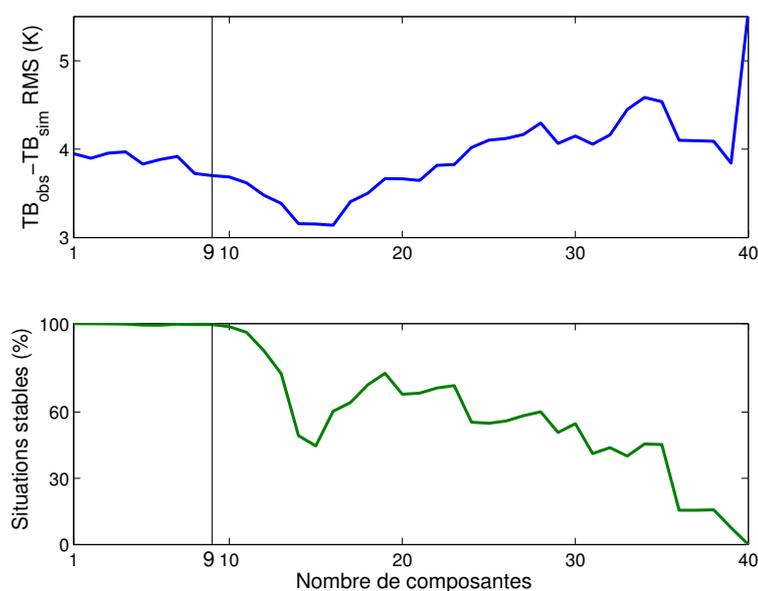


FIGURE 2.12 – Sélection du nombre de composantes pour la restitution : la courbe bleue dans la partie haute correspond à la racine de l'erreur quadratique moyenne entre les températures de brillance observées et celles calculées grâce aux variables de surface restituées. La courbe verte dans la partie basse correspond au pourcentage total des situations où la restitution est stable. Ces deux courbes sont tracées en fonction du nombre de composantes utilisées.

Afin de s'assurer de la parfaite stabilité de l'algorithme, nous avons décidé de ne prendre en compte que 9 composantes d'émissivité dans la restitution pour limiter au maximum le nombre de situations instables. Les quelques situations instables restantes sont liées à des erreurs d'estimation de la couverture nuageuse. En effet, les re-analyses de l'ECMWF étant disponibles que toutes les 6 heures, une mauvaise caractérisation des nuages est probable. Cette mauvaise caractérisation entraîne des différences importantes entre les radiances cal-

culées et observées et implique alors une importante modification de la première ébauche, ce qui explique cette différence<sup>6</sup>.

### 2.3.2.5 Conclusion

Au final, l'algorithme résout un système de 512 équations (une pour chaque canal "fenêtre" sélectionné) et 10 inconnues (les 9 composantes de l'émissivité et la température de surface). Le problème est donc surdéterminé, ce qui facilite sa stabilité. La redondance de l'information dans les observations diminue la sensibilité de l'algorithme au bruit.

### 2.3.3 Influence de la première ébauche

Afin de tester la sensibilité de l'algorithme à la qualité de la première ébauche, nous avons testé plusieurs inversions utilisant plusieurs premières ébauches différentes (dans la limite du bruit considéré sur la première ébauche). Les résultats similaires obtenus nous permettent de confirmer la stabilité de l'inversion. La gamme de bruit considérée sur la première ébauche, tant sur les composantes de l'émissivité que sur la température de surface, permet une assez large modification de la première ébauche. Cependant, une telle modification ne change pas le résultat de l'inversion.

Nous avons essayé de lancer l'algorithme en utilisant une émissivité constante de 0,98 comme première ébauche (approximation fréquemment effectuée par le passé, dans les NWP). Cette fois ci, l'algorithme a perdu en stabilité et les résultats obtenus étaient différents de ceux espérés. Ce comportement peut être dû à deux principales causes. Premièrement, il est difficile de représenter un spectre parfaitement lisse avec les composantes de l'ACP que l'on a utilisées. Les composantes ayant été calculées sur une base de spectres mesurés en laboratoire, il n'y avait aucun spectre parfaitement lisse. L'ACP a alors beaucoup de mal à représenter un tel spectre et la projection sur l'espace des composantes entraîne une grande instabilité, avec des valeurs de composantes importantes qui essaient de se compenser les unes les autres. Deuxièmement, l'algorithme en lui-même nécessite une première ébauche robuste et assez proche de la réalité pour que la linéarisation locale du transfert radiatif soit réaliste. Une première ébauche trop différente de la réalité entraîne donc des instabilités.

Un autre test de stabilité de l'algorithme a été mené. Il s'agit d'itérer l'inversion plusieurs fois et d'examiner l'évolution des variables restituées. Concrètement, nous relançons l'algorithme en utilisant les données précédemment restituées comme première ébauche. Le résultat que nous obtenons est quasi identique à celui obtenu à la première itération. Ceci démontre que le schéma d'inversion a exploité au mieux, dès la première itération, l'in-

---

6. Prendre en compte 9 composantes nous permet de nous affranchir de presque toutes les situations instables. Dans le cas des situations stables, l'algorithme pourrait utiliser 15 composantes afin de mieux représenter les spectres d'émissivité. Cela impliquerait une différenciation de l'algorithme suivant la situation. Nous n'avons pas fait ce choix ici.

formation qui était disponible (les observations IASI, la première ébauche et les données auxiliaires sur l'état de l'atmosphère).

## 2.3.4 Résultats et évaluation

### 2.3.4.1 Conditions expérimentales

Comme les mesures *in situ* de l'émissivité et de la température de surface sont limitées sur le globe, tant spatialement que temporellement, il est difficile d'effectuer des comparaisons directes avec les produits que l'on restitue. De plus, l'étendue du champ de vue de l'instrument IASI au sol est d'au moins 12 km. Les mesures de l'émissivité depuis le sol sont plus localisées. D'importantes différences peuvent alors apparaître entre les différentes mesures, dues à des inhomogénéités de l'émissivité dans le champ de vue de IASI. Un moyen d'éviter ces sources d'erreur dans la validation des émissivités restituées est de limiter l'étude à des zones du globe particulièrement homogènes (par exemple les désert de Namib ou de Kalahari) (Zhou et al. 2011; Hulley and Hook 2009). De telles études valident en partie les résultats obtenus, mais de façon très localisée. Nous cherchons ici à valider notre algorithme d'inversion de façon globale.

Une autre méthode de validation est mise en place par Li et al. (2010). Il s'agit de comparer les radiances observées aux radiances calculées par des simulations du transfert radiatif en utilisant les émissivités restituées. Au lieu de comparer directement les radiances canal par canal, ce sont des variations entre canaux qui sont comparées. Ceci permet, théoriquement, de supprimer les autres sources d'erreur, comme la mauvaise caractérisation de l'atmosphère (uniquement dans la comparaison, pas dans la restitution). Nous avons décidé de ne pas utiliser cette méthode car nous souhaitons conserver les erreurs liées à une mauvaise caractérisation de l'atmosphère. En effet, dans les zones spectrales où l'atmosphère est opaque, les émissivités et la température de surface n'influencent pas les radiances calculées. La différence entre les radiances observées et calculées provient alors uniquement des erreurs dans le code de transfert radiatif et dans la mauvaise caractérisation de l'atmosphère. La connaissance de ces erreurs permet alors de vérifier que l'algorithme ne les a pas compensées en modifiant artificiellement la température de surface ou l'émissivité. Les profils que nous utilisons ici sont tous issus des re-analyses de l'ECMWF : on s'attend, dès lors, à retrouver des erreurs systématiques sur la caractérisation atmosphérique.

Nous avons donc décidé de comparer directement les spectres simulés grâce au code de transfert radiatif utilisant les émissivités et la température de surface restituées à partir des mesures IASI, aux mesures IASI elles-mêmes. Ici encore, plutôt que de comparer les radiances qui ont des valeurs très différentes d'un canal à l'autre, nous comparons les températures de brillance, qui ont des différences du même ordre d'un canal à l'autre. Pour plus de simplicité, nous appelons  $TB_{obs}$  le spectre de température de brillance de IASI réel et  $TB_{sim}$  le spectre simulé grâce à RTTOV en utilisant les variables de surface restituées. À

l'image de ce qui a été fait par [Vogel et al. \(2011\)](#), nous préférons utiliser cette comparaison directe pour sa simplicité et la facilité de son interprétation. Les statistiques de la différence ( $TB_{obs} - TB_{sim}$ ) donnent directement accès à l'impact spectral et spatial de l'utilisation des émissivités et de la température de surface restituées.

Nous utilisons les quatre semaines de données réelles IASI décrites à la Section 2.3 (page 60, il s'agit des premières semaines de janvier, avril, juillet et octobre de l'année 2008). Les informations de surface, du profil atmosphérique et de la couverture nuageuse des re-analyses de l'ECMWF sont colocalisées avec chaque sondage IASI. Afin de mesurer l'apport des émissivités restituées, la comparaison est effectuée en utilisant différentes sources d'émissivités :

- Les différentes bases d'émissivités décrites à la Section 2.1.2 (page 42) : UWIREMIS qui est également liée à MODIS (voir Section 2.1.2.3, page 46) et la base de la NASA qui consiste en des moyennes mensuelles de restitutions IASI (voir Section 2.1.2.4, page 48) ;
- Les émissivités de la première ébauche issue de l'interpolateur (voir Section 2.2.5, page 59), que l'on dénomme émissivités interpolées par la suite ;
- Les émissivités restituées à partir des observations IASI.

Les émissivités interpolées sont basées sur les données MODIS de 2007, de même que la base UWIREMIS. Les émissivités de la NASA sont des moyennes mensuelles d'émissivités restituées à partir de IASI sur l'année 2008 (ce qui "favorise" ces émissivités, par rapport aux autres estimations). Les émissivités restituées correspondent au spectre de température de brillance de IASI considéré. Toutes les statistiques sont effectuées sur la base de données de 4 semaines de mesures IASI en 2008. Nous ne présentons pas les statistiques de la base ARA car cette dernière est limitée à la ceinture tropicale (*i.e.*, entre  $-30$  et  $+30^\circ$ )<sup>7</sup>.

#### 2.3.4.2 Analyse spectrale

La Figure 2.13 présente les statistiques des différences entre les températures de brillance calculées et les réelles. Nous utilisons la racine carrée de l'erreur quadratique moyenne (appelée RMS de l'erreur). Cette erreur correspond à la racine de la somme du biais au carré et de la variance de l'erreur au carré, elle inclut donc les erreurs de biais et d'écart-type. Utiliser la RMS de l'erreur permet de caractériser l'écart au mieux avec une seule mesure. Ces statistiques sont effectuées sur les 4 semaines de mesures IASI considérées, en ciel clair.

Les courbes colorées correspondent à la différence entre les températures de brillance réelles observées par IASI et les températures de brillance calculées par RTTOV, en utilisant l'atmosphère et la température de surface issus des profils ECMWF et différentes

---

7. De ce fait, le nombre de sondage IASI considérés est différent, les atmosphères sont différents, les résultats obtenus sont donc légèrement différents. Cependant, en effectuant des statistiques uniquement sur la ceinture tropicale, on constate que les résultats obtenus avec cette base de données sont similaires à ceux obtenus avec UWIREMIS, les données de la NASA ou les données interpolées.

informations pour l'émissivité de surface. Seule l'émissivité de surface est modifiée d'une courbe à l'autre. La courbe bleue correspond au calcul utilisant les émissivités interpolées, la courbe verte correspond au calcul utilisant les émissivités UWIREMIS et la courbe rouge au calcul utilisant les émissivités de la NASA. Pour obtenir les statistiques de la courbe noire, nous avons modifié à la fois l'émissivité utilisée dans le calcul de transfert radiatif mais également la température de surface. Cette courbe correspond donc aux statistiques d'écart entre les températures de brillance réelles et celles calculées à partir des profils atmosphériques ECMWF et de la température de surface et de l'émissivité restituée.

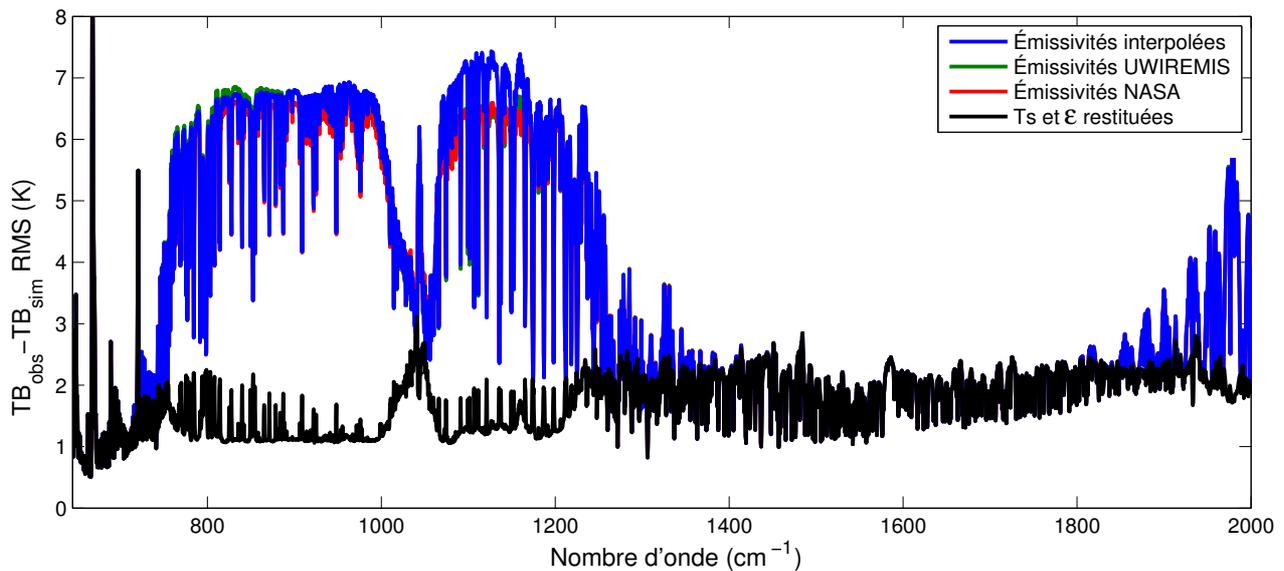


FIGURE 2.13 – Chaque courbe correspond à la racine de l'erreur quadratique moyenne entre le spectre de température de brillance observé par IASI et un spectre de température de brillance calculé à partir de la température de surface donnée par l'ECMWF et les émissivités interpolées (courbe bleue), les émissivités UWIREMIS (courbe verte), les émissivités de la NASA (courbe rouge). La courbe noire correspond au même calcul en utilisant la température de surface et les émissivités restituées.

Nous ne présentons ici que les bandes 1 et 2 de IASI. En effet, la bande 3 de 2000 à 2760  $\text{cm}^{-1}$  n'est pas montrée car elle présente un bruit instrumental trop important, ainsi que de la contamination solaire. Les quelques pics que l'on peut remarquer entre 660 et 700  $\text{cm}^{-1}$  sont dus à des composants atmosphériques mal considérés dans les re-analyses de l'ECMWF. Ces composants présentent de fines bandes d'absorption qui ne sont pas calculées par RTTOV puisque ils n'apparaissent pas dans les profils atmosphériques. Ces pics ne sont pas liés à des propriétés de surface, il ne seront donc pas discutés plus en avant ici<sup>8</sup>.

On peut également voir très clairement apparaître sur cette figure la bande d'absorption

8. Un calcul de transfert radiatif utilisant des profils atmosphériques restitués à partir des observations IASI devrait réduire ces pics.

de l’ozone entre 1000 et 1100  $\text{cm}^{-1}$ . Quelques soient les propriétés de surface considérées, toutes les courbes sont très proches dans cette zone spectrale. Ceci est dû à des erreurs dans la caractérisation du profil d’ozone, dans les re-analyses de l’ECMWF. Comme avec les composants mineurs de l’atmosphère qui entraînaient des pics aux alentours de 700  $\text{cm}^{-1}$ , la mauvaise caractérisation de l’ozone entraîne une erreur systématique entre les températures de brillance observées et calculées. L’erreur entre les températures de brillance est plus faible dans cette zone car l’atmosphère y est plus opaque, les erreurs de caractérisation de la surface n’entrent donc plus en jeu. L’erreur est uniquement liée à une mauvaise estimation de l’atmosphère. Les valeurs que l’on y retrouve sont proches des valeurs observées entre 1400 et 1800  $\text{cm}^{-1}$ . Dans cette région spectrale, l’atmosphère est opaque, c’est donc uniquement une mauvaise caractérisation de la contribution atmosphérique au rayonnement mesuré par IASI qui entraîne cette erreur avoisinant les 2 K.

On peut constater que les trois courbes colorées (bleue, verte et rouge) sont très proches les unes des autres. Ces trois courbes correspondent à des calculs utilisant la température de surface de première ébauche et des émissivités moyennées mensuellement. Les différentes estimations mensuelles de l’émissivité sont donc relativement équivalentes. On peut tout de même noter de légères différences très localisées. Les émissivités de la NASA semblent être légèrement meilleures aux alentours de 1150  $\text{cm}^{-1}$ . Ce résultat n’est pas surprenant car les émissivités de la NASA sont certes des moyennes mensuelles, mais calculées à partir de restitutions de sondages IASI pour l’année 2008, qui est celle considérée ici. Les autres bases d’émissivité sont liées aux données MODIS (mauvaise résolution spectrale) et pour l’année 2007 (seules les restitutions que nous effectuons utilisent les données de 2008).

L’avantage d’avoir une première ébauche issue de MODIS est que nous n’utiliserons qu’une seule fois les données IASI. Utiliser ces dernières deux fois (dans l’inversion et l’estimation de la première ébauche) pourrait entraîner des erreurs plus importantes.

Cependant, les différences observées sont minimales par rapport à la courbe noire. Il faut tout de même garder à l’esprit que la courbe noire est la seule à avoir été obtenue à partir de calculs de transfert radiatif utilisant la température restituée. L’importante différence entre cette courbe et les autres provient, partiellement, de la température de surface qui est mieux caractérisée une fois restituée, que dans la première ébauche. Sur une grande portion du spectre IASI, l’erreur peut être diminuée de 6 à seulement 2 K lorsque l’on caractérise mieux les caractéristiques de surface.

Ainsi que l’ont noté [Vogel et al. \(2011\)](#), il est très difficile de comparer des bases de données d’émissivités sans une connaissance précise de la température de surface. Nous avons donc mené une deuxième étude afin de mesurer précisément l’apport de la restitution d’émissivité. La Figure 2.14 présente à nouveau des statistiques de la différence ( $TB_{obs} - TB_{sim}$ ). La courbe bleue présente les statistiques en utilisant l’émissivité et la température de surface de première ébauche (courbe similaire à la bleue précédente), la courbe verte présente les statistiques utilisant les émissivités interpolées et la température de surface restituées et la

courbe rouge présente les statistiques utilisant les émissivités et la température de surface restituées. Les émissivités UWIREMIS et celles de la NASA donnent des résultats similaires avec la température de surface restituée, leurs courbes ne sont pas affichées par souci de clarté.

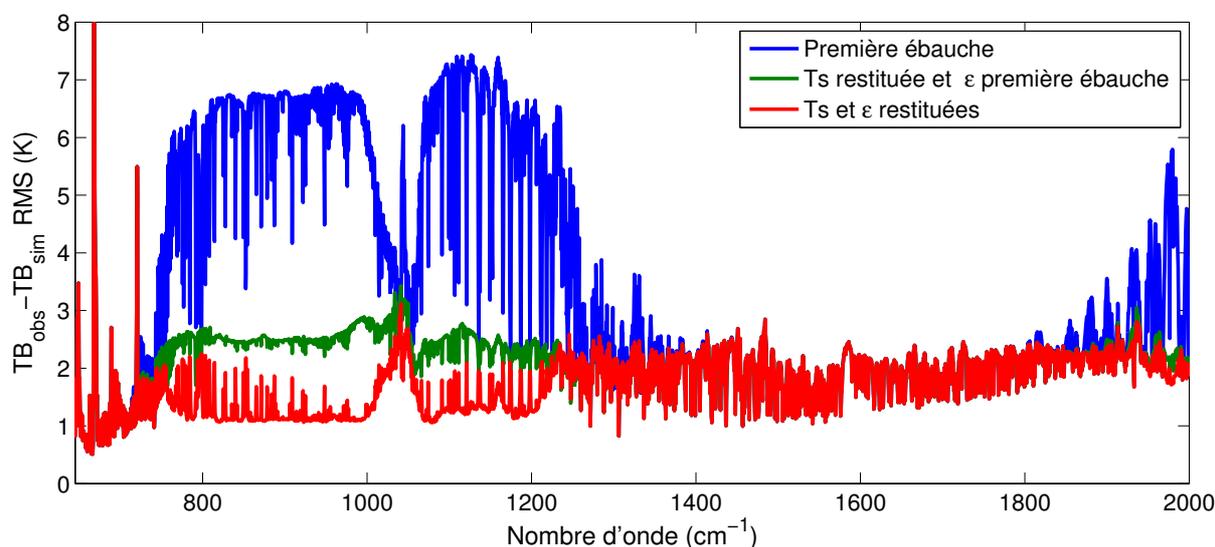


FIGURE 2.14 – RMS de la différence entre le spectre de température de brillance observée par IASI et un autre calculé à partir de la température de surface restituée et une émissivité constante à 0,98 (courbe bleue), les émissivités interpolées (courbe verte) et les émissivités restituées (courbe rouge).

On peut mesurer directement l’apport de la température de surface en comparant les courbes bleue et verte, et l’apport des émissivités restituées en comparant les courbes verte et rouge. Il y a une amélioration de 4 K liée uniquement à une meilleure estimation de la température de surface. L’amélioration de l’émissivité de la surface diminue de 1,5 K l’erreur par rapport à l’estimation faite avec une climatologie mensuelle. Bien qu’une mauvaise caractérisation de la température de surface est équivalente à une mauvaise caractérisation de l’émissivité, la plus grande variabilité de la température de surface implique de plus grande modification dans le calcul de transfert radiatif. Il est donc nécessaire de considérer les valeurs restituées de la température et de l’émissivité pour caractériser au mieux la surface.

On remarque ici encore les pics aux alentours de  $660\text{ cm}^{-1}$  liés à la mauvaise caractérisation de certains composants mineurs de l’atmosphère. On note également encore la bande d’absorption de l’ozone. Cependant le fait le plus remarquable est l’écart qui apparaît d’une courbe à l’autre. Les statistiques obtenues en utilisant une émissivité constante sont assez proches de ce qu’on avait précédemment en utilisant la température de surface de première ébauche.

L’erreur que l’on retrouve sur les canaux en prenant en compte les caractéristiques

de surface restituées est semblable à celle qui est présente dans les zones spectrales où l’atmosphère est opaque (notamment entre 1400 et 1800  $\text{cm}^{-1}$ ). Il est communément admis que notre connaissance de l’atmosphère est meilleure que celle des propriétés de surface. Le fait que l’erreur résiduelle ne descende pas en dessous nous permet de vérifier que la restitution des paramètres de surface n’a pas essayé de compenser artificiellement les erreurs de caractérisation de l’atmosphère pour approcher au mieux les températures de brillance réelles.

Restituer l’émissivité nous permet d’obtenir un produit plus proche de la réalité. On prend ainsi en compte toutes les variations de l’émissivité, tant temporellement que spatialement. [Mira et al. \(2007\)](#) ont montré que le changement de l’humidité d’un matériau pouvait modifier son émissivité jusqu’à près de 15%. De telles variations, liées aux intempéries ou à la végétation, ne sont pas prises en compte dans les climatologies mensuelles d’émissivité. Cette étude montre qu’il est nécessaire de bien caractériser la température et l’émissivité de surface pour réussir à interpréter les observations satellites. Les erreurs restantes entre les observations et les températures de brillance simulées seront ensuite particulièrement importantes pour restituer les profils atmosphériques.

### 2.3.4.3 Évaluation de la température de surface

Afin d’analyser la structure spatiale de la restitution et de mesurer l’erreur commise sur l’évaluation de la température de surface par notre algorithme, nous avons décidé de comparer la température restituée à une autre base de températures de surface indépendante.

Nous nous intéressons ici à la base restituée par le LSA/SAF (Land Surface Analysis/Satellite Applications Facility) à partir de mesures de SEVIRI (Spinning Enhanced Visible and InfraRed Imager). SEVIRI est un radiomètre à balayage sur la plateforme MSG (Meteosat Second Generation). Il s’agit d’un satellite géostationnaire au-dessus de l’Europe. SEVIRI fournit des mesures allant de l’est de l’Amérique du sud à l’est de la péninsule arabique et du sud de l’Afrique à la Norvège. À partir de ces mesures, un algorithme “split-window” permet de restituer la température de surface. Il s’agit de déduire la température de la surface directement à partir des différences de températures de brillance entre les canaux. Ces différences de températures de brillance permettent de s’affranchir des erreurs liées à une mauvaise caractérisation de l’atmosphère. Ils utilisent une estimation des paramètres atmosphériques à partir d’une base de données fixe et d’une émissivité de surface qu’ils restituent directement à partir de l’indice de végétation.

Les températures de surface issues de SEVIRI sont disponibles de 2007 à nos jours, avec une mesure toute les 15 minutes ([Trigo et al. 2011](#); [Caselles et al. 1997](#); [Wan and Dozier 1996](#)). Chaque sondage IASI de la base de 4 semaines de 2008 considérée précédemment est mis en coïncidence avec la mesure SEVIRI la plus proche. Nous ne gardons, évidemment, que les situations situées dans le champ de vision de SEVIRI.

Tant du point de vue des méthodes que des données utilisées, les températures de surface

du LSA/SAF et celles que nous avons restituées sont indépendantes. Les comparer entre elles permettra de valider l'une comme l'autre.

Nous calculons la RMS de l'erreur entre les deux bases de données ainsi que l'écart-type de l'erreur et le biais. Pour comparaison, nous effectuons le même travail avec la base de températures de surface de l'ECMWF, utilisée comme première ébauche dans notre restitution. Le tableau ci-après présente les résultats obtenus :

Bases de données	RMS	Biais	Écart-type
$T_s$ (ECMWF - LSA/SAF)	8,7	-6,0	6,3
$T_s$ (restituée - LSA/SAF)	4,2	2,3	3,5

La comparaison entre les températures de surface restituées et celles du LSA/SAF présente un biais et un écart-type plus faibles (de respectivement 62% et 44%) que la comparaison entre les températures de surface de l'ECMWF et du LSA/SAF. On retrouve une RMS de l'erreur plus faible (de 52%). On a donc largement amélioré la première ébauche de l'ECMWF avec notre inversion qui se rapproche d'une base de températures de surface totalement indépendante.

La Figure 2.15 présente l'écart entre la température de surface restituée et celle du LSA/SAF (dans la partie haute) et l'écart entre la température de surface de l'ECMWF et celle du LSA/SAF (dans la partie basse). Toutes les données sont compilées sur cette carte (les premières semaines de janvier, avril, juillet et octobre 2008) afin d'avoir une carte plus complète. Certaines zones peuvent alors ne pas être couvertes par la base de données, car, que ce soit les températures de surface que l'on a restituées ou celles fournies par le LSA/SAF, elles ne sont disponibles qu'en ciel clair.

On peut observer sur cette carte la géométrie de sondage de l'instrument SEVIRI. On devine un cercle centré sur l'Afrique, s'étendant de l'Amérique du sud à la péninsule arabe.

Un biais négatif important est clairement visible entre les températures de surface de l'ECMWF et celles du LSA/SAF, conformément aux calculs globaux que nous avons faits précédemment. Certaines structures spatiales apparaissent entre les températures de l'ECMWF et celles du LSA/SAF, et on ne les retrouve pas avec les températures restituées. On peut citer par exemple le biais positif au-dessus de l'Europe du nord que l'on ne retrouve pas sur la carte du haut. Les importantes différences, entre les températures de l'ECMWF et celles du LSA/SAF au-dessus du Sahara, sont nettement plus faibles (en valeur absolue) avec les températures restituées.

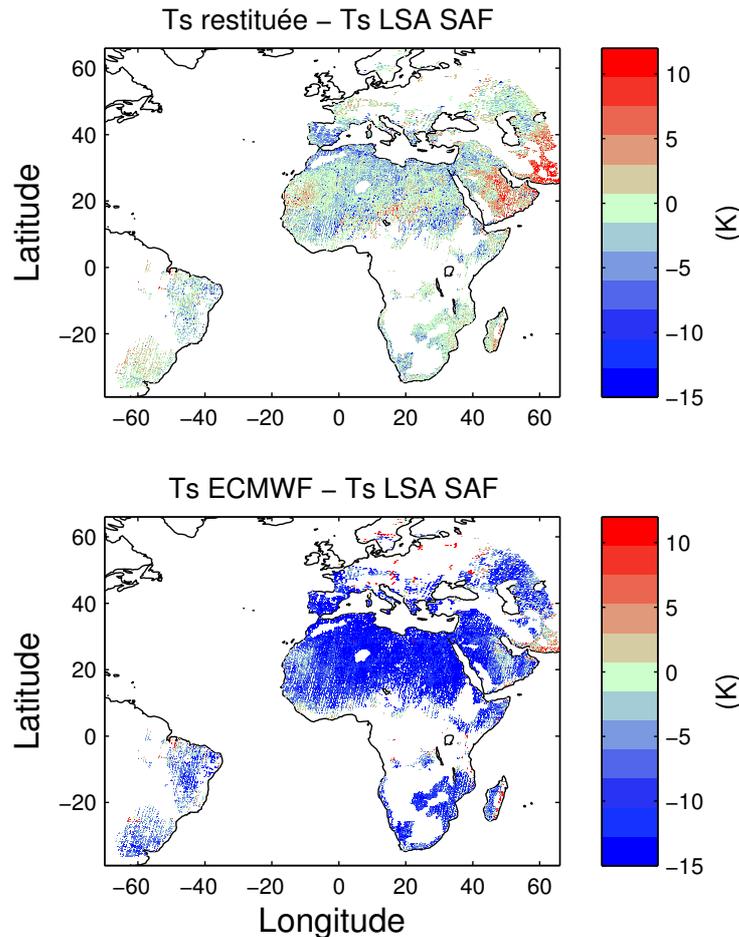


FIGURE 2.15 – Carte d'écart entre les températures de surface fournies par le LSA/SAF et les températures restituées (en haut) et celles de l'ECMWF (en bas), pour les quatre premières semaines de janvier, avril, juillet et octobre de 2008.

Au niveau des bords du disque de mesure de SEVIRI, plus particulièrement au niveau de la péninsule arabique et de l'Iran, d'importantes différences entre les températures de surface restituées et celles du LSA/SAF peuvent être notées. Ces erreurs sont liées, en partie, à des imprécisions dans les mesures de SEVIRI pour des angles de visée élevés. L'erreur sur l'estimation de la température de surface dans ces régions (notamment la péninsule arabique) est supérieure à 3 K (Freitas et al. 2010).

Aucune autre structure spatiale claire ne peut être identifiée sur la carte d'écart entre les températures de surface restituées et celles du LSA/SAF. Ceci nous indique que la restitution des paramètres de surface que l'on a mis en place n'est pas dépendante du type de surface ou de la structure géologique du sol.

### Triple colocalisation

Afin d'aller plus loin encore dans la comparaison de ces bases de données de tempéra-

tures de surface (ECMWF, LSA/SAF et celles restituées), nous avons décidé d'utiliser une méthode de validation supplémentaire appelée la triple colocalisation (Janssen et al. 2007; Stoffelen 1998). Cette méthode permet de facilement comparer trois bases de données et d'estimer leurs erreurs respectives sans connaître la valeur de la variable à estimer.

Soit  $X$ ,  $Y$  et  $Z$  trois estimations d'une variable  $T$ . La triple colocalisation fait l'hypothèse que ces estimations peuvent s'écrire comme la variable réelle multipliée par un coefficient d'erreur plus un terme d'erreur :

$$\begin{aligned} X &= \beta_X \cdot T + \epsilon_X \\ Y &= \beta_Y \cdot T + \epsilon_Y \\ Z &= \beta_Z \cdot T + \epsilon_Z \end{aligned}$$

où  $\beta_X$ ,  $\beta_Y$  et  $\beta_Z$  sont trois coefficients de calibration et  $\epsilon_X$ ,  $\epsilon_Y$  et  $\epsilon_Z$  sont les erreurs sur l'estimation de la variable  $T$ <sup>9</sup>.

On suppose que les erreurs de ces trois estimations sont non-corrélées. Afin de ne plus se soucier des coefficients de calibration, on utilise les variables  $X' = \frac{X}{\beta_X}$ ,  $Y' = \frac{Y}{\beta_Y}$  et  $Z' = \frac{Z}{\beta_Z}$ . Les coefficients de calibration seront calculés dans un deuxième temps. On a alors :

$$\begin{aligned} X' &= T + \epsilon_{X'} \\ Y' &= T + \epsilon_{Y'} \\ Z' &= T + \epsilon_{Z'} \end{aligned}$$

Afin d'estimer les erreurs de chacune des estimations, il faut transformer ces équations pour ne plus avoir la variable à estimer que l'on ne connaît pas. Il nous suffit pour cela de considérer les différences entre les différentes estimations :

$$\begin{aligned} X' - Y' &= \epsilon_{X'} - \epsilon_{Y'} \\ X' - Z' &= \epsilon_{X'} - \epsilon_{Z'} \\ Y' - Z' &= \epsilon_{Y'} - \epsilon_{Z'} \end{aligned} \tag{2.10}$$

On peut alors multiplier la première équation par la deuxième et calculer la moyenne de ce

---

9. Les hypothèses effectuées ici sont assez fortes. On suppose que la relation est linéaire, alors qu'il peut y avoir des cas de saturations. On suppose de plus que l'erreur est indépendante de la valeur de la température en elle-même, ce qui peut ne pas être le cas. On suppose également que les erreurs sont gaussiennes, ce qui nous permet de ne pas considérer les moments d'ordres supérieurs de l'erreur.

terme, on obtient :

$$\begin{aligned} \overline{(X' - Y') \cdot (X' - Z')} &= \overline{(\epsilon_{X'} - \epsilon_{Y'}) \cdot (\epsilon_{X'} - \epsilon_{Z'})} \\ &= \overline{\epsilon_{X'}^2 - \epsilon_{X'} \cdot \epsilon_{Z'} - \epsilon_{X'} \cdot \epsilon_{Y'} + \epsilon_{Y'} \cdot \epsilon_{Z'}} \\ &= \overline{\epsilon_{X'}^2 - \overline{\epsilon_{X'} \cdot \epsilon_{Z'}} - \overline{\epsilon_{X'} \cdot \epsilon_{Y'}} + \overline{\epsilon_{Y'} \cdot \epsilon_{Z'}}} \end{aligned}$$

Or on a supposé précédemment que les erreurs de chacune des estimations étaient décorré-  
lées, donc  $\overline{\epsilon_{X'} \cdot \epsilon_{Z'}} = \overline{\epsilon_{X'} \cdot \epsilon_{Y'}} = \overline{\epsilon_{Y'} \cdot \epsilon_{Z'}} = 0$ . On a alors :

$$\overline{(X' - Y') \cdot (X' - Z')} = \overline{\epsilon_{X'}^2}$$

On peut effectuer le même raisonnement en multipliant entre elles les autres équations du  
système (2.10), on obtient alors les erreurs des trois estimateurs :

$$\begin{aligned} \overline{\epsilon_{X'}^2} &= \overline{(X' - Y') \cdot (X' - Z')} \\ \overline{\epsilon_{Y'}^2} &= \overline{(Y' - X') \cdot (Y' - Z')} \\ \overline{\epsilon_{Z'}^2} &= \overline{(Z' - Y') \cdot (Z' - Y')} \end{aligned} \quad (2.11)$$

On peut donc directement estimer les trois erreurs à partir des moyennes des différences  
entre les trois bases de données si leurs erreurs ne sont pas corrélées.

Il faut désormais estimer les coefficients de calibration. La variable réelle  $T$  n'est pas  
connue, il faut donc considérer que l'un des coefficients de calibration vaut 1, prenons par  
exemple  $X' = X$ . Le résultat que l'on obtient ne dépend pas du choix de la variable dont le  
coefficient de calibration est fixé à l'avance. On connaît les erreurs  $\epsilon_{X'}$ ,  $\epsilon_{Y'}$  et  $\epsilon_{Z'}$ , on peut  
donc calculer directement les coefficients de calibration de  $Y$  et  $Z$  grâce à une régression  
neutre telle que présentée par Marsden (1999), on obtient :

$$\beta_Y = \frac{B + \sqrt{B^2 - 4 \cdot A \cdot C}}{2 \cdot A}$$

où  $A = \frac{\overline{\epsilon_X^2}}{\overline{\epsilon_Y^2}} \cdot \overline{X \cdot Y}$ ,  $B = \overline{X^2} - \frac{\overline{\epsilon_X^2}}{\overline{\epsilon_Y^2}} \cdot \overline{Y^2}$  et  $C = -\overline{X \cdot Y}$ . On obtient le même résultat pour  $\beta_Z$   
en remplaçant  $Y$  par  $Z$ .

Il faut connaître les différentes erreurs pour pouvoir déterminer les coefficients de calibra-  
tion. Il est donc courant d'effectuer ce calcul par itérations. On calcule les erreurs  $\epsilon_X$ ,  $\epsilon_Y$  et  
 $\epsilon_Z$  en considérant que les coefficients de calibration valent tous l'unité ( $\beta_X = \beta_Y = \beta_Z = 1$ ).  
À partir de cette première estimation des erreurs, on calcule les deux coefficients de cali-  
bration (le troisième étant fixé à un). On peut désormais réévaluer la valeur des erreurs en  
prenant en compte ces coefficients de calibration (*i.e.*, on considère les estimations  $Y' = \frac{Y}{\beta_Y}$   
et  $Z' = \frac{Z}{\beta_Z}$ ). Ces étapes sont répétées jusqu'à convergence des coefficients  $\beta$ .

Les coefficients de calibration obtenus par ce calcul sont relatifs, car on ne connaît pas

*T a priori*. Il faut fixer un coefficient pour connaître les autres. Cela veut donc dire qu'on ne pourra pas obtenir *T* au final, mais juste une estimation des erreurs.

Avant d'effectuer ces calculs sur les bases de données de températures de surface que l'on considère, il faut faire en sorte qu'elles ne soient pas biaisées. Pour cela, comme les températures de surface restituées par le LSA/SAF sont assez largement utilisées par la communauté, nous considérons que cette base n'est pas biaisée. Conformément aux calculs de biais précédent, on ajoute donc 6,0 K aux températures de l'ECMWF et on en enlève 2,3 à celles restituées. On a alors trois bases de données de températures de surface non-biaisées.

La première itération de la méthode de triple colocalisation nous donne une mesure des erreurs :

$$\begin{aligned}\epsilon_{LSA/SAF} &= 3,1 \\ \epsilon_{ECMWF} &= 5,4 \\ \epsilon_{Restituée} &= 1,5\end{aligned}$$

Le calcul des coefficients de calibration des températures de l'ECMWF ou celles restituées donne 1 à  $10^{-2}$  près. On considère donc que les coefficients de calibration valent tous l'unité. On retrouve le même résultat en considérant que la base de l'ECMWF a un coefficient de calibration unitaire et en calculant les deux autres, idem avec les températures de surface restituées. Le fait que ces trois coefficients de calibration soient unitaires confirme une nouvelle fois la consistance entre ces trois bases de données.

Les erreurs estimées à la première itération sont donc finales dans notre processus. Le fait que les températures de surface restituées soient plus proches de celles du LSA/SAF que de celles de l'ECMWF entraîne une erreur supérieure pour les températures de l'ECMWF. L'erreur très faible obtenue pour les températures de surface restituées est sans doute sous-estimée. Elle reste néanmoins nettement inférieure au deux autres. Du fait de la corrélation entre les températures restituées et celles de l'ECMWF (utilisées comme première ébauche), l'erreur calculée ici est inférieure à l'erreur réelle. Regarder les différentes corrélations entre ces bases permet d'arriver à mieux comprendre les différents liens entre elles.

$$\begin{aligned}Cor(ECMWF,LSA/SAF) &= 0,78 \\ Cor(ECMWF,Restituées) &= 0,83 \\ Cor(Restituées,LSA/SAF) &= 0,94\end{aligned}$$

Bien que les températures de surface de l'ECMWF aient servi de première ébauche à la restitution, les températures restituées sont plus corrélées avec celles du LSA/SAF. Les températures de surface restituées présentent les meilleures corrélations avec l'une et l'autre

base, ce qui explique que la triple colocalisation ait favorisé cette base<sup>10</sup>.

## 2.4 Conclusion

Nous avons créé un schéma d'inversion de la température et des émissivités de surface à partir des observations IASI, d'une première ébauche de l'émissivité de surface (issue d'une climatologie), d'une première ébauche de la température de surface et d'informations sur l'état de l'atmosphère.

Nous avons montré que la prise en compte des caractéristiques de surface précises permet de réduire l'écart entre les observations simulées et calculées de 5,5 K, ce qui permet de se rapprocher des erreurs liées à une mauvaise estimation de l'atmosphère.

Les inversions de température que nous avons effectuées ici, à partir des re-analyses de l'ECMWF, se sont fortement rapprochées des données du LSA/SAF, qui sont pourtant indépendantes. Nous allons désormais chercher à rendre notre algorithme plus opérationnel, car il dépend des re-analyses de l'ECMWF qui ne sont pas disponibles en temps réel, par définition. Il faut également valider les émissivités infrarouges hyperspectrales que nous avons restituées. Nous avons effectué une validation à partir de calculs de transfert radiatif, il faut désormais valider les émissivités de façon plus physique.

---

10. Les erreurs calculées ici par la méthode de triple colocalisation ne doivent pas être considérées comme absolues. Les fortes hypothèses de départ (non-corrélation des erreurs ; linéarité du modèle d'erreur) ne sont pas forcément exactes. Cependant, ce calcul montre la cohérence entre ces trois bases et tend à indiquer une meilleure qualité pour les températures de surface que l'on a restituées et celles du LSA/SAF.




---

# EXPLOITATION D'UNE CHAÎNE OPÉRATIONNELLE DE RESTITU- TION DE LA TEMPÉRATURE ET DE L'ÉMISSIVITÉ INFRAROUGE DE LA SURFACE

## Sommaire

<b>3.1</b>	<b>Mise en place d'une chaîne opérationnelle</b>	<b>84</b>
<b>3.2</b>	<b>Restitutions en pseudo temps réel</b>	<b>84</b>
3.2.1	Étude spatiale	85
3.2.2	Étude spectrale	88
<b>3.3</b>	<b>Moyennes mensuelles</b>	<b>90</b>
3.3.1	Variabilité spatiale	91
3.3.2	Variabilité temporelle	95
<b>3.4</b>	<b>Comparaison avec des radiosondages au-dessus du dôme C</b>	<b>98</b>
3.4.1	Données disponibles et objectif	98
3.4.2	Méthode	98
3.4.3	Résultats	99
<b>3.5</b>	<b>Perspectives</b>	<b>104</b>
<b>3.6</b>	<b>Conclusion</b>	<b>105</b>

---

Nous avons donc montré la robustesse de l'algorithme que nous avons mis en place. L'objectif est désormais de faire de cet algorithme un outil utile, modulable et capable de fonctionner en "pseudo temps réel". Nous utilisons ici le terme pseudo temps réel car l'algorithme nécessite un certain nombre de données en entrée (re-analyse de l'ECMWF), auxquelles nous n'avons pas accès, seules les données satellites sont disponibles en temps réel.

Ce chapitre a été volontairement séparé du chapitre précédent, malgré une thématique commune. Cette séparation permet au lecteur de bien assimiler le changement de données

utilisées en entrées de l'algorithme. Il s'agit dans ce chapitre d'adapter l'algorithme décrit et validé précédemment pour le rendre opérationnel. Nous analyserons ensuite les résultats obtenus.

### 3.1 Mise en place d'une chaîne opérationnelle

Nous avons utilisé précédemment des re-analyses de l'ECMWF comme données de première ébauche pour la température de surface et l'état de l'atmosphère. Ces données sont issues de révisions des sorties de modèles de prévision grâce aux observations *in situ*. Nous cherchons à effectuer des restitutions les plus indépendantes possible du modèle, afin de pouvoir effectuer des comparaisons par la suite.

Dans cette optique, nous avons opté pour les L2 IASI fournies par l'EUMETSAT (EUropean organisation for the exploitation of METeorological SATellites). Il s'agit de restitutions en temps réel, à partir notamment des observations IASI d'octobre 2007 à nos jours. Une description plus poussée de ces données, de la façon dont elles sont restituées et de leur contenu exact est fourni à la Section 6.2.1 (page 160). Nous soulignerons simplement ici que ces données nous fournissent une première ébauche en température de surface et pour les profils atmosphériques. Ils sont indépendants des modèles de prévision, contrairement aux re-analyses utilisées précédemment.

Notre choix s'est porté sur ces données car elles sont disponibles facilement sur un serveur, en coïncidence directe avec les mesures de IASI. En effet, chacune des situations issues des L2 est une restitution d'un sondage de IASI (appelé L1). Les deux jeux de données sont donc en coïncidence temporelle et spatiale exactes.

Nous avons dû faire face à de nombreux problèmes, liés notamment à la qualité du flag nuageux. Une mauvaise identification des situations claires peut entraîner d'importantes erreurs et une instabilité de la restitution. Inversement, un flag nuageux trop restrictif diminuerait le nombre de restitutions effectuées et donc la couverture globale de nos restitutions.

Il a fallu également adapter notre algorithme afin de le rendre compatible avec différents formats de données, car ces dernières ont évoluées au cours du temps au sein même des L2 d'EUMETSAT (nombre de niveaux de pression des profils différents, unité des différentes variables...).

### 3.2 Restitutions en pseudo temps réel

Nous disposons au final de restitutions de températures et d'émissivités de surface d'octobre 2007 à nos jours. Chaque situation identifiée comme étant en ciel clair et au-dessus des surfaces continentales est considérée pour la restitution. Nos restitutions sont disponibles par fichiers journaliers, contenant, pour chaque situation restituée, la température de surface, l'émissivité ainsi que la matrice de covariance d'erreur des données restituées.

### 3.2.1 Étude spatiale

La Figure 3.1 présente des restitutions obtenues pour la journée du 1 mars 2013. Les deux cartes de la partie gauche de la figure représentent l'émissivité de la surface à  $1100\text{ cm}^{-1}$ , ce qui correspond, d'après la Section 2.1.2 (page 42), à une des bandes d'absorption des silicates. Les cartes de droite représentent quant à elles l'émissivité à  $1500\text{ cm}^{-1}$ , ce qui correspond à une bande d'absorption des carbonates. Les deux cartes dans la partie supérieure correspondent à l'émissivité de première ébauche interpolée à partir de MODIS et utilisée en entrée de l'algorithme. Les deux cartes dans la partie basse de l'image correspondent à l'émissivité restituée.

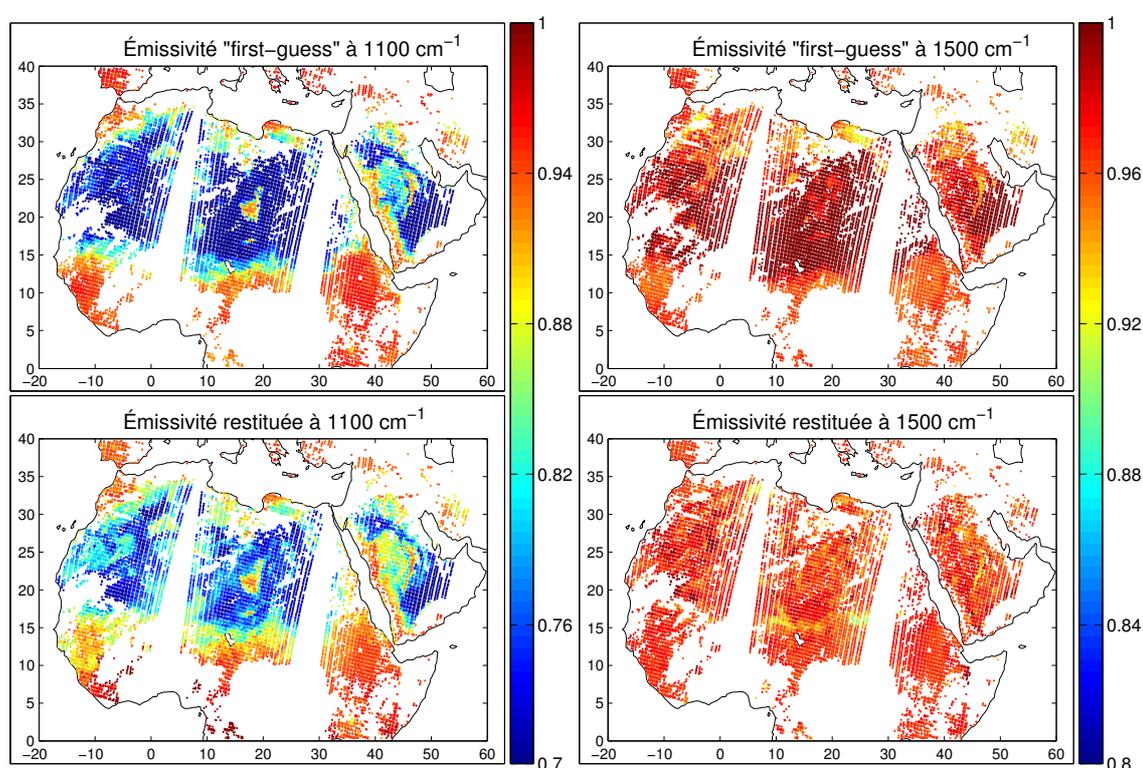


FIGURE 3.1 – Carte d'émissivités de première ébauche (partie haute) et restituées (partie basse) à  $1100\text{ cm}^{-1}$ (gauche) et  $1500\text{ cm}^{-1}$ (droite).

Chaque point sur la carte correspond à une observation IASI. On peut constater qu'il y a eu trois survols de la zone. Le premier survol, au-dessus de la péninsule arabique, a eu lieu à environ 6h30 heure universelle ( $\approx 9\text{h}30$  heure solaire locale). Le deuxième, au-dessus de la Libye, du Niger et du Tchad, a eu lieu à environ 8h15 heure universelle ( $\approx 9\text{h}30$  en heure solaire locale, MetOp est un satellite héliosynchrone donc il passe à la même heure en heure solaire locale dans son orbite descendante). Le troisième survol, sur l'ouest de l'Afrique, a eu lieu à environ 10h00 heure universelle (soit encore  $\approx 9\text{h}30$ ). Ces trois survols ont duré environ 7 minutes chacun. Le temps qui sépare chacun de ces survols correspond au temps de

rotation de MetOp autour de la Terre (1h40 environ). Il y a eu d'autres survols de l'Afrique du nord plus tard dans la journée (vers 18h00, 19h30 et 21h00, vers 21h en heure locale, il s'agit d'orbites montantes) mais on ne les représente pas ici pour ne pas superposer les différentes orbites.

Les cartes de droite représentent l'émissivité de la surface au niveau d'une bande d'absorption des carbonates. Cette bande est assez peu marquée, c'est pourquoi il est plus compliqué de voir ici des signatures claires (l'échelle n'est pas la même que pour les cartes de droite). Les régions connues pour être composées en majeure partie de carbonates sont (Jiménez et al. 2010) :

- La pointe sud de la péninsule arabique, non visible ici car considérée comme nuageuse (par le flag nuage d'EUMETSAT) ;
- Une virgule en travers de la péninsule arabique visible sur les deux cartes ;
- Une petite zone dans le Darfour, entre le Soudan et le Tchad. Cette zone n'est pas visible sur la carte d'émissivité de première ébauche, mais sur la carte des émissivités restituées qui présentent un minimum local exactement sur cette zone (en bas à droite de l'orbite centrale).

Les zones déjà marquées par la signature des carbonates ont donc été conservées dans les émissivités restituées. Mais certaines zones non identifiées par les émissivités de première ébauche (interpolées à partir de MODIS) apparaissent clairement sur la carte des émissivités restituées. Cette signature spectrale n'a pas été fournie en entrée à l'algorithme d'inversion qui arrive pourtant à la créer. Cela constitue donc une bonne validation de nos restitutions.

Les cartes de gauche, correspondant à l'émissivité dans une bande d'absorption des silicates, présentent des structures très marquées. En effet, le sable du Sahara, composé en grande partie de silicates, montre une forte signature dans cette bande. On peut remarquer sur ces cartes l'importante différence entre l'émissivité du sable avoisinant 0,7 et celle des zones végétalisées (plus au sud) avoisinant 0,95. Le mont Tibesti, au milieu du Tchad, est également très visible sur ces cartes. On retrouve aussi la "virgule" de carbonates sur la péninsule arabique. Elle apparaît ici comme une zone à l'émissivité supérieure, car elle ne présente pas de signature à  $1100\text{ cm}^{-1}$ .

On note cependant de légères différences entre ces deux cartes. L'ouest africain, par exemple, présente de nombreuses structures sur la carte du bas (émissivité restituée). Ces structures ne sont pas visibles sur les émissivités de première ébauche. Afin de déterminer quelles structures sont les plus représentatives de la surface réelle, deux solutions s'offrent à nous. Nous pouvons comparer ces motifs spatiaux aux structures géologiques connues, ce qui a été fait précédemment avec les carbonates en utilisant les informations de Jiménez et al. (2010). Nous pouvons également regarder une photographie de la région dans le spatial afin d'identifier finement les structures. Cette méthode fonctionne particulièrement bien avec les silicates car ils composent la majeure partie des roches. On peut alors facilement identifier visuellement les zones composées de silicates.

La Figure 3.2 représente une imagerie satellite de l'Afrique du nord issue de Google Earth. Sur cette carte dans le visible, nous avons superposé les émissivités restituées à  $1100\text{ cm}^{-1}$  présentées précédemment. Il s'agit donc des restitutions de la matinée du 1<sup>er</sup> mars 2013. L'échelle utilisée pour le code couleur est la même que précédemment, le bleu correspondant à une émissivité de 0,7 et le rouge à une émissivité de 1.

On peut remarquer sur cette carte la coïncidence entre les motifs spatiaux observables dans le visible et les structures d'émissivité. Le mont Tibesti est encore une fois bien visible. Les structures spatiales sur la péninsule arabique coïncident également. Ces dernières étaient similaires dans les émissivités restituées et les émissivités de la première ébauche. Pour évaluer l'algorithme de restitution, les structures qu'il nous faut valider sont celles au-dessus de l'Afrique de l'ouest qui différaient entre les deux cartes précédentes. Ici encore, on note une très bonne corrélation entre les structures visibles et celles d'émissivité. Les structures présentes sur la carte d'émissivités restituées à  $1100\text{ cm}^{-1}$  sont plus réalistes que celles de la première ébauche. On valide ainsi une nouvelle fois l'apport de notre algorithme.

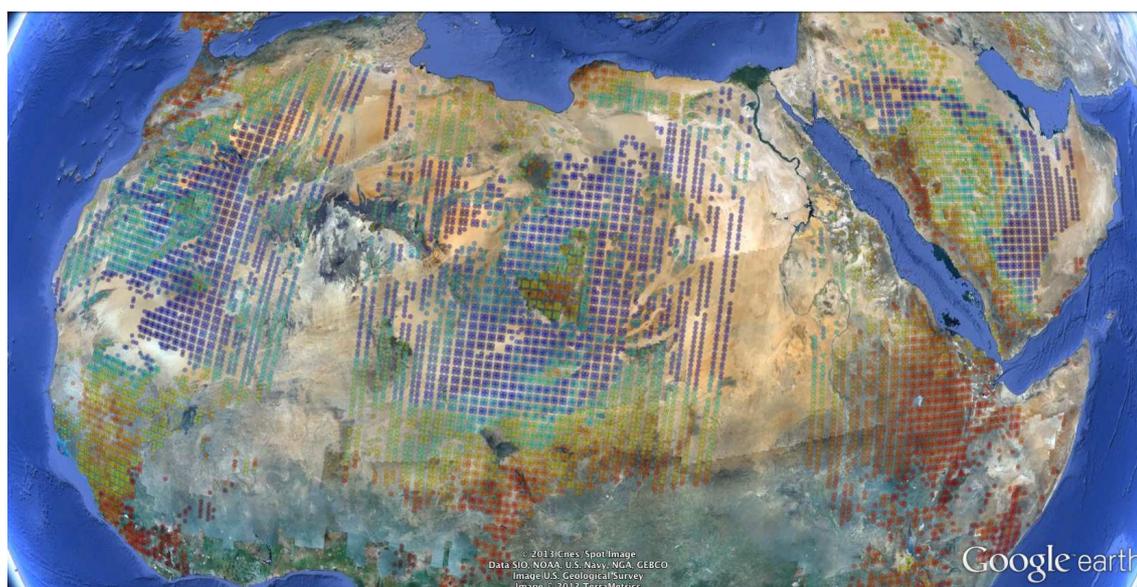


FIGURE 3.2 – Représentation des émissivités restituées à  $1100\text{ cm}^{-1}$  pour le matin du 1<sup>er</sup> mars 2013 au-dessus de l'Afrique du Nord. © Google Earth pour la photo.

Ces cartes correspondent à des restitutions directes à partir de sondages IASI. La géométrie de l'instrument, qui sonde en balayage cross-track, est visible. Chaque bande cross-track représente un balayage (par séries de 4 pixels). Il est intéressant de noter qu'il ne semble pas y avoir de dépendance de l'émissivité restituée à l'angle de visée du satellite. Les émissivités au bord de la fauchée sont semblables à celle au nadir. Cette observation vient une nouvelle fois valider notre algorithme.

### 3.2.2 Étude spectrale

Étudier les caractéristiques spectrales des émissivités restituées est complexe. Chaque spectre est constitué de 8461 valeurs, il est donc difficile de les représenter en totalité sur une carte. Calculer des statistiques d'écart entre différentes bases de données d'émissivités ou entre des calculs de transfert radiatif utilisant plusieurs bases d'émissivité en entrée (comme il a été fait précédemment) ne suffit pas. Nous avons donc décidé de comparer les émissivités de plusieurs bases de données en parallèle en prenant en compte le type de surface.

Pour cela nous utilisons la classification de surface en 17 classes de l'IGBP (voir Section 2.2.2, page 51). Les bases d'émissivité que nous comparons sont :

- Les émissivités restituées à partir de MODIS (voir Section 2.1.2.2, page 44) de l'année 2007, on utilise les moyennes mensuelles de janvier, avril, juillet et octobre ;
- Les émissivités interpolées à partir de celles de MODIS (voir Section 2.2.5, page 59), on utilise les mêmes mois ;
- Les moyennes mensuelles d'émissivités restituées par la NASA à partir de sondages IASI de l'année 2008, ici encore on utilise les 4 mois répartis dans l'année ;
- Les émissivités restituées par notre algorithme, afin d'avoir un nombre comparable de situations et une variabilité annuelle approchant celles des autres bases, nous utilisons ici toutes les données restituées pour la première semaine de juillet 2012 et la première semaine de janvier 2013.

Au final chaque base de données est constituée d'environ 800.000 spectres d'émissivité. Si les trois dernières bases sont hyperspectrales, la première ne dispose que de 6 valeurs d'émissivité à 833,3 ; 909,1 ; 1162,8 ; 2500 ; 2564 et 2631,6  $\text{cm}^{-1}$ . Nous utilisons cette base de données car c'est une base largement utilisée, reconnue et validée dans la communauté. Nous ne pouvons donc comparer les différentes bases de données qu'au niveau de ces 6 canaux.

Nous choisissons de comparer ces bases d'émissivité sur les canaux à 833.3 et 1162.8  $\text{cm}^{-1}$ . Le premier canal correspond à un canal où l'émissivité reste élevée, tandis que le deuxième est dans une bande spectrale des silicates. La comparaison entre les valeurs de ces deux émissivités sur les quatre bases sera intéressante, puisque nous pourrions analyser leur stabilité respective (premier canal) ainsi que leur capacité à identifier les structures les plus importantes (deuxième canal).

Nous allons calculer des histogrammes de dispersion de l'émissivité à ces deux nombres d'onde, pour les quatre bases de données, en fonction du type de surface. Présenter 17 histogrammes correspondant aux 17 classes de l'IGBP serait rébarbatif et complexe à analyser. D'autant que certaines classes sont très semblables, du point de vue de l'émissivité infra-rouge. Nous considérons alors une nouvelle classification en regroupant certaines classes de l'IGBP. Nous avons choisi de considérer que trois catégories de surface :

- La neige et la glace, correspondant à la catégorie neige et glace de l'IGBP (16) ;
- Les surfaces végétalisées, correspondant aux catégories forêt résineuse (2), forêt de feuillus à feuilles persistantes (3), forêt de conifères à feuilles persistantes (4), forêt de

feuillus à feuilles caduques (5), forêt mixte (6), broussaille dense (7), savane boisée (9), savane (10), prairie (11), marécage permanent (12), terre cultivée (13) et mosaïque de terres cultivées ou naturelles (15) de l'IGBP ;

- Les surfaces arides ou semi-arides correspondant aux catégories broussaille clairsemée (8) et désertique ou quasi désertique (17) de l'IGBP.

Pour chacun de ces trois types de surfaces, nous calculons les histogrammes de répartition de l'émissivité aux deux longueurs d'onde considérées. Les résultats obtenus sont présentés sur la Figure 3.3.

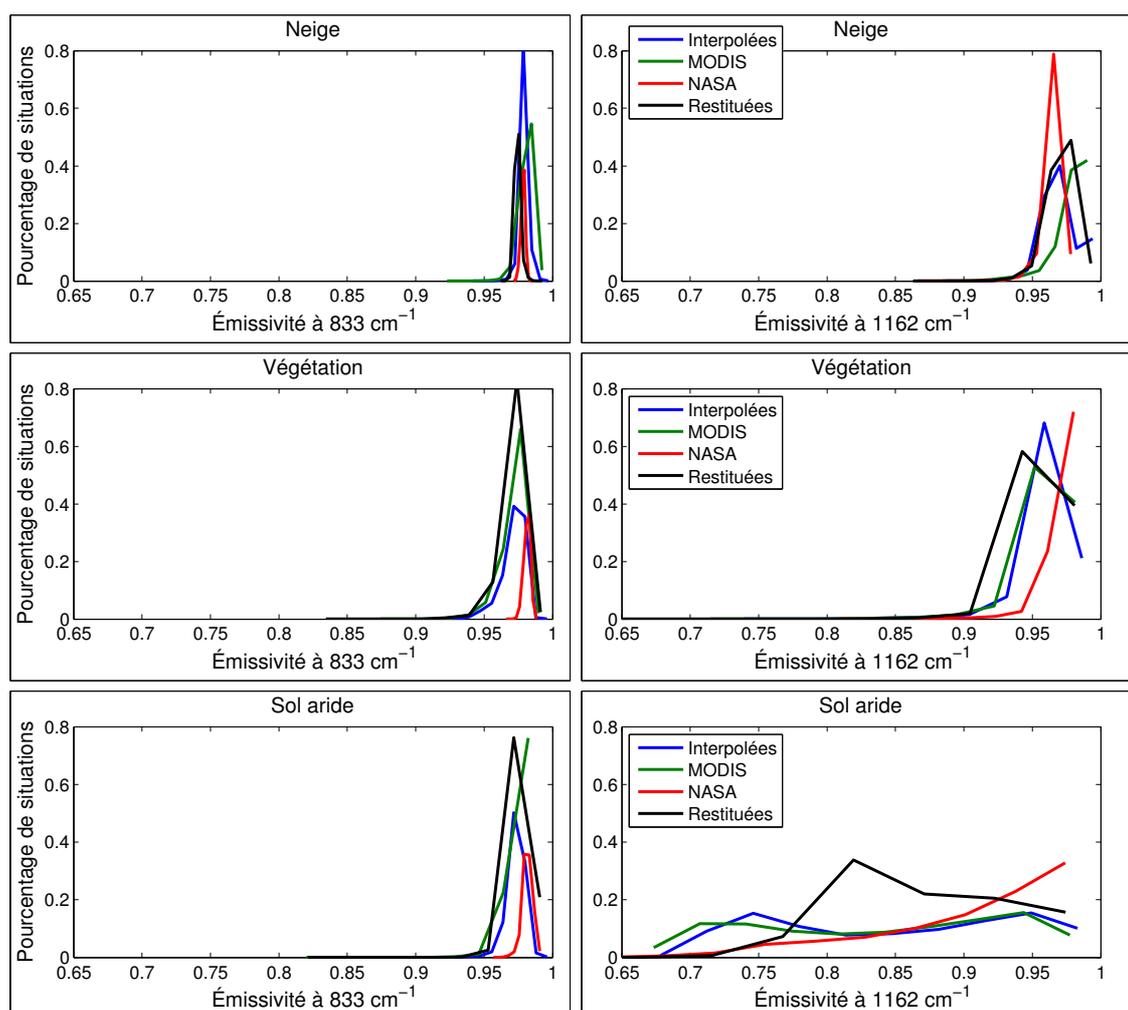


FIGURE 3.3 – Histogrammes des émissivités à 833 (à gauche) et 1160  $\text{cm}^{-1}$  (à droite), pour des surfaces neigeuses ou glacées (tout en haut), des surfaces végétalisées (au milieu) et arides (en bas). Les statistiques sont présentées pour les émissivités restituées à partir de MODIS (en vert), interpolées à partir de MODIS (en bleu), celles restituées par la NASA (en rouge) et celles que nous restituons (en noir).

Nous avons volontairement utilisé la même échelle pour tous les histogrammes afin que le lecteur prenne conscience de la faible dispersion des émissivités sur certaines surfaces ou à

certaines longueurs d'onde. Utiliser une échelle réduite entrainerait l'apparition de certaines différences qui sont d'un ordre très inférieur à la dispersion globale des bases de données.

On peut remarquer la très faible dispersion des émissivités à  $830\text{ cm}^{-1}$ , quelque soit la surface considérée. C'est une zone du spectre très stable, où peu de composantes du sol viennent interférer, il y a donc quasi systématiquement une émissivité de 0,98. Le fait que les émissivités que l'on restitue présentent également cette caractéristique est encore une fois un signe de leur validité.

Les émissivités à  $1160\text{ cm}^{-1}$  sont plus variées. Elles présentent une très faible variabilité sur les zones recouvertes de neige ou de glace, pour les quatre bases de données. Leur variabilité est supérieure mais tout de même faible au-dessus des zones végétalisées. Encore une fois, on note une grande similitude entre les quatre bases de données. La base de données de la NASA semble tout de même présenter des émissivités plus élevées, avoisinant plutôt les 0,97, alors que les autres restent autour de 0,95. Ceci est certainement dû à la méthode de restitution des émissivités qui, à l'aide d'une transformation logarithmique, contraint les émissivités à un intervalle de variation restreint.

Au-dessus des sols arides, on note des différences dans la répartition des émissivités à  $1162\text{ cm}^{-1}$ . Les émissivités de MODIS et celles interpolées à partir de ces dernières sont par construction très proches et ceci se retrouve sur l'histogramme. Les données de la NASA semblent ici encore présenter des valeurs nettement supérieures, pour les mêmes raisons que précédemment (*i.e.*, la contrainte logarithmique dans leur algorithme). Les émissivités que nous restituons présentent des valeurs légèrement supérieures également. Cette différence plus faible peu s'expliquer par la différence entre une émissivité issue de climatologie (moyennes mensuelles) et une émissivité restituée directement à partir d'un sondage. En effet, une variation de l'humidité du sol augmentera nettement son émissivité (jusqu'à atteindre les valeurs correspondant aux sols végétalisés). Toutefois, les différences observées ici restent faibles et leur dispersion semble équivalente.

Cette étude vient confirmer plus encore la validité des émissivités restituées par notre algorithme à partir des sondages IASI.

### 3.3 Moyennes mensuelles

À partir des restitutions d'émissivités de chaque sondage IASI en ciel clair au-dessus des continents, nous cherchons à créer une base de données d'émissivités infrarouges hyperspectrales, facilement exploitable. L'intérêt est de pouvoir étudier les variations spatiales, temporelles et spectrales de l'émissivité des surfaces.

Une telle étude est compliquée à mener sur la base de données en pseudo temps réel. La présence éventuelle de nuages entraîne l'absence de données. La grande quantité de données quotidiennes rend ardues les études spectrales et spatiales simultanées. Nous avons donc décidé de compiler toutes les données disponibles depuis le 01/10/2007 jusqu'à aujourd'hui.

### 3.3.1 Variabilité spatiale

Pour rassembler ces données, la base d'émissivités journalières a été projetée sur la grille "equal-area" à  $0,25^\circ \times 0,25^\circ$ , déjà utilisée précédemment. Pour chaque pixel continental, les moyennes mensuelles globales et annuelles des restitutions de l'émissivité sont calculées. On obtient ainsi une climatologie mensuelle de l'émissivité ainsi que les matrices de covariance d'erreur théoriques associées. On calcule également la dispersion intra-mensuelle des différentes restitutions sur ce pixel.

La Figure 3.4 présente l'émissivité restituée à  $1120 \text{ cm}^{-1}$ , moyennée sur les mois de janvier de 2008 à 2013. La partie supérieure de la figure montre la première ébauche utilisée dans les restitutions. On voit que l'on se situe au niveau d'une bande d'absorption des silicates. Le graphique du bas de la figure présente un spectre d'absorption typique d'un silicate (en bleu), ainsi qu'un spectre d'absorption d'un carbonate (en vert) et un spectre de végétation (en rouge). La barre verticale noire indique la zone du spectre où l'on se situe. On remarque que les silicates sont les seuls à présenter une signature spectrale nette dans cette zone.

Sur la carte supérieure, la limite entre les données interpolées à partir des restitutions de MODIS et celles restituées par la NASA est bien visible au nord. En effet, comme expliqué à la Section 2.2.5 (page 59), pour tous les pixels pour lesquels les restitutions de MODIS n'étaient pas disponibles, nous avons complété notre carte de première ébauche avec les restitutions de la NASA. Au mois de janvier, le pôle nord reste dans la nuit, il n'y a pas d'alternance jour/nuit, donc pas de restitutions possibles pour MODIS. Il en est de même pour l'Antarctique, même si la limite entre les deux bases de données est moins visible.

La partie basse de la figure présente les moyennes par pixel d'émissivité restituée. Sur cette carte, la limite nette entre les deux types de premières ébauches utilisés n'est plus visible au nord. Aucune contrainte spatiale n'est utilisée dans les restitutions, cette meilleure cohérence spatiale de la restitution vient une fois de plus valider notre algorithme. La discontinuité que l'on peut remarquer sur l'Antarctique n'est plus vraiment en rapport avec la discontinuité de la première ébauche. Elle correspond plutôt à la variation en altitude de l'Antarctique<sup>1</sup>.

Nous ne reviendrons pas ici sur les structures spatiales de silicates, plus en accord avec les structures géologiques pour les émissivités restituées, car ce point a déjà été traité précédemment. Il est tout de même intéressant de noter les structures qui apparaissent au-dessus de l'Amazonie. Dans la première ébauche, il semble y avoir un léger bruit, avec des valeurs assez faibles de l'émissivité sur certaines zones. Ici encore, la complétion des zones nuageuses par les émissivités de la NASA engendre des structures étranges sur la base de première ébauche. Les émissivités restituées présentent, quant à elles, une forte valeur, il en est de même au-dessus de la forêt tropicale africaine et de l'Indonésie. Ces maxima localisés au niveau des tropiques (zones particulièrement nuageuses) semblent correspondre à

1. Le niveau de l'Antarctique est très bas proche de l'océan et monte rapidement à plus de 2000 m

une pollution nuageuse. En effet, comme il a été expliqué précédemment, le flag nuageux d'EUMETSAT utilisé dans notre schéma n'est pas idéal. Certaines situations sont restituées à tort, ces zones semblent en pâtir. Ce flag a été amélioré depuis, il serait intéressant de mesurer l'impact de ce dernier sur les restitutions.

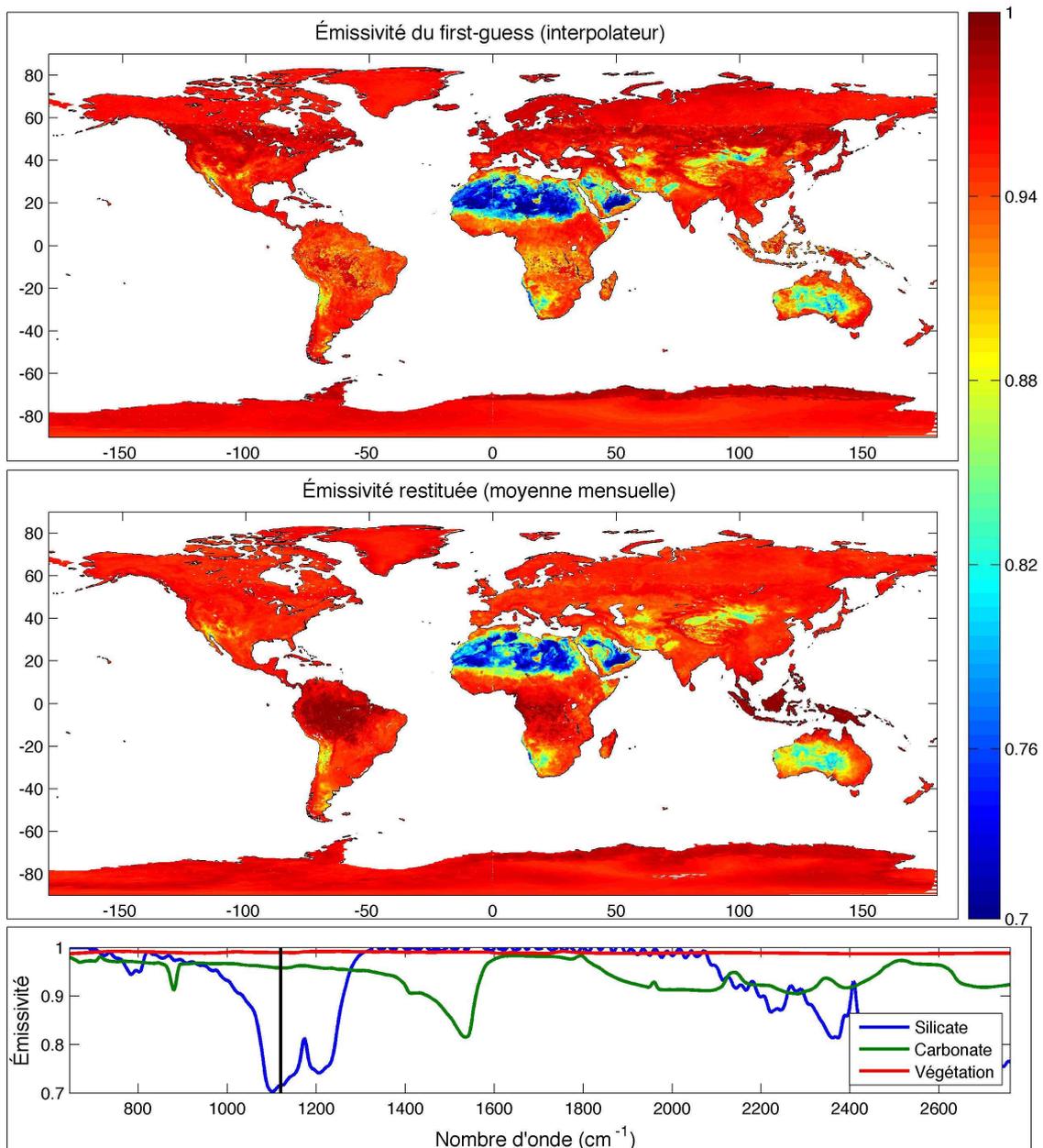


FIGURE 3.4 – Émissivité moyenne en janvier à  $1120 \text{ cm}^{-1}$ , la première ébauche interpolée à partir des émissivités de MODIS est en haut, les émissivités restituées sont au centre. Les trois spectres du bas représentent des spectres d'émissivité correspondant à des sols arides (silicates en bleu et carbonates en vert) et à un sol végétalisé (en rouge).

Afin de mieux étudier ces contaminations nuageuses, il faut se pencher sur l'étude de l'erreur de restitution théorique (disponible avec les restitutions, voir l'équation (2.9)) et sur la dispersion inter-annuelle de l'émissivité. La Figure 3.5 présente, dans sa partie haute, la carte de l'erreur théorique de restitution à  $1120\text{ cm}^{-1}$  moyennée sur les mois de janvier de 2008 à 2012. Les spectres de silicates, carbonates et de végétation sont tracés dans la zone inférieure de la figure afin de garder à l'esprit la zone du spectre étudiée. Les erreurs théoriques de restitution sont très faibles, de l'ordre de  $10^{-2}$  sur tout le globe. Certaines zones semblent présenter des statistiques d'erreur supérieures, notamment l'Inde, le sub-Sahel, le nord de la Russie et le Zimbabwe. Les erreurs restent cependant relativement faibles.

La partie basse de la figure présente la dispersion mensuelle des émissivités au mois de janvier sur les 5 années. Cette dispersion mensuelle prend en compte plusieurs choses : la variabilité réelle de l'émissivité intra-mensuelle et intra-annuelle, les erreurs de restitution liées à l'algorithme et également les erreurs de restitution liées à une mauvaise caractérisation de l'atmosphère dans les L2 d'EUMETSAT (vapeur d'eau, aérosols, nuages...).

Les valeurs sont ici supérieures aux erreurs de restitution (l'échelle est différente). Les zones que l'on avait identifiées comme étant potentiellement contaminées par les nuages (*i.e.*, l'Amazonie, la forêt équatoriale africaine et l'Indonésie) ressortent très bien. Cette importante dispersion de l'émissivité ( $\approx 5 \times 10^{-2}$ ) dans des régions très végétalisées et donc devant présenter une grande stabilité de l'émissivité, met en évidence la pollution nuageuse. En effet, une mauvaise identification des nuages entraîne un écart important entre les radiances simulées et celles mesurées par le satellite. L'algorithme de restitution peut alors donner des résultats incorrects. Les valeurs que prend l'émissivité dans ces régions sont tout de même proches des valeurs attendues (1 contre 0,98).

Certaines régions du globe présentent des statistiques de dispersion élevées sans être dans des zones particulièrement nuageuses. On peut noter par exemple l'ouest africain ou l'Australie. La dispersion de l'émissivité dans ces régions ne s'explique donc pas par des erreurs de caractérisation des nuages. La variation de l'humidité de la surface peut entraîner une modification pouvant atteindre 15% de son émissivité (Mira et al. 2007). Donc, si pendant le mois considéré, un orage éclate, la valeur de l'émissivité après le passage des nuages peut être modifiée. De plus, l'humidité du sol est différente d'une année sur l'autre. Une étude de l'interdépendance de l'humidité des sols avec l'émissivité est ainsi utile dans ce cas pour mieux expliquer ces dispersions. C'est l'objet de la prochaine section.

L'étude a été menée sur une seule longueur d'onde. Notre choix s'est porté sur  $1120\text{ cm}^{-1}$  car c'est le nombre d'onde auquel le plus de structures tant spectrales que spatiales sont identifiables<sup>2</sup>. Il faut noter que la dispersion mensuelle des émissivités est importante uniquement autour de la bande d'absorption des silicates entre  $1100$  et  $1200\text{ cm}^{-1}$  et dans la bande 3. Le

---

2. Les cartes correspondant aux autres longueurs d'onde ne sont pas montrées ici par souci de concision mais la base sera distribuée à la communauté.

reste du spectre reste assez peu variable. Les erreurs théoriques de restitution présentent les mêmes caractéristiques, avec des valeurs nettement inférieures. Les variations de l'émissivité dans la bande 3, ou autour des bandes d'absorption des silicates, peuvent s'expliquer par la différence entre des spectres d'émissivité de sol arides et végétalisés (ou humides). Afin de vérifier cela, nous allons conduire une étude temporelle des émissivités au cours des 5 années considérées.

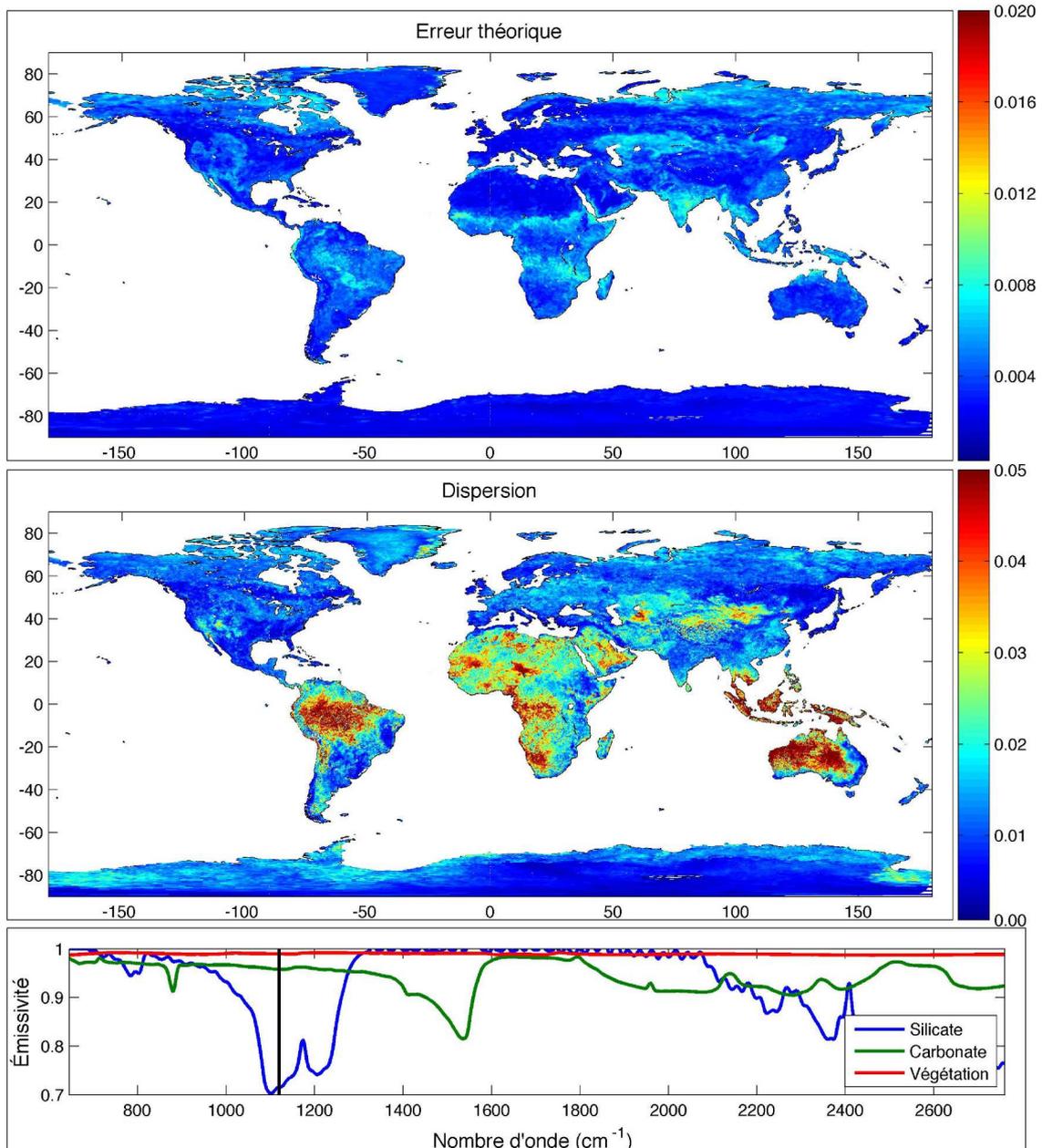


FIGURE 3.5 – Erreur théorique (en haut) et dispersion temporelle moyenne (en bas) de l'émissivité restituée à  $1120 \text{ cm}^{-1}$ , pour les mois de janvier de 2008 à 2012.

### 3.3.2 Variabilité temporelle

L'objectif de cette section est de vérifier l'hypothèse apparue dans la section précédente, à savoir la dépendance de l'émissivité à la variation du couvert végétal des sols. Cette mesure est complexe à effectuer. Nous avons donc décidé d'étudier les liens entre la pluviométrie et l'émissivité de la surface. En effet, les épisodes pluvieux augmentent l'humidité du sol et donc son couvert végétal et impactent son émissivité. Les zones très végétalisées ne seront pas modifiées par l'arrivée de la pluie. À l'inverse, les zones semi-arides seront grandement modifiées en cas de pluies. Il faut donc séparer cette étude suivant le type de surface considéré afin d'éviter de conclure sur de mauvaises raisons.

Pour cette étude, nous utiliserons la base de données de mesures de pluviométrie, issue de TRMM (Tropical Rainfall Measuring Mission) (Kummerow et al. 1998). Il s'agit d'une mission américano-japonaise visant à étudier les précipitations tropicales et leur impact sur le flux de chaleur latente. Elle inclut un radar dans le micro-onde à 18,83 GHz (PR, Precipitation Radar), un radiomètre micro-onde à 9 canaux (TMI, TRMM Microwave Image) et un radiomètre dans le visible et l'infrarouge à 5 canaux (VIRS, Visible and InfraRed Scanner). La base de données de mesures de pluviométrie que l'on considère ici ne provient pas uniquement de ces mesures. Sa construction prend en compte les mesures de nombreux autres instruments satellites : AMSR-E (Advanced Microwave Scanning Radiometer for Earth observing systems), SSM/I (Special Sensor Microwave Imager), SSMI/S (Special Sensor Microwave Imager/Sounder), AMSU (Advanced Microwave Sounding Unit), MHS (Microwave Humidity Sounder) et certains sondeurs infrarouges géostationnaires.

À ces données satellites sont ajoutées des mesures de pluie *in situ*. Ces mesures *in situ* proviennent de toutes sortes de capteurs. Il y a des capteurs cylindriques simples à graduation recueillant la pluie (mesures visuelles), des capteurs de poids (mesures grâce à des traces d'encre sur des papiers à la manière des sismographes), des détecteurs de gouttes électroniques... À l'inverse des mesures satellites dont les erreurs sont relativement constantes, les mesures *in situ* sont très hétérogènes. Leurs erreurs peuvent dépendre de l'instrument, du vent éventuel, de la position et des habitudes d'observation de l'utilisateur, et également des moyens de transmission de données des responsables locaux ou gouvernementaux. Toute la base de données de précipitations ne contient pas la même quantité d'informations *in situ*, compte tenu de leur répartition limitée sur le globe.

Au final, nous disposons d'un produit mensuel de précipitations, depuis 2007 jusqu'à nos jours, entre  $\pm 50^\circ$  de latitude. Ces mesures sont des moyennes mensuelles de millimètres de pluie par jour.

Nous allons donc comparer ce produit mensuel de pluviométrie à une moyenne mensuelle par années de l'émissivité sur trois surfaces différentes :

- Une première région dense en végétation tout au long de l'année, située dans le parc national de Taï en Côte d'Ivoire (sud de l'Afrique de l'ouest, coordonnées GPS :  $5,3750^\circ\text{N}-7,1548^\circ\text{E}$ );

- Une deuxième région plus variable, alternativement végétalisée ou aride suivant la saison, située à côté de Diandioumé à la frontière sud entre le Mali et la Mauritanie (coordonnées GPS : 15, 3750°N–9, 2075°E) ;
- Une troisième région très aride tout au long de l’année, située en plein désert mauritanien (coordonnées GPS : 20, 8750°N–9, 7695°E).

Ces trois régions ont été choisies pour leur localisation relativement proche les unes des autres mais leur caractéristiques différentes. Chacune de ces zones correspond au centre d’un pixel de la grille “equal-area” à  $0,25^\circ \times 0,25^\circ$ . On dispose donc, pour chacune de ces trois zones, de moyennes mensuelles de l’émissivité, pour chaque mois d’octobre 2007 à mars 2013. La première et la deuxième région sont relativement bien instrumentées pour la mesure de la pluviométrie, environ 52% des mesures proviennent de capteurs *in situ*. La troisième région (dans le désert) est bien moins instrumentée. La faible pluviométrie dans cette région permet de s’affranchir de ces données *in situ*, car il ne nous est pas nécessaire d’avoir une grande précision sur de si faibles valeurs, par rapport à ce que l’on souhaite montrer ici (il ne pleut quasiment pas dans cette région).

La Figure 3.6 représente la variation annuelle des moyennes mensuelles de l’émissivité à 845 (en bleu) et  $1120 \text{ cm}^{-1}$  (en vert), ainsi que de la pluviométrie (en rouge). Chaque cadre correspond à une des zones précédemment présentées : de haut en bas, la zone très végétée (le parc national Taï), la zone variable (frontière sud entre la Mauritanie et le Mali) et la zone aride (désert en Mauritanie). Nous avons choisi d’examiner l’émissivité à ces deux nombres d’onde car ils correspondent à un canal très stable en émissivité ( $845 \text{ cm}^{-1}$ ) et un autre qui varie fortement en fonction du couvert de végétation (émissivité élevée) et de la sécheresse du sol, entraînant la visibilité de roches silicates (émissivité faible).

Le suivi de la courbe de pluviométrie d’octobre 2007 à mars 2013 illustre bien le climat qui règne dans ces zones. Les mois de décembre à mai correspondent généralement à la saison sèche. Si seul le mois de janvier semble être plus sec que les autres dans la zone très végétalisée, cette saison sèche est bien visible sur la zone variable. Il ne pleut presque pas pendant toute une partie l’année. La diminution des réserves hydriques entraîne une désertification ponctuelle de certaines régions. La zone aride reste sèche toute l’année, on ne voit donc pas d’évolution de son émissivité. De même, la zone très végétée n’évolue pas beaucoup au cours de l’année, la quantité de pluie reste assez importante, on ne voit pas d’évolution nette de son émissivité. Dans la zone variable, la désertification du sol avec la saison sèche est apparente sur l’évolution de son émissivité. En effet, au long de la saison sèche, quelque soit l’année considérée, son émissivité à  $1120 \text{ cm}^{-1}$  diminue. Ceci signifie qu’il y a de plus en plus de silicates visibles sur le sol, il y a donc de moins en moins de végétation à la surface qui recouvre les roches. La stabilité du canal à  $845 \text{ cm}^{-1}$  nous permet de vérifier que cette variation de l’émissivité n’est pas simplement due au bruit.

Le mois de juin correspond au début de la saison des pluies. On peut remarquer (surtout sur la zone variable, au milieu) que cela correspond à un pic de pluie. Avec ces pluies,

l'activité végétale redémarre dans le sub-Sahel. La surface de la zone variable se recouvrira alors de végétation, même si cette dernière reste relativement clairsemée. On constate alors une augmentation de l'émissivité du sol, signe d'une végétation plus importante par rapport aux silicates en surface.

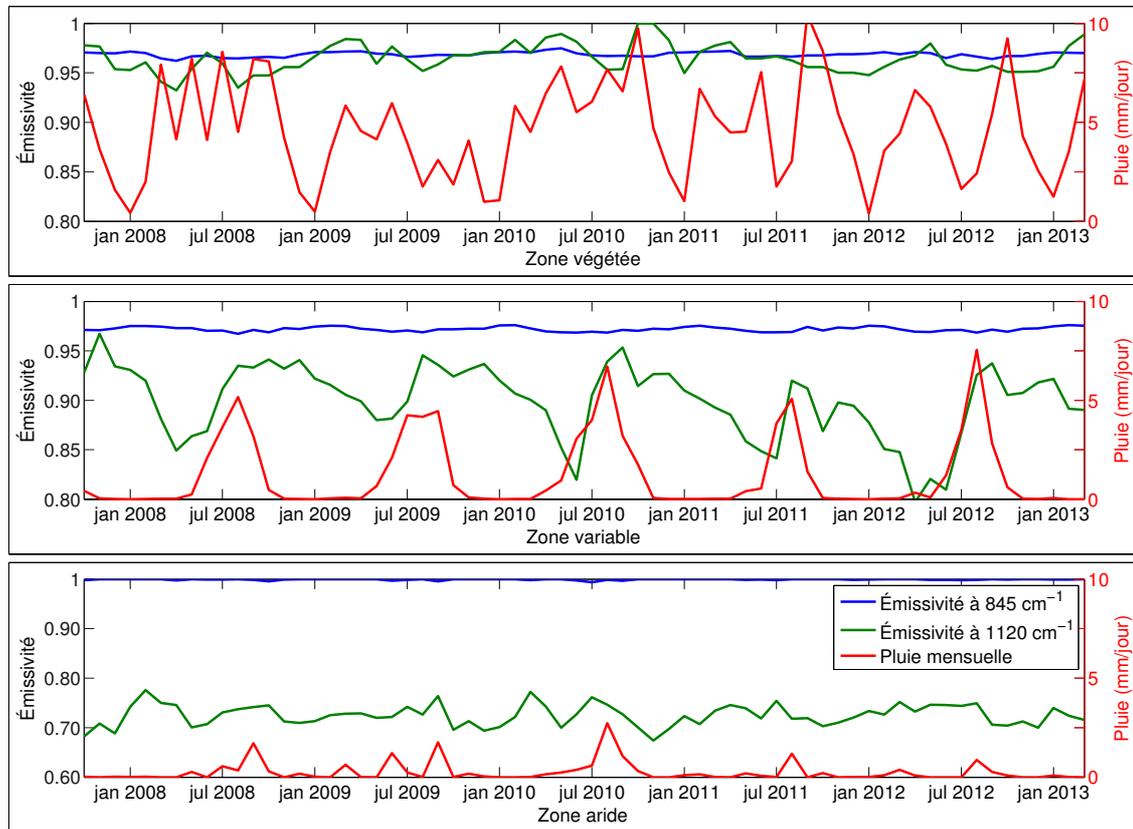


FIGURE 3.6 – Variation temporelle de l'émissivité mensuelle moyenne à 845 (en bleu) et 1120  $\text{cm}^{-1}$  (en vert), ainsi que de la pluviométrie mensuelle moyenne (en rouge), pour trois régions du globe : une à forte densité de végétation toute l'année (en haut) ; une variable suivant la saison (au milieu) ; et une aride (en bas).

La coïncidence temporelle entre les événements pluvieux et l'émissivité de la surface, vient valider la cohérence temporelle des restitutions d'émissivités. La dispersion de ces données, mise en avant précédemment, est liée à la quantité de végétation recouvrant les sols.

## 3.4 Comparaison avec des radiosondages au-dessus du dôme C

### 3.4.1 Données disponibles et objectif

Nous disposons, dans cette partie, de radiosondages effectués au-dessus du Dôme C dans l'Antarctique. Il s'agit de mesures des profils atmosphériques grâce à différents capteurs installés sur un ballon qui s'élève dans l'air. Ces sondages sont effectués de façon quotidienne à midi et couvrent toute l'année 2010. Ils contiennent la pression, la température et la vapeur d'eau sur tout le profil atmosphérique.

Nous disposons également d'une base de restitutions de profils de température et de vapeur d'eau au-dessus du Dôme C, à partir des mesures de l'instrument HAMSTRAD (H<sub>2</sub>O Antarctica Microwave Stratospheric and Tropospheric Radiometer). Il s'agit d'un radiomètre micro-onde qui possède deux canaux, un à 183 GHz (utilisé pour les restitutions de vapeur d'eau) et un à 60 GHz (utilisé pour les restitutions de température). Les profils restitués vont de 0 à 10 km d'altitude. Ces profils sont destinés à être comparés à des restitutions IASI, des radiosondages, des re-analyses... (Ricaud et al. 2012; Tremblin et al. 2011; Ricaud et al. 2010b,a).

Notre objectif dans cette partie est de tester la stabilité de notre algorithme au changement des données de première ébauche utilisées en entrée. Cela nous permettra également de voir si une meilleure qualité des premières ébauches peut affiner les restitutions, ou si l'information de première ébauche d'EUMETSAT est suffisante pour l'inversion, et ce dans un environnement particulièrement complexe : l'Antarctique. Pour cela nous allons effectuer des restitutions à l'aide de notre algorithme en modifiant les données utilisées en entrée. Nous utiliserons les L2 d'EUMETSAT comme précédemment, mais également les profils des radiosondages et les restitutions d'HAMSTRAD.

### 3.4.2 Méthode

Il a fallu au préalable colocaliser les sondages IASI avec les autres mesures. Les L2 d'EUMETSAT étant des restitutions directes des sondages IASI, les profils atmosphériques sont déjà parfaitement colocalisés. Les radiosondages et les mesures d'HAMSTRAD sont effectués au-dessus du Dôme C.

La Figure 3.7 représente tous les sondages IASI de l'année 2010 situés à moins de 18 km du Dôme C et qui sont considérés comme étant des situations de ciel clair. Sur cette figure, chaque sondage IASI est représenté par un point noir. Ces points noirs ont normalement un diamètre de 12 km, correspondant à champ de vue d'un sondage IASI (voir Section 1.3.1, page 29). Afin de pouvoir distinguer les différents points, nous avons réduits la représentation de ces points. Il faut cependant garder à l'esprit qu'ils recouvrent un disque de diamètre équivalent à la distance du centre (*i.e.*, le Dôme C) à la zone bleue sombre. Nous avons dé-

cidé de fixer à 18 km la limite de distance pour la colocalisation afin de conserver un nombre de points suffisant. Il reste alors 301 profils pour toute l'année 2010. Les autres carrés colorés indiquent qu'une restriction de la distance au Dôme C diminuerait drastiquement le nombre de situations restantes.

La colocalisation temporelle entre les sondages IASI et les restitutions HAMSTRAD est assez aisée puisque ces dernières sont disponibles toutes les 15 minutes. Les radiosondages sont disponibles de façon quotidienne à midi. Il faut donc garder à l'esprit qu'il peut y avoir un écart en temps important entre les radiosondages et les autres mesures.

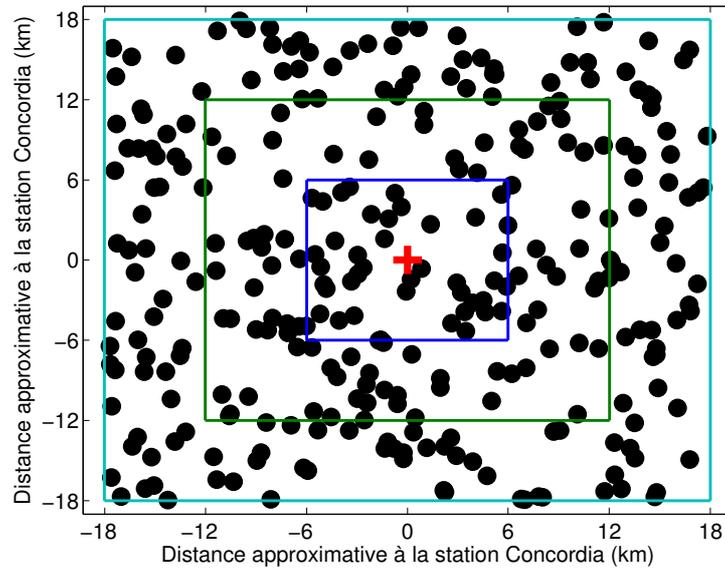


FIGURE 3.7 – Colocalisation des sondages IASI au-dessus du Dôme C (représenté par une croix rouge) pour toute l'année 2010. Une tolérance de 18 km est considérée (carré bleu clair), les autres carrés correspondent à des tolérances plus strictes en distance et donc moins de points. Chaque point noir représente un sondage IASI. L'orientation est vers le Nord.

### 3.4.3 Résultats

La Figure 3.8 représente en traits continus les températures de surface restituées en utilisant comme première ébauche les informations issues des L2 d'EUMETSAT (en bleu), des restitutions HAMSTRAD (en vert) et des radiosondages (en rouge). Pour information, les premières ébauches utilisées dans chacun des cas sont aussi tracées en pointillés avec le même code couleur.

On constate sur cette figure que les premières ébauches issues des L2 et des radiosondages sont assez proches, avec toutefois quelques différences. Les températures de surface restituées utilisant les L2 et les radiosondages sont, elles, quasi-identiques. Les zones où la courbe bleue n'est pas visible sont les zones où les deux courbes sont superposées. Même

lorsque les premières ébauches semblent différentes, les températures restituées sont similaires, ce qui veut dire que l'inversion arrive à corriger les erreurs de la première ébauche et est stable. La première ébauche en température de surface issue des restitutions HAMSTRAD est sensiblement différente des autres premières ébauches. Cet écart se traduit par des différences entre les températures de surface restituées<sup>3</sup>. Certaines restitutions semblent même avoir divergé et atteint des valeurs aberrantes. Ces erreurs sont également liées à la qualité du profil atmosphérique donné en première ébauche. Ainsi, une mauvaise caractérisation de l'atmosphère entraîne des erreurs importantes pour la restitution de la température de surface. Certains profils (les situations 250 à 290) ne sont même pas restitués par notre algorithme car les profils atmosphériques fournis par les restitutions HAMSTRAD n'appartiennent pas à l'intervalle d'existence de profils acceptés par RTTOV. Un filtrage des données d'HAMSTRAD semble donc nécessaire.

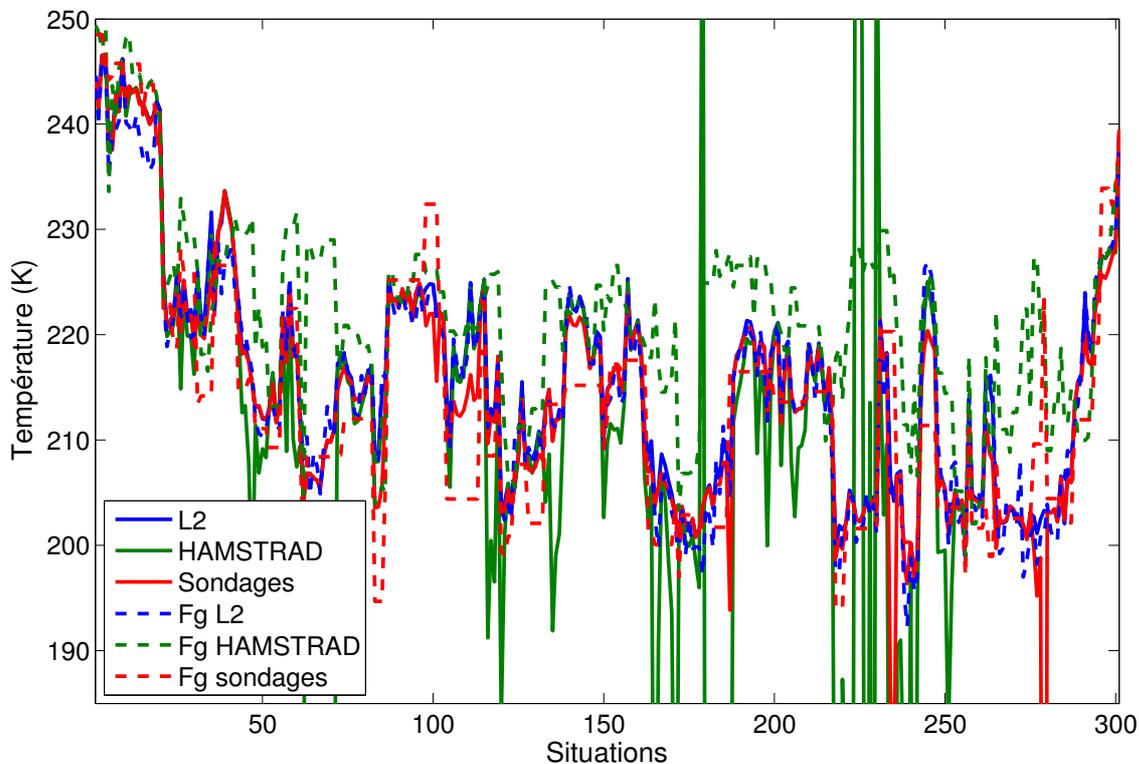


FIGURE 3.8 – Comparaison des températures de surface restituées au-dessus de l'Antarctique. Les courbes en pointillés correspondent aux premières ébauches utilisées dans la restitution, celles en lignes continues correspondent aux températures restituées. Les courbes bleues sont celles correspondant aux L2 d'EUMETSAT, celles en vert correspondent aux restitutions utilisant HAMSTRAD et celles en rouge correspondent aux radiosondages.

Le tableau ci-dessous montre les différents coefficients de corrélation temporelle entre les

3. Cet écart est en partie dû à la différence entre la température de surface et la température de l'air proche de la surface. En effet, HAMSTRAD ne fournit pas de température de surface à proprement parler, contrairement aux autres bases de données. La première ébauche en température de surface est alors erronée.

### 3.4. COMPARAISON AVEC DES RADIOSONDAGES AU-DESSUS DU DÔME C

6 bases de données de températures de surface (ces statistiques sont calculées uniquement sur les situations où la restitution utilisant HAMSTRAD a fonctionné, soit 259 situations sur les 301). On retrouve ici ce que l'on expliquait précédemment. Les premières ébauches en température de surface issues des L2 et des radio-sondages sont très proches (corrélation de 0,89) et les températures restituées à partir de ces données sont quasiment identiques (corrélation de 0,96). Les L2 d'EUMETSAT semblent donc être suffisantes pour notre algorithme, au moins pour cette localisation spatiale. Les données issues des restitutions HAMSTRAD sont, elles, un peu plus éloignées, avec des corrélations de 0,72 avec les L2 et 0,8 avec les données de radiosondages. Les températures restituées à partir des données HAMSTRAD sont, elles, vraiment éloignées des autres températures, avec des corrélations presque toutes inférieures à 0,6.

$T_s$	Fg L2	Fg HAMSTRAD	Fg Sondages	L2	HAMSTRAD	Sondages
Fg L2	1	0,72	0,89	0,99	0,58	0,96
Fg HAMSTRAD	0,72	1	0,80	0,74	0,26	0,73
Fg Sondages	0,89	0,80	1	0,91	0,49	0,89
L2	0,99	0,74	0,91	1	0,60	0,97
HAMSTRAD	0,58	0,26	0,49	0,60	1	0,62
Sondages	0,96	0,73	0,89	0,97	0,62	1

En effectuant les inversions, nous sommes passés d'une corrélation de 0,96 à 0,97. Ceci peut paraître peu, mais cela veut dire que les L2 d'EUMETSAT ne représentaient qu'une partie de la variabilité de la température de surface et que notre algorithme a été capable d'améliorer un produit déjà de bonne qualité.

Le fait que les restitutions, utilisant les températures d'HAMSTRAD, soient encore plus éloignées des autres que les températures d'HAMSTRAD elles-mêmes montre que les erreurs de restitution de la température de surface, en utilisant HAMSTRAD, ne sont surement pas entièrement dues à une mauvaise caractérisation de la température de surface utilisée en première ébauche. Les profils atmosphériques utilisés viennent augmenter l'erreur commise dans le calcul de transfert radiatif. L'écart trop élevé entre les radiances simulées et réelles entraîne de fortes erreurs dans l'estimation des paramètres de surface.

Pour mettre en évidence ces écarts dans la caractérisation de l'atmosphère, la Figure 3.9 montre la RMS de l'erreur entre les profils de température et de vapeur d'eau issus des L2 d'EUMETSAT, des restitutions d'HAMSTRAD et des radiosondages. Sur la partie gauche de la figure, nous avons ajouté au profil de température les écarts en température de surface issues des trois bases de données (celles utilisées comme premières ébauches dans les restitutions). Les courbes bleues correspondent aux RMS des erreurs entre les profils de température (et la température de surface) et d'humidité relative provenant des L2 d'EU-

METSAT et des radiosondages. Les courbes vertes correspondent aux mêmes statistiques, mais entre les profils des L2 d'EUMETSAT et ceux issus d'HAMSTRAD. Les courbes rouges correspondent aux RMS d'erreur entre les profils d'HAMSTRAD et les radiosondages. Les différents profils sont disponibles sur des niveaux de pression différents, voire même variables d'une situation à l'autre pour les radiosondages ou les profils issus d'HAMSTRAD. Nous avons projeté tous les profils sur les 90 niveaux auxquels les mesures de ballon sonde sont effectuées. C'est pourquoi l'échelle en ordonnées n'est pas présentée en pression, celle-ci pouvant être variable.

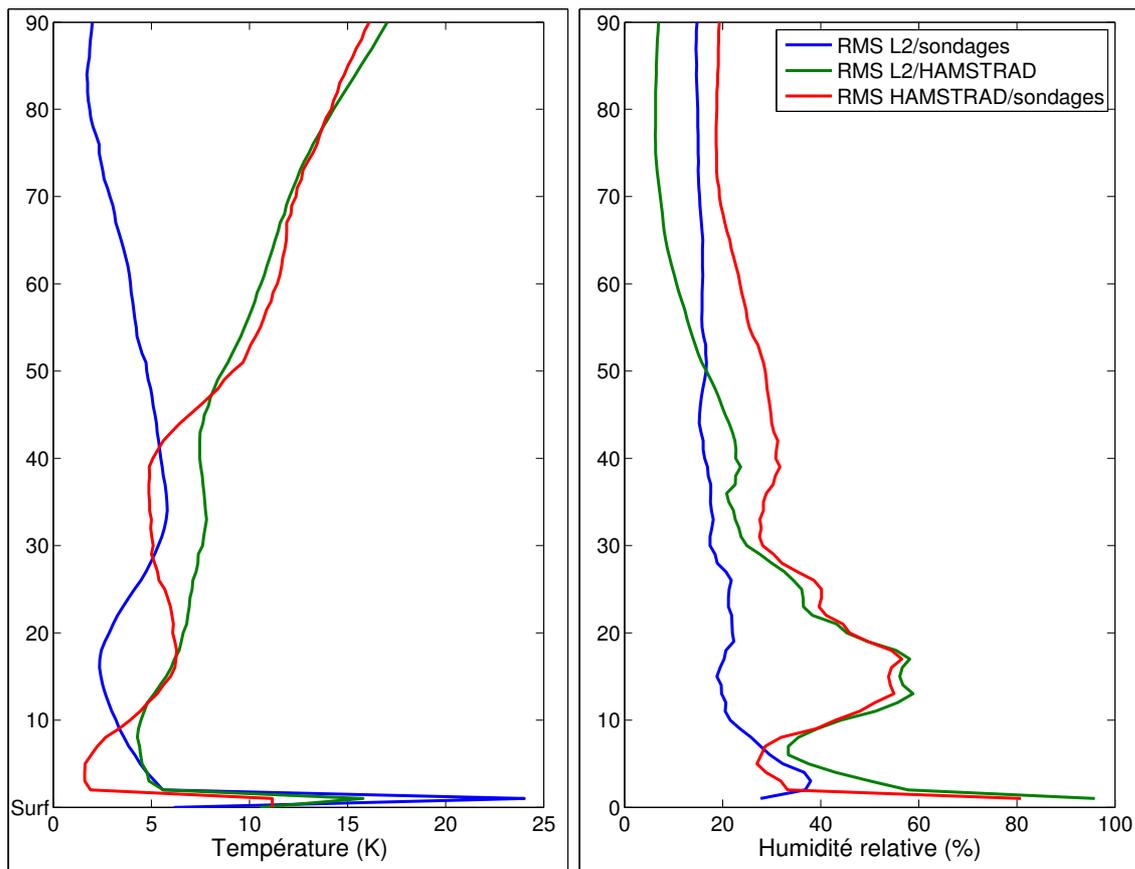


FIGURE 3.9 – Comparaison des profils atmosphériques des L2 d'EUMETSAT, des restitutions d'HAMSTRAD et des radiosondages. L'échelle verticale correspond à des niveaux de mesures des radi-sondages sur lesquels sont projetés les autres profils.

On remarque ici que les profils de température issus d'HAMSTRAD semblent être sensiblement différents des profils des radiosondages et des L2. La courbe bleue présente, en effet, des statistiques inférieures aux autres courbes, notamment dans le haut de l'atmosphère. Si la température dans les basses couches de l'atmosphère issue des L2 d'EUMETSAT présente des statistiques éloignées de celle issue des radiosondages, ce n'est plus le cas au niveau de la température de surface.

Des observations similaires peuvent être faites sur les profils d'humidité relative. La plus basse couche de l'atmosphère des profils d'humidité relative issue d'HAMSTRAD présente des écarts importants avec les autres bases de données (courbes vertes et rouges). On retrouve ces écarts importants vers le milieu de l'atmosphère.

Ces mauvaises caractérisations, tant de l'atmosphère que de la température de surface, entraînent donc des erreurs dans la restitution de la température de surface. La Figure 3.10 présente les courbes d'émissivité restituée sur toute l'année (301 situations). Du fait de la faible variabilité des profils d'émissivité de la neige et de la glace, seule la valeur moyenne de l'émissivité est représentée ici. La courbe noire correspond à la première ébauche en émissivité. Les courbes bleue, verte et rouge correspondent respectivement aux émissivités restituées en utilisant les L2 d'EUMETSAT, les données d'HAMSTRAD et les radiosondages.

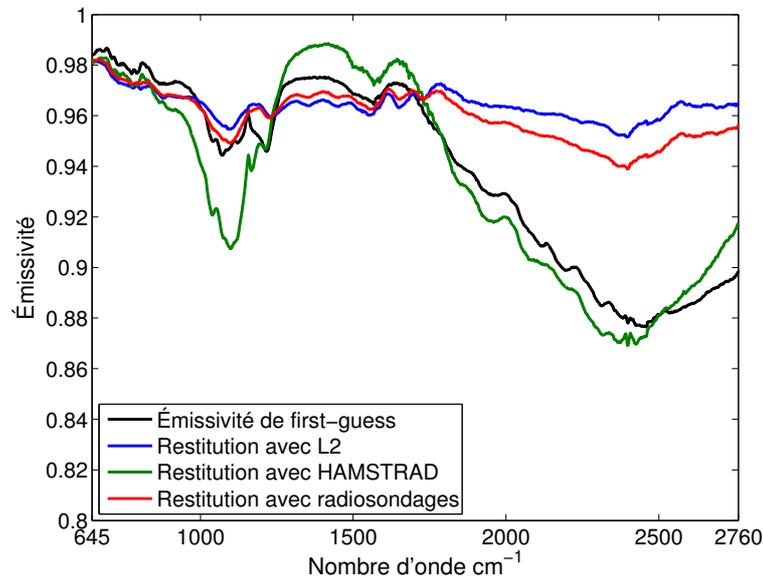


FIGURE 3.10 – Émissivités moyennes de première ébauche (courbe noire) et restituées au-dessus du Dôme C en utilisant les L2 d'EUMETSAT (courbe bleue), les données d'HAMSTRAD (courbe verte) et les radiosondages (courbe rouge).

L'émissivité de première ébauche présente une importante décroissance dans la bande 3. Cette décroissance de l'émissivité de glace ou de neige dans la bande 3 est rarement observée en laboratoire (Snyder et al. 1998). Il est donc intéressant de noter qu'on ne la retrouve pas avec les profils restitués utilisant les radiosondages ou les L2 d'EUMETSAT. Ces deux profils sont encore une fois très semblables et donnent une émissivité relativement constante autour de 0,96. Notre algorithme a donc réussi à corriger la première ébauche en émissivité. L'émissivité restituée utilisant les profils d'HAMSTRAD présente, quant à elle, cette décroissance dans la bande 3. On peut noter également un pic autour de 1100  $\text{cm}^{-1}$ . Ceci est dû à l'instabilité de la restitution qui entraîne des modifications importantes des compo-

santes de l'émissivité. La variabilité globale (spatialement) de l'émissivité à  $1100\text{ cm}^{-1}$  fait qu'elle est très marquée sur les composantes. Les instabilités sur les composantes entraînent cette surreprésentation de la bande à  $1100\text{ cm}^{-1}$ , correspondant à des silicates et qui n'a donc pas de sens physique au-dessus de la glace.

Ainsi, l'algorithme que nous avons développé ici est relativement stable. Les restitutions utilisant les L2 d'EUMETSAT et les radiosondages mènent au même résultat. Cependant, de trop grandes incertitudes sur la caractérisation initiale, tant de la surface que de l'atmosphère, peuvent entraîner des instabilités dans l'inversion et des résultats erronés. Il faut donc garder à l'esprit que l'algorithme que nous avons développé nécessite une première ébauche raisonnable. Utiliser diverses premières ébauches de bonne qualité ne modifie pas les inversions de notre algorithme.

### 3.5 Perspectives

Nous avons validé notre algorithme de différentes façons (comparaisons de températures de brillance simulées et réelles, comparaisons globales de température de surface, stabilité à la première ébauche, variation saisonnière de l'émissivité). Il serait intéressant de continuer à comparer notre base d'émissivités à d'autres bases de données. Pour cela, il est nécessaire de mener une large étude d'intercomparaisons des produits d'émissivité. On peut citer à ce jour 6 équipes travaillant sur cette thématique : la NASA, le groupe ARA au LMD, EUMETSAT, le MetOffice (équivalent anglais de Météo-France), l'Université de Basilicate et nous mêmes (Zhou et al. 2011; Capelle et al. 2012; Schlüssel et al. 2005; August et al. 2012; Thelen et al. 2009; Masiello and Serio 2013; Paul et al. 2012). Nous avons décidé de joindre nos efforts, sous la direction d'EUMETSAT, afin de comparer les différents produits disponibles.

La première partie de l'étude consistera en des comparaisons sur des sites localisés. 10 sites ont été sélectionnés, majoritairement au-dessus de zones arides, puisque ce sont les zones présentant le plus de variabilité spectrale. Nous recommandons, à plus long terme, de mener une étude à l'échelle planétaire de ces émissivités hyperspectrales. Il serait intéressant de se pencher sur une interprétation plus physique des émissivités de surface (corrélations avec l'humidité du sol, sa composition...). Une validation des différentes émissivités dans l'espace des températures de brillance pourrait être menée, en prenant en compte les différentes spécificités d'erreur des bases.

Les délais impartis à de telles collaborations internationales dépassant malheureusement le cadre du travail de thèse mené ici, les résultats de cette étude ne peuvent être fournis. La diversité des équipes travaillant sur ce projet rend complexe le partage d'informations à des formats différents. Un long travail de formatage et de coïncidence des données est nécessaire.

### 3.6 Conclusion

Nous avons donc créé dans ce chapitre un outil capable de restituer simultanément la température de surface et son émissivité hyperspectrale dans l'infrarouge. Cette inversion est effectuée à partir des observations de IASI et d'une bonne première ébauche sur l'état de l'atmosphère et de la surface. La comparaison des résultats obtenus avec de nombreuses bases différentes, la variabilité spatiale, spectrale et temporelle des produits restitués, l'application de notre schéma à des données indépendantes ont permis de valider les caractéristiques de surface que nous avons restituées. La connaissance précise de la température de surface et de son émissivité à tous les canaux de IASI, associée aux mêmes informations issues des instruments micro-ondes, devrait nous permettre (dans le chapitre suivant) de mieux exploiter les observations satellites et donc de restituer de façon plus précise les profils atmosphériques au-dessus des continents.



---

# SYNERGIE : PRINCIPES GÉNÉ- RAUX



## Sommaire

<b>4.1 Principe</b> . . . . .	<b>108</b>
<b>4.2 Exemple sans synergie</b> . . . . .	<b>111</b>
<b>4.3 Synergie additive</b> . . . . .	<b>111</b>
<b>4.4 Synergie indirecte</b> . . . . .	<b>114</b>
<b>4.5 Synergie de débruitage</b> . . . . .	<b>117</b>
<b>4.6 Généralisation</b> . . . . .	<b>118</b>
<b>4.7 Conclusion</b> . . . . .	<b>119</b>

---

L'autre volet de notre travail, au delà de la caractérisation des surfaces, est d'étudier l'utilisation combinée de différentes informations à différentes longueurs d'ondes. La combinaison de l'infrarouge et du micro-onde, en particulier, permet-elle de mieux restituer les variables atmosphériques? Historiquement, des méthodes ont été développées pour utiliser une seule gamme de longueur d'onde pour restituer une variable atmosphérique donnée. Pour supprimer les ambiguïtés liées à la contribution des autres paramètres, ces algorithmes exploitent, par exemple, la complémentarité des bandes de fréquences proches, mesurées par les mêmes instruments, ou utilisent des informations auxiliaires sur l'état de l'atmosphère, provenant de sources indépendantes. Cette première approche, utilisant des informations de sources différentes successivement ou hiérarchiquement, est utile car elle combine plusieurs informations pour la restitution. Toutefois, elle ne semble pas optimale car elle n'exploite pas aux mieux les synergies potentielles. La fusion d'observations simultanées de l'état de l'atmosphère dans des gammes de longueurs d'onde différentes pourrait aider à séparer les différentes contributions en vue d'obtenir de meilleures estimations d'une variable donnée.

Le terme synergie vient du grec *syn* signifiant "avec" et *ergazomai* signifiant "travailler". La notion de travail conjoint est mise en valeur dans ce mot. Il apporte également une notion implicite d'amélioration liée au travail collectif. La synergie est le principe selon lequel, en combinant deux ou plus agents différents, on obtient un résultat plus précis qu'en utilisant les résultats de ces différents agents séparément.

Dans le cadre de la télédétection, qui nous intéresse ici, la synergie correspond au principe selon lequel la restitution d'un profil atmosphérique, en utilisant plusieurs sources d'informations complémentaires, sera plus précise que la meilleure restitution indépendante. Cette idée semble facile à accepter, cependant ce principe suggère que lorsque l'on dispose d'une information de qualité sur une variable, on peut améliorer sa précision en y adjoignant une information de moins bonne qualité.

Il est possible de mettre en avant trois types distincts de synergie (Aires 2011) :

- La synergie additive ;
- La synergie indirecte ;
- La synergie de débruitage.

Ce chapitre détaille ces différentes synergies et les met en évidence grâce à des exemples simples. La synergie est, dans cette thèse, appliquée aux mesures atmosphériques par satellite. Le principe de synergie est valable pour de nombreuses autres applications, où plusieurs observations sont sensibles à une même variable. C'est pourquoi nous choisissons ici de présenter la synergie sous une forme très générale.

## 4.1 Principe

La synergie concerne un algorithme qui utilise de façon simultanée ou hiérarchique les observations d'au moins deux différentes observations pour obtenir un résultat plus précis que les restitutions obtenues avec chacune des observations séparément. Dans le cadre de la télédétection, nous l'appliquerons en utilisant des algorithmes de restitution de variables atmosphériques. Les observations que nous utiliserons seront les mesures satellites dans les différentes gammes de fréquence (micro-onde et infrarouge). Nous définissons un facteur de synergie  $F_{syn}$  d'un algorithme utilisant  $n$  sources d'informations :  $(x_1, x_2 \dots x_n)$  (chacune pouvant être multiple, comme les différents canaux d'un capteur), comme :

$$F_{syn} = \frac{\min_{i=1,2,\dots,n} (E_i)}{E_{1,2,\dots,n}}$$

où  $E_i$  est l'erreur quadratique moyenne de restitution en utilisant seulement l'observation  $x_i$  et  $E_{1,2,\dots,n}$  est l'erreur de restitution en utilisant toutes les informations combinées. Plus  $F_{syn}$  est élevé plus la synergie est efficace. Il permet de mesurer directement le gain en terme de pourcentage d'erreur dû à l'utilisation de la synergie. Si ce facteur est supérieur à 100%, l'erreur de restitution, en utilisant toutes les observations, est inférieure à la meilleure des restitutions indépendantes. L'algorithme est alors à même d'exploiter la synergie entre les observations. Inversement si le facteur est inférieur à 100%, l'inversion indépendante est plus efficace que l'inversion en utilisant les différentes observations simultanément, l'algorithme n'est donc pas synergique. Ce facteur permet de quantifier la synergie d'un algorithme.

Afin de mettre en évidence les différentes synergies qui peuvent exister il est plus simple

d'utiliser un schéma de restitution donné. Nous utiliserons ici un schéma avec seulement deux variables à restituer et deux observations d'entrée afin de faciliter les calculs et la démonstration. Cette méthode s'appliquerait de la même façon si il y avait plus d'observations ou de variables à restituer.

Soient  $f_1$  et  $f_2$  deux variables (dans le cadre de la télédétection atmosphériques, ce sont deux variables atmosphériques comme la vapeur d'eau et la température),  $B_1$  et  $B_2$  deux observations différentes (une température de brillance dans le micro-onde et l'autre dans l'infrarouge.). On modélise leur relation sous une forme simple linéaire :

$$\begin{pmatrix} B_1 \\ B_2 \end{pmatrix} = A \cdot \begin{pmatrix} f_1 \\ f_2 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix}$$

où  $A = \begin{pmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{pmatrix}$  représente, dans le cadre du sondage atmosphérique une linéarisation de l'équation de transfert radiatif.  $\varepsilon_1$  et  $\varepsilon_2$  correspondent aux erreurs instrumentales des observations  $B_1$  et  $B_2$ . On suppose que ces deux erreurs suivent une distribution gaussienne sans biais. On représente la matrice de covariance d'erreur par :  $S_\varepsilon = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ . Les valeurs 1 et 2 pour la variance des deux erreurs sont choisies arbitrairement pour l'exemple traité dans ce chapitre. Elles correspondent à un écart-type d'erreur de 1 et  $\sqrt{2}$ , ce qui est cohérent avec l'erreur instrumentale des instruments satellites.

On suppose également que l'on dispose d'une première estimation (appelée première ébauche ou "first-guess") de  $f_1$  et  $f_2$  (issue, par exemple, des NWP (voir Section 1.5.2, page 34)) et que leurs erreurs sont caractérisées par une distribution gaussienne de biais nul et de matrice de covariance  $S_f$ . La solution de ce problème est alors (Annexe B.2.1, page 213) :

$$f = f_g + \left( A^t \cdot S_\varepsilon^{-1} \cdot A + S_f^{-1} \right)^{-1} \cdot A^t \cdot S_\varepsilon^{-1} \cdot (F_\varepsilon - F_g)$$

$F_\varepsilon$  correspond aux observations  $B_1$  et  $B_2$ , c'est-à-dire aux mesures satellites.  $F_g = A \cdot f_g$  est la simulation de transfert radiatif (voir Section 1.4, page 33) sur la première ébauche  $f_g$  des variables  $f_1$  et  $f_2$ . L'erreur commise sur la restitution des variables est caractérisée par (directement d'après l'Annexe B.2.1, page 213) :

$$Q = \left( A^t \cdot S_\varepsilon^{-1} \cdot A + S_f^{-1} \right)^{-1} \quad (4.1)$$

On obtient ainsi la matrice de covariance d'erreur  $Q$  sur la restitution des variables  $f_1$  et  $f_2$ . Nous considérons, dans ce cas, que l'erreur de restitution n'est pas biaisée. On peut donc obtenir l'erreur quadratique moyenne de restitution directement en prenant la racine carrée des termes diagonaux de la matrice de covariance d'erreur.

Cette méthode peut s'appliquer de la même façon en utilisant qu'une seule des deux

observations  $B_1$  ou  $B_2$ . On peut aisément comparer les erreurs de restitution lorsqu'une ou plusieurs observations sont utilisées, calculer  $F_{syn}$  et conclure sur la présence ou non de synergie. Différents types de synergie sont présentés ici. Afin de simplifier la méthode, ils sont clairement dissociés. L'objectif est de prouver l'existence de la synergie de façon mathématique et d'illustrer ses différentes formes afin de justifier l'utilisation qu'il en est faite dans les chapitres suivant.

La Figure 4.1 présente le schéma de restitution et illustre les différentes synergies. Comme présenté dans les sections suivantes, chaque type de synergie s'appuie sur une certaine relation entre les variables ou les observations.

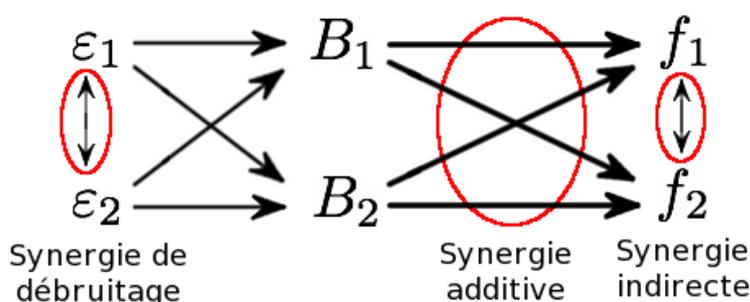


FIGURE 4.1 – Schéma de restitution.  $B_1$  et  $B_2$  sont les observations connues,  $f_1$  et  $f_2$  sont les variables à restituer,  $\varepsilon_1$  et  $\varepsilon_2$  sont les bruits instrumentaux sur les observations. Les différents types de synergie sont détaillés dans le texte.

La synergie additive est la plus simple : il s'agit tout simplement d'avoir plusieurs sources d'informations sur une ou plusieurs variables à restituer. La synergie indirecte repose sur une relation entre les deux variables à restituer : une meilleure connaissance d'une variable entraîne une meilleure connaissance de l'autre de part les corrélations entre ces variables. La synergie de débruitage est plus complexe : il s'agit d'une relation entre les bruits des deux observations. L'ajout d'une observation permet alors de mieux connaître le bruit de l'autre et ainsi améliorer la restitution.

En partant d'une situation de départ, pour chaque type de synergie, on décrit les différentes matrices de sensibilité. Les résultats théoriques obtenus seront comparés (à l'aide de la matrice  $Q$ ) avec les premiers résultats obtenus dans un cadre sans synergie. Nous nous intéresserons particulièrement à la restitution de  $f_1$ , les résultats pour  $f_2$  étant similaires.

Afin de mentionner plus simplement les différents termes des matrices, ils seront chacun désigné par la lettre minuscule de la matrice correspondante indexée de la position du terme dans la matrice comme par exemple :  $A = \begin{pmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{pmatrix}$ .

## 4.2 Exemple sans synergie

Dans cet exemple, nous définissons la matrice  $A$  de l'équation (4.1) par :  $\begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}$ . Ainsi, chacune des observations  $B_1$  et  $B_2$  donne respectivement des informations uniquement sur la variable  $f_1$  ou  $f_2$ . On garde la matrice de covariance d'erreur sur les observations  $S_\varepsilon = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$  comme définie précédemment. On définit la matrice de covariance d'erreur sur la première ébauche par  $S_f = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$ . On a ainsi deux variables  $f_1$  et  $f_2$  qui présentent respectivement un bruit blanc gaussien de variance 2 et 3. On peut alors, d'après l'équation (4.1), calculer la matrice de covariance d'erreur de restitution des variables  $f_1$  et  $f_2$  :

$$Q = \left( \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}^t \cdot \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix} + \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}^{-1} \right)^{-1} \approx \begin{pmatrix} 0,76 & 0 \\ 0 & 1,73 \end{pmatrix}$$

Si les deux observations et les deux variables sont décorrélées, et que leur matrice de covariance d'erreur est diagonale, alors il n'y a pas de synergie mise en jeu dans la restitution de  $f_1$  et  $f_2$ . Dans ce cas, l'erreur de restitution sur  $f_1$  est de  $\sqrt{0,76} = 0,87$  et de  $\sqrt{1,73} = 1,32$  sur  $f_2$ . Les cas, qui sont présentés dans les sections suivantes, seront comparés à ce cas simple sans synergie. En ajoutant des termes non-diagonaux dans les différentes matrices de ce problème, on pourra présenter l'apport d'un type de synergie particulier.

## 4.3 Synergie additive

La synergie additive correspond au cas où deux observations contiennent de l'information sur plusieurs variables. Prenons un cas simple avec une seule variable à restituer,  $f_1$ . Nous avons vu précédemment que l'erreur de restitution de  $f_1$  en utilisant seulement  $B_1$  est de 0,87. Considérons que  $B_2$  est également sensible à  $f_1$ , on pose  $A = \begin{pmatrix} 0,9 \\ 0,7 \end{pmatrix}$ . En gardant les mêmes caractéristiques d'erreur que précédemment, c'est-à-dire  $S_\varepsilon = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$  et  $S_f = 2$ , on a :

$$Q = \left( \begin{pmatrix} 0,9 \\ 0,7 \end{pmatrix}^t \cdot \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 0,9 \\ 0,7 \end{pmatrix} + (2)^{-1} \right)^{-1} \approx (0,64)$$

L'erreur de restitution de  $f_1$  est alors de  $\sqrt{0,64} = 0,8$ , ce qui est inférieur à ce que l'on avait précédemment. On peut alors calculer le facteur de synergie  $F_{syn} = 109\%$ . L'utilisation de  $B_2$  dans la restitution améliore donc la restitution de 9%, ce qui est intéressant, la synergie

additive est bien présente ici.

Cette synergie est l'équivalent de la loi des grands nombres, plus il y a d'observations disponibles, plus les estimateurs convergeront vers la bonne solution.

Afin de mettre en valeur cette synergie, une étude de l'influence de la sensibilité de  $B_2$  à  $f_1$  (notée précédemment  $a_{2,1}$  dans la matrice 1) sur la restitution de  $f_1$  est présentée sur la Figure 4.2.

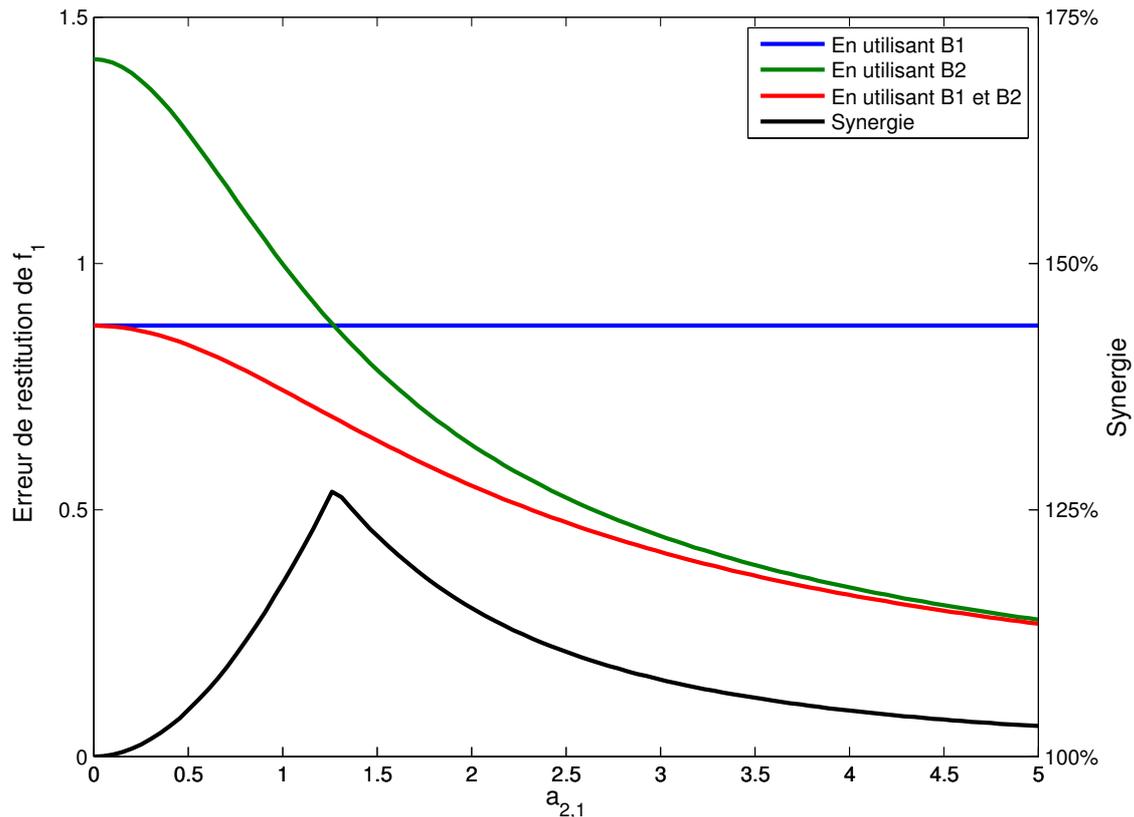


FIGURE 4.2 – Étude de la variation de la synergie additive en fonction du jacobien de  $B_2$  par rapport à la variable  $f_1$ . Les courbes représentent l'erreur de restitution en utilisant seulement  $B_1$  (bleu), seulement  $B_2$  (vert) et les deux conjointement (rouge). La synergie de la restitution est tracée en noir par rapport à un axe des ordonnées différent représenté à droite.

L'erreur de restitution de  $f_1$  est ici représentée en fonction de  $a_{2,1}$ . Les différentes courbes de couleur correspondent aux schémas utilisant différentes observations : la courbe bleue correspond à une restitution en utilisant uniquement  $B_1$ , la courbe verte en utilisant seulement  $B_2$  et la courbe rouge correspond à une restitution utilisant les deux observations.

On constate que l'erreur de restitution en utilisant seulement  $B_1$  est fixe. En effet, cette dernière ne dépend pas de  $a_{2,1}$ , elle reste constante et égale à  $\sqrt{0,76} = 0,87$ . L'erreur de restitution en utilisant uniquement  $B_2$  est une fonction décroissante de la sensibilité de  $f_1$  à  $B_2$ . Ce qui est normal, car plus le jacobien de  $B_2$ , par rapport à  $f_1$ , est élevé, plus  $B_2$

contient une information sur  $f_1$  et plus l'erreur de restitution est faible.

Ce qu'il faut noter ici, c'est que la courbe rouge est systématiquement située au-dessous des deux autres courbes colorées. Cela signifie que l'erreur de restitution de  $f_1$ , en utilisant conjointement  $B_1$  et  $B_2$ , est toujours plus faible que l'erreur de restitution en utilisant l'une ou l'autre des observations de façon indépendante. Ceci est également mis en exergue par la courbe noire. En effet, le facteur de synergie est toujours au-dessus de 100%, ce qui signifie que la restitution simultanée est meilleure que la meilleure restitution utilisant les deux observations séparément.

Lorsque  $a_{2,1}$  est faible (à gauche sur le graphique), l'apport de  $B_2$  à la restitution est faible, la synergie l'est donc également. Plus  $a_{2,1}$  augmente, plus l'apport de  $B_2$  est important. À l'autre extrémité du graphique,  $a_{2,1}$  est si élevé que  $B_2$  est nettement plus informatif que  $B_1$  et l'apport de  $B_1$  devient de plus en plus faible. Ceci explique la forme de la courbe de synergie, qui est faible au début, présente un pic quand les deux observations sont autant informatives, puis diminue.

La synergie additive est donc toujours présente lorsque plusieurs observations sont sensibles à une même variable.

L'impact du bruit instrumental de l'observation  $B_2$  sur la restitution simultanée est également analysée pour la synergie additive. En effet, on pourrait être amené à penser que rajouter une observation fortement bruitée n'apporte rien à la restitution. La Figure 4.3 représente l'erreur de restitution de  $f_1$  en fonction de l'erreur instrumentale sur  $B_2$ . Comme sur la figure précédente, les erreurs de restitution de  $f_1$  en utilisant seulement  $B_1$ , seulement  $B_2$  ou les deux observations simultanément sont tracées.

Ici encore, l'erreur de restitution de  $f_1$  en utilisant seulement  $B_1$  ne dépend pas du bruit sur  $B_2$ , il est donc logique de retrouver une fonction constante. L'erreur de restitution en utilisant seulement  $B_2$  est une fonction croissante du bruit instrumental sur  $B_2$ . Plus  $B_2$  est bruitée, plus il est dur d'en extraire de l'information. La courbe rouge quant à elle reste sous les deux autres quel que soit le bruit instrumental sur  $B_2$ . Ainsi l'erreur de restitution en utilisant simultanément les deux observations  $B_1$  et  $B_2$  est toujours plus faible que l'erreur commise en n'utilisant qu'une des deux observations. La courbe noire, qui représente le facteur de synergie de la restitution, montre qu'il est préférable d'utiliser les deux observations simultanément. La synergie est toujours supérieure à 100%, même si  $B_2$  est fortement bruitée. La synergie additive est encore une fois bien présente.

Cette synergie est très simple à comprendre. Puisqu'il s'agit de deux observations qui apportent chacune de l'information, il faut prendre en compte les deux pour avoir le plus possible d'informations. Le cas présenté ici est un cas où tous les paramètres de l'équation sont maîtrisés (le modèle direct  $A$ , les différents bruits...). Dans la pratique, ce n'est pas toujours le cas et il se peut qu'en fusionnant de l'information on dégrade la restitution. Dans ce cas, l'algorithme utilisé n'est pas adapté à l'utilisation de la synergie et il faut revoir la modélisation du problème.

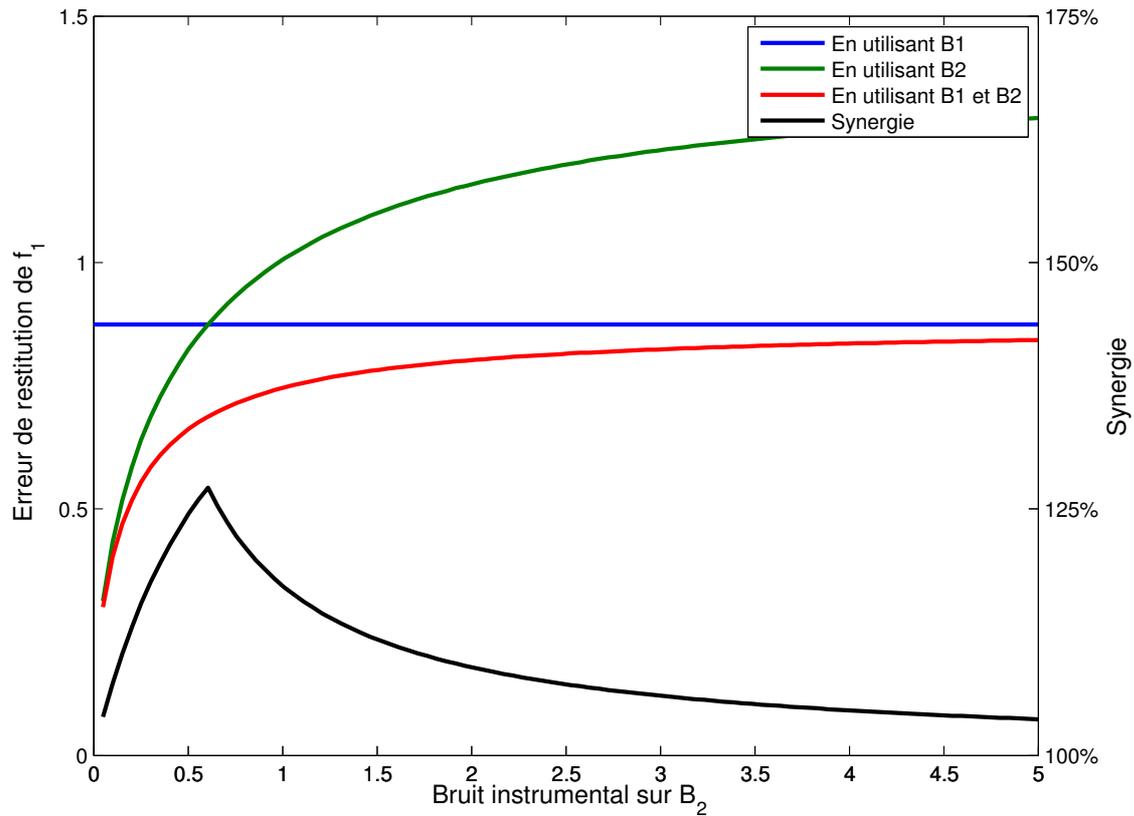


FIGURE 4.3 – Étude de la variation de la synergie additive en fonction du bruit instrumental sur l’observation  $B_2$ . Les courbes représentent l’erreur de restitution en utilisant seulement  $B_1$  (bleu), seulement  $B_2$  (vert) et les deux conjointement (rouge). La synergie de la restitution est tracée en noir par rapport à un axe des ordonnées différent représenté à droite.

Les types de synergie présentés par la suite sont plus complexes.

#### 4.4 Synergie indirecte

Dans cet exemple, la matrice  $A$  de l’équation (4.1) vaut  $\begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}$ . C’est-à-dire que chaque observation est sensible uniquement à une variable. La matrice  $S_\varepsilon$  vaut toujours  $\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ . La matrice qui est modifiée est la matrice de covariance d’erreur de la première ébauche  $S_f$ , nous la définissons ici par  $\begin{pmatrix} 2 & 1,5 \\ 1,5 & 3 \end{pmatrix}$ . L’erreur sur les premières ébauches est donc corrélée. C’est cette corrélation qui induit une synergie indirecte.

Nous pouvons alors calculer la matrice de covariance d'erreur de restitution  $Q$  :

$$Q = \left( \left( \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}^t \cdot \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix} + \begin{pmatrix} 2 & 1,5 \\ 1,5 & 3 \end{pmatrix}^{-1} \right)^{-1}$$

$$Q \approx \begin{pmatrix} 0,71 & 0,37 \\ 0,37 & 1,47 \end{pmatrix}$$

En comparant aux résultats obtenus à la Section 4.2 (page 111), l'erreur de restitution sur les deux variables  $f_1$  et  $f_2$  a diminué. Les erreurs commises valent désormais respectivement  $\sqrt{0,71} = 0,84$  et  $\sqrt{1,47} = 1,21$ . Le facteur de synergie est très simple à calculer puisque pour chacune des deux variables, il suffit de calculer le rapport entre l'erreur de restitution obtenue à la Section 4.2 (page 111) avec celle obtenue ici. En effet, considérer qu'il n'y a pas de synergie (*i.e.*, définir toutes les matrices d'erreur diagonales) revient à restituer chaque variable indépendamment. L'erreur de restitution minimum commise en restituant séparément chaque variable est directement celle présentée dans la Section 4.2 (page 111), soit : 0,87 sur la variable  $f_1$  et 1,32 pour  $f_2$ . On a alors un facteur synergique pour chaque variable qui vaut :

$$F_{syn1} \approx 103\%$$

$$F_{syn2} \approx 108\%$$

Il y a donc de la synergie pour les deux variables. Elle est plus forte pour  $f_2$ , ce qui n'est pas surprenant puisque l'erreur de première ébauche sur  $f_2$  est plus élevée. Il est alors plus "facile" de réduire l'erreur de restitution sur  $f_2$  en ajoutant de l'information que l'erreur de restitution sur  $f_1$ .

Pour aller plus loin dans l'étude de la synergie, nous avons étudié la variation de l'erreur de restitution de  $f_1$  en fonction de la covariance d'erreur entre  $f_1$  et  $f_2$ . Nous pouvons calculer l'erreur de restitution de  $f_1$  en utilisant seulement  $B_2$ , même si  $B_2$  n'est pas sensible à  $f_1$ . D'après l'équation (4.1), les termes non diagonaux de  $S_f$  (considérés comme non nuls ici) entraînent des termes en  $q_{1,1}$ , même si  $a_{1,1} = 0$ . Cette restitution n'a pas beaucoup de sens, mais comme nous le verrons à la Section 4.6 (page 118), cette considération, indépendante des différents types de synergie, n'existe pas, et si  $S_f$  n'est pas diagonale,  $A$  ne peut l'être. Si les variables à restituer sont corrélées, alors une observation qui est sensible à une variable sera sensible à l'autre. Nous nous plaçons ici dans un cas purement théorique en restituant  $f_1$  avec  $B_2$  afin de mettre en valeur la synergie indirecte.

La Figure 4.4 présente la variation de l'erreur de restitution de  $f_1$  en utilisant les observations de façon séparée ou conjointe, en fonction de la covariance d'erreur de première ébauche des variables  $f_1$  et  $f_2$ . La courbe bleue représente l'erreur de restitution de  $f_1$ , en ne prenant en compte que  $B_1$  (*i.e.*,  $a_{2,2} = 0$  dans l'équation (4.1)), la courbe verte en

utilisant que  $B_2$  (*i.e.*,  $a_{1,1} = 0$  dans l'équation (4.1)) et la courbe rouge en les utilisant conjointement. On représente également en noir le facteur de synergie qui correspond au rapport entre l'erreur de restitution conjointe (courbe rouge) sur le minimum des erreurs de restitutions séparées (courbe bleue et verte).

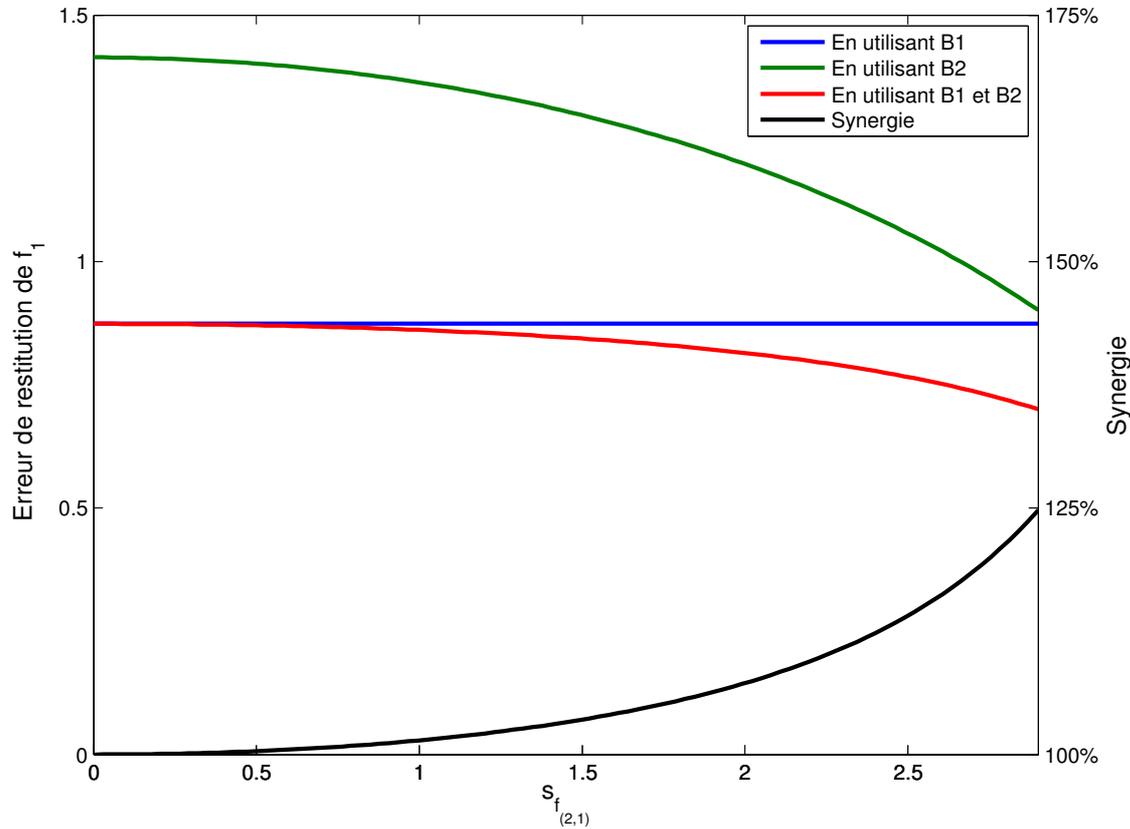


FIGURE 4.4 – Étude de la variation de la synergie indirecte, en fonction de la matrice de covariance d'erreur  $S_F$ . Les courbes représentent l'erreur de restitution en utilisant seulement  $B_1$  (bleu), seulement  $B_2$  (vert) et les deux conjointement (rouge). La synergie de la restitution est tracée en noir par rapport à un axe des ordonnées différent représenté à droite.

Le facteur de synergie est tout le temps supérieur à 100%, et la courbe rouge est tout le temps sous les courbes bleue et verte. Plus la covariance d'erreur entre les variables  $f_1$  et  $f_2$  augmente, plus la synergie augmente et plus l'erreur de restitution, en utilisant les observations simultanément, diminue. Cette variation est normale, car plus les deux variables à restituer seront corrélées, plus il sera facile d'exploiter les informations sur l'autre variable pour en restituer un. Ici, l'augmentation de la covariance d'erreur entre  $f_1$  et  $f_2$  rend l'utilisation de  $B_2$  de plus en plus utile et diminue ainsi l'erreur de restitution de  $f_1$  en utilisant uniquement  $B_2$  mais aussi en utilisant les deux observations conjointement.

La synergie indirecte existe bien et peut être utilisée pour réduire les erreurs de restitution, même si l'observation que l'on ajoute à la restitution a une erreur indépendante de

restitution plus élevée.

## 4.5 Synergie de débruitage

Dans le cadre de la synergie de débruitage, la matrice, qui est considérée non diagonale dans l'équation (4.1), est la matrice de covariance d'erreur des observations  $S_\varepsilon$ . Dans cet exemple, nous prenons  $S_\varepsilon = \begin{pmatrix} 1 & 0,4 \\ 0,4 & 2 \end{pmatrix}$ . Les autres matrices de l'équation ne sont pas modifiées :  $A = \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}$  et  $S_f = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$ . On obtient alors :

$$Q = \left( \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix}^t \cdot \begin{pmatrix} 1 & 0,4 \\ 0,4 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 0,9 & 0 \\ 0 & 0,7 \end{pmatrix} + \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}^{-1} \right)^{-1}$$

$$Q \approx \begin{pmatrix} 0,74 & 0,17 \\ 0,17 & 1,70 \end{pmatrix}$$

L'erreur de restitution de la variable  $f_1$  est alors de  $\sqrt{0,74} = 0,86$  et de  $\sqrt{1,70} = 1,30$  pour  $f_2$ . En comparant aux résultats obtenus en Section 4.2 (page 111), on obtient un facteur de synergie de  $F_{syn} = 101\%$  pour les deux variables  $f_1$  et  $f_2$ . Il y a donc une synergie de débruitage. On a choisi cette dénomination pour ce phénomène de synergie car c'est la covariance d'erreur entre les observations, qui permet de réduire la covariance d'erreur sur les variables. C'est la corrélation entre les bruits instrumentaux qui diminue l'erreur de restitution.

L'étude de la variation de l'erreur de restitution de  $f_1$ , en fonction de la covariance d'erreur entre les observations  $s_{\varepsilon_{2,1}}$ , permet de mieux caractériser la synergie de débruitage. Dans ce cas, on ne peut pas restituer  $f_1$  en utilisant uniquement  $B_2$ , car le fait que la matrice  $A$  soit diagonale fait que les termes non-diagonaux de  $S_\varepsilon$  ne sont pas pris en compte dans  $q_{1,1}$ , si  $a_{1,1}$  est nul. L'erreur de restitution de  $f_1$ , en utilisant seulement  $B_2$ , serait directement l'erreur de première ébauche de  $f_1$ .

La Figure 4.5 présente les résultats obtenus. La courbe bleue correspond à l'erreur de restitution de  $f_1$  en utilisant uniquement  $B_1$  et la courbe rouge en utilisant les deux observations simultanément. Ici encore on représente en noir le facteur de synergie qui correspond au quotient entre la courbe bleue et la courbe rouge.

On remarque, sur cette figure, que la restitution, en utilisant les observations conjointement, est toujours meilleure. La synergie est toujours supérieure à 100%. Il y a bien une synergie de débruitage, qui permet de réduire l'erreur de restitution des variables si la covariance des erreurs d'observations est non nulle. Comme attendu, la synergie est une fonction croissante de la corrélation d'erreur sur les observations. Plus les erreurs seront corrélées plus il sera aisé d'extraire l'information des observations.

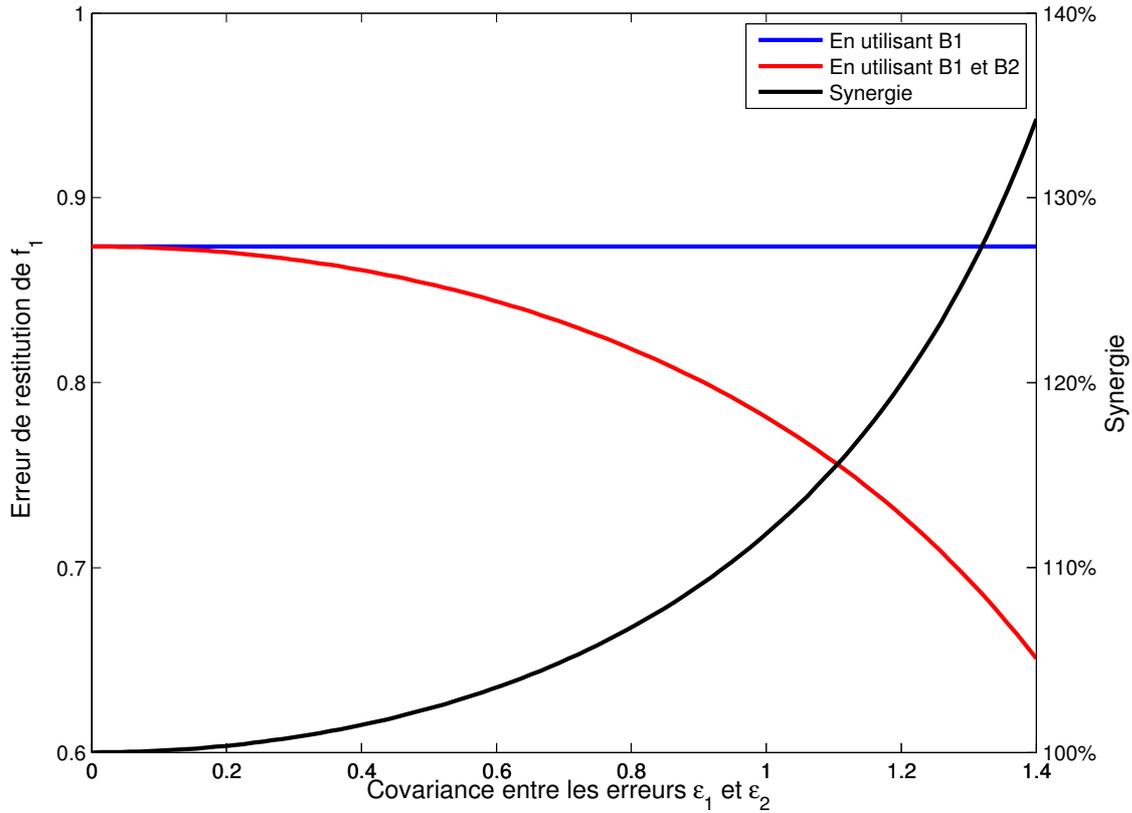


FIGURE 4.5 – Étude de la variation de la synergie de débruitage en fonction de la matrice de covariance d’erreur  $S_\epsilon$  des observations  $B_1$  et  $B_2$ . Les courbes représentent l’erreur de restitution en utilisant seulement  $B_1$  (bleu) et  $B_1$  et  $B_2$  conjointement (rouge). La synergie de la restitution est tracée en noir par rapport à un axe des ordonnées différent représenté à droite.

## 4.6 Généralisation

Dans le cas général, il est difficile de dissocier les différentes synergies. Il est rare que les matrices présentes dans l’équation (4.1) soient diagonales. Cependant, l’approximation sur la matrice de covariance d’erreur des observations  $S_\epsilon$  est souvent effectuée afin de simplifier les calculs. Cette approximation entraîne également une augmentation de l’erreur de restitution comme nous l’avons vu précédemment avec les différents types de synergie.

Nous allons ici considérer le cas général, où aucune des matrices n’est diagonale. Pour reprendre les exemples précédents, on prend :  $A = \begin{pmatrix} 0,9 & 0,7 \\ 0,9 & 0,7 \end{pmatrix}$ ,  $S_\epsilon = \begin{pmatrix} 1 & 0,4 \\ 0,4 & 2 \end{pmatrix}$  et

$S_f = \begin{pmatrix} 2 & 1,5 \\ 1,5 & 3 \end{pmatrix}$ . On a alors :

$$Q = \left( \begin{pmatrix} 0,9 & 0,7 \\ 0,9 & 0,7 \end{pmatrix}^t \cdot \begin{pmatrix} 1 & 0,4 \\ 0,4 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 0,9 & 0,7 \\ 0,9 & 0,7 \end{pmatrix} + \begin{pmatrix} 2 & 1,5 \\ 1,5 & 3 \end{pmatrix}^{-1} \right)^{-1}$$

$$Q \approx \begin{pmatrix} 0,60 & -0,19 \\ -0,19 & 0,95 \end{pmatrix}$$

Les erreurs de restitution de  $f_1$  et  $f_2$  valent alors respectivement  $\sqrt{0,60} = 0,77$  et  $\sqrt{0,95} = 0,97$ , ce qui est inférieur à ce que l'on obtenait précédemment. On peut calculer le facteur synergique pour chaque variable :

$$F_{syn_1} \approx 113\%$$

$$F_{syn_2} \approx 135\%$$

La synergie dans le cas général est une combinaison des différentes formes de synergie. Le facteur de synergie est alors élevé. L'erreur de restitution est diminuée de 35% pour  $f_2$ . Il est donc important de prendre en compte toutes ces différentes synergies dans les algorithmes de restitution.

## 4.7 Conclusion

Ce chapitre peut paraître très théorique, utilisant des exemples synthétiques, mais les conclusions à en tirer sont, elles, capitales et pratiques pour les algorithmes de restitution.

Nous pouvons résumer les différentes synergies mises en avant ici :

- **La synergie additive** : Si les observations disponibles sont sensibles à une même variable, il faut les utiliser conjointement et simultanément dans l'algorithme de restitution. La combinaison *a posteriori* des restitutions indépendantes, ainsi que l'utilisation hiérarchique (à la suite les unes des autres) ne sont pas des solutions optimales.
- **La synergie indirecte** : Si des corrélations existent entre plusieurs variables à restituer, alors la méthode optimale consiste en des restitutions simultanées qui utilisent ces informations de corrélation.
- **La synergie de débruitage** : Si on dispose d'une information sur la corrélation entre les bruits des observations, il faut utiliser cette information.

Comme présenté dans l'Annexe B (page 211), il y a plusieurs sortes d'algorithmes d'inversion. Les algorithmes statistiques, comme les plus proches voisins, la régression linéaire et les réseaux de neurones, prennent naturellement en compte la synergie par construction. Il suffit, pour cela, que la synergie existe dans la base d'apprentissage, elle sera alors utilisée pour paramétrer la restitution.

Les algorithmes analytiques, qui utilisent du calcul matriciel direct, doivent prendre en compte la synergie de façon explicite. Comme nous l'avons vu dans ce chapitre, il est important de préciser les termes non-diagonaux des différentes matrices qui sont malheureusement souvent considérés comme nuls. Cette approximation augmente l'erreur de restitution car elle néglige la synergie existante entre les observations ou les variables.

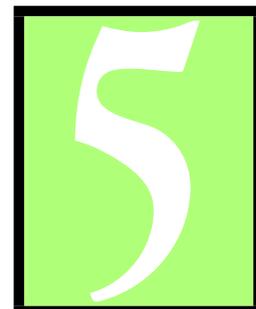
Nous allons utiliser par la suite des méthodes statistiques pour illustrer nos propos dans le cadre de la télédétection atmosphérique. Nous n'utiliserons pas de méthodes analytiques pour la restitution des variables atmosphériques, car, comme nous l'avons montré ici, la synergie est alors à définir et à mettre en oeuvre. Il faut paramétrer les différentes matrices afin de bien prendre en compte les différentes synergies qui entrent en jeu. Notre choix se porte sur les méthodes statistiques car il est plus simple de créer une base d'apprentissage bien représentative, que de modéliser analytiquement les différentes synergies.

Dans ce chapitre nous avons détaillé les cas les plus simples de la synergie afin de mettre en avant les mécanismes synergiques. On peut également considérer des synergies plus complexes, voire non linéaires. Dans le cadre de la télédétection atmosphérique cela peut se traduire par des instruments sensibles à certaines variables mais qui peuvent saturer si d'autres variables dépassent certains seuils. Il se peut également qu'une observation soit une fonction linéaire d'un produit plus ou moins complexe de variables. Il est alors compliqué de prendre en compte les synergies dans un modèle analytique.

L'étude menée ici met en avant la synergie sans pour autant couvrir la totalité des synergies existantes. L'essentiel est de montrer qu'il faut prendre en compte les éléments non-diagonaux des différentes matrices mises en jeu. Dans certains cas, l'écriture matricielle de l'inversion est plus ou moins complexe. La simplification des matrices de covariance entraîne certes une simplification du problème inverse, mais également une perte d'information et notamment une perte d'une partie de la synergie qui pourrait être utilisée afin de réduire l'erreur de restitution.

---

# SYNERGIE INFRAROUGE ET MICRO-ONDE AU-DESSUS DES OCÉANS, EN CIEL CLAIR



## Sommaire

---

<b>5.1</b>	<b>Base d'apprentissage</b>	<b>123</b>
5.1.1	Données atmosphériques utilisées	123
5.1.2	Échantillonnage de la base	124
5.1.3	Simulation de transfert radiatif et bruit instrumental	127
5.1.4	Réduction de la dimension des données IASI	129
5.1.4.1	Sélection de canaux	129
5.1.4.2	Compression de canaux	130
5.1.4.3	Statistiques de compression	131
5.1.4.4	Statistiques de débruitage	133
<b>5.2</b>	<b>Restitutions atmosphériques</b>	<b>134</b>
5.2.1	Corrélations entre les variables et les mesures	136
5.2.2	Restitutions par $k$ -plus proches voisins	139
5.2.2.1	Méthode	139
5.2.2.2	Résultats	141
5.2.3	Restitutions par régression linéaire	145
5.2.3.1	Méthode	145
5.2.3.2	Résultats	145
5.2.4	Restitutions grâce à un réseau de neurones	147
5.2.4.1	Méthode	147
5.2.4.2	Résultats	148
5.2.5	Comparaison des trois méthodes	150
5.2.6	Évaluation de la synergie	152
<b>5.3</b>	<b>Conclusion</b>	<b>154</b>

---

Après avoir présenté dans le chapitre précédent les principes de la synergie d'un point de vue théorique, l'objectif est maintenant de montrer leur utilisation dans un cas concret. Le sondage atmosphérique au-dessus des surfaces continentales est plus complexe qu'au-dessus des surfaces océaniques. Comme présenté dans la Section 2.1 (page 40), les caractéristiques de surface sont peu variables au-dessus des surfaces océaniques. Il est donc plus simple

d'effectuer des restitutions de profils atmosphériques au-dessus des océans. Dans ce chapitre, seules les situations au-dessus des océans sont considérées. Ce choix a été effectué, dans un premier temps, afin de mettre en avant la synergie dans le cas le plus simple possible et d'accroître la complexité des restitutions étapes par étapes.

Dans la section 1.2 (page 18), nous avons vu que les rayonnements infrarouges et micro-ondes sont complémentaires. Notamment pour les situations nuageuses :

- le sondage dans l'infrarouge donne accès aux informations au-dessus des nuages et aux nuages en eux-même (car ces derniers absorbent le rayonnement infrarouge) ;
- Le sondage dans le micro-onde donne accès aux informations sous les nuages.

Cette complémentarité n'est pas utilisée dans ce chapitre. Toujours dans l'optique d'étudier le cas le plus simple possible pour mettre en avant la synergie, nous considérons uniquement les situations en ciel clair. Cependant, la sensibilité des deux rayonnements tant à la température de l'atmosphère qu'à la vapeur d'eau suggère déjà une possible synergie entre les rayonnements qui permettrait d'améliorer les restitutions.

Certaines études sur la synergie ont déjà été menées. On peut citer, par exemple, le travail de [Cho and Staelin \(2006\)](#), qui ont utilisé les données d'AMSU pour s'affranchir des nuages dans l'étude des mesures infrarouges de AIRS (Atmospheric InfraRed Sounder). [Li et al. \(2004\)](#) ont mis au point un schéma de restitution variationnel (similaire à ce qui est utilisé dans les NWP, voir Section 1.5.2, page 34) des paramètres nuageux à partir des données de MODIS et de AIRS. [Susskind et al. \(2003\)](#) combinent quant à eux les données de AIRS, AMSU et HSB (Humidity Sounder for Brazil, similaire à MHS) pour restituer des variables tant atmosphériques que de surface. Cependant, la plupart du temps, les algorithmes, qui utilisent la synergie entre les différents capteurs, se contentent de combiner *a posteriori* les produits restitués. Une combinaison de plusieurs produits semblables, en connaissant leur matrice de covariance d'erreurs, permet de diminuer l'incertitude sur ces produits restitués. Cette façon de faire reste moins efficace qu'une combinaison des données en amont, ce qui permet à l'algorithme de restitution d'exploiter tout le potentiel synergique (voire Chapitre 4, page 107) ([Aires et al. 2012](#)).

Dans l'étude qui suit, nous comparons l'apport synergique de l'utilisation simultanée des données de AMSU-A, de MHS et de IASI, pour la restitution de profils de température et de vapeur d'eau. Comme nous l'avons vu à la Section 1.3.1 (page 29), ces capteurs sont particulièrement sensibles à la vapeur d'eau et à la température, donc propices à une telle étude. Nous utiliserons pour cela trois algorithmes de restitutions distincts : les  $k$ -plus proches voisins, une régression linéaire et un réseau de neurones. Il s'agit de trois algorithmes statistiques, qui vont donc nécessiter la construction d'une base d'apprentissage. Ce chapitre a fait l'objet d'une publication : [Aires et al. \(2011a\)](#).

## 5.1 Base d'apprentissage

Afin de pouvoir paramétrer les différents algorithmes utilisés par la suite, il est nécessaire de construire une base d'apprentissage. Celle-ci doit être robuste, avec un nombre de situations raisonnable, afin de ne pas nécessiter trop de calculs, mais également représentative de la variabilité naturelle des situations atmosphériques.

### 5.1.1 Données atmosphériques utilisées

Les algorithmes statistiques de restitution nécessitent une base de variables atmosphériques (température de surface et de l'atmosphère, profils d'ozone et de vapeur d'eau, etc.) ainsi que les observations satellites correspondantes. Dans le cadre de cette étude au-dessus des océans, nous utiliserons les re-analyses de l'ECMWF (European Center for Medium-range Weather Forecast) ([Simmons and Gibson 2000](#); [Uppala et al. 2005](#); [Dee et al. 2011](#); [Berrisford et al. 2011](#)). Il s'agit de re-analyses des données atmosphériques globales effectuées toutes les 6 heures. Ces profils atmosphériques ont été retravaillés pour présenter une certaine cohérence spatiale et temporelle. Nous utilisons ici les données ERA-40 de l'année 2007 (la version plus récente des données ERA-Interim n'était pas encore disponible au début de cette étude). Afin de pouvoir effectuer des simulations de transfert radiatif, nous ne conservons que certaines variables parmi toutes les données disponibles : la température, la vapeur d'eau (l'humidité relative en %) et l'ozone. Ces profils sont disponibles sur 43 niveaux de pression fixe. Ces 43 niveaux seront donc ceux considérés dans ce chapitre, ils sont présentés sur le tableau ci-après. Les profils restitués seront également projetés sur ces niveaux. Nous conservons aussi pour les simulations de transfert radiatif les informations de vent à 10m, de pression et de température à 2m et de température de surface.

Pression (hPa)	0,10	0,29	0,69	1,42	2,611	4,407	6,95	10,37
Niveau	1	2	3	4	5	6	7	8
Pression (hPa)	14,81	20,40	27,26	35,51	45,29	56,73	69,97	85,18
Niveau	9	10	11	12	13	14	15	16
Pression (hPa)	102,05	122,04	143,84	167,95	194,36	222,94	253,71	286,60
Niveau	17	18	19	20	21	22	23	24
Pression (hPa)	321,50	358,28	396,81	436,95	478,54	521,46	565,54	610,60
Niveau	25	26	27	28	29	30	31	32
Pression (hPa)	656,43	702,73	749,12	795,09	839,95	882,80	922,46	957,44
Niveau	33	34	35	36	37	38	39	40
Pression (hPa)	985,88	1005,43	1013,25					
Niveau	41	42	43					

Nous disposons ainsi d'une base de données de plusieurs millions de situations. Ces situations sont filtrées une première fois car nous nous intéressons uniquement aux situations au-dessus des océans en ciel clair. Il ne reste alors qu'environ un million de situations.

### 5.1.2 Échantillonnage de la base

Cette base de données de profils atmosphériques est encore trop volumineuse. Utiliser une telle base pour effectuer un apprentissage est difficile, en particulier en termes de temps de calcul pour la simulation des observations IASI (8461 canaux). Nous souhaitons cependant conserver la variabilité spatiale et temporelle de notre base de données. Une sélection aléatoire de situations parmi cette base risque de nous faire perdre de l'information qu'elle contient. Il faut alors mettre en place un algorithme d'échantillonnage. Aires and Prigent (2007) ont prouvé qu'une méthode de  $k$ -moyennes était un moyen efficace d'effectuer un échantillonnage multivarié. En nous basant sur ces conclusions, nous utilisons donc un algorithme de  $k$ -moyennes pour extraire une base d'apprentissage raisonnable de la base de données complète.

L'algorithme des  $k$ -moyennes (Lloyd 1982) est un algorithme largement répandu pour faire de la classification non supervisée et du regroupement. Son fonctionnement est simple, ainsi que sa mise en place. De plus, il est assez rapide à converger.

L'objectif de cet algorithme est de résumer une base de données à un nombre donné de prototypes. Le nombre de prototypes voulus dépend des objectifs recherchés. Il s'articule en quatre étapes :

1. Les prototypes initiaux sont sélectionnés au préalable, généralement de façon aléatoire, parmi les échantillons de la base de données.
2. Chaque échantillon de la base de données est associé au prototype le plus proche (la distance généralement utilisée ici est une simple distance euclidienne).
3. Une moyenne des échantillons est effectuée sur chaque groupe associé à un prototype. Ces moyennes constituent les nouveaux prototypes.
4. Les étapes 2 et 3 sont répétées jusqu'à ce qu'un critère de convergence soit atteint.

Il s'agit de minimiser la variance intra-classes, c'est-à-dire que tous les spectres associés à un prototype soient les plus semblables possibles, mais également de maximiser la variance inter-classe afin de bien dissocier les différents prototypes.

Cet algorithme est très sensible à la distance utilisée à l'étape 2. Généralement la distance euclidienne est choisie par souci de simplicité, mais il peut être plus intéressant, dans certains cas, d'utiliser une autre distance comme le minimum, le maximum ou la distance de Mahalanobis (voir Section 5.2.2.1, page 139).

Le choix des prototypes initiaux peut également influencer le résultat. Cependant un choix aléatoire de prototype est en général satisfaisant. Des tirages aléatoires multiples permettent de s'assurer de la convergence de l'algorithme.

Le choix du critère de convergence est aussi important. Le critère le plus communément utilisé consiste à mesurer la différence entre les prototypes à chaque itération. Une fois que ces prototypes sont stables (*i.e.*, ils ne changent pas d'une itération à l'autre), la convergence est atteinte. On peut également chercher à minimiser le nombre d'échantillons qui changent de groupe.

Cet algorithme est très couteux en temps de calcul et plus le nombre de prototypes recherchés est grand plus le temps de calcul est grand.

L'avantage de cette méthode d'échantillonnage est qu'elle permet de conserver la variabilité de la base de données, en particulier les propriétés statistiques des distributions de probabilité dans la base d'origine. On parlera ici d'un échantillonnage "statistique". Nous utiliserons par la suite d'autres méthodes qui permettent au contraire de conserver les extrema de la base. On parlera alors d'échantillonnage "uniforme". Chaque méthode présente des avantages dont il est possible de tirer parti, suivant ce que l'on veut faire de la base ainsi échantillonnée. Nous avons fait le choix dans cette étude de conserver la variabilité de la base ce qui permet de mettre sur un pied d'égalité les différents algorithmes utilisés<sup>1</sup>.

Le problème de l'algorithme des  $k$ -moyennes est le temps qu'il met à converger. Le critère de convergence, que nous avons choisi ici, est le nombre de changement des prototypes sélectionnés d'une itération à l'autre. Nous cherchons à extraire 10.000 prototypes de la base complète d'un million de profils. Appliquer directement l'algorithme serait trop long, nous avons donc fait le choix de le faire en deux étapes. Dans un premier temps, nous extrayons de la base de données complète 100 prototypes, dits de première génération. Dans un deuxième temps, nous considérons les profils associés à chacun de ces prototypes comme une nouvelle base de données que nous échantillonons une nouvelle fois à l'aide de l'algorithme des  $k$ -moyennes pour en obtenir 100 prototypes, dits de deuxième génération. Nous avons donc au final  $100 \times 100$  prototypes, soit 10.000 prototypes. Cette hiérarchisation est pratique car elle permet, entre autres, une recherche plus rapide du plus proche voisin. Il suffit de chercher le plus proche voisin parmi les 100 premiers prototypes, puis parmi les 100 prototypes de deuxième génération associés au prototype de première génération choisi. On calcule alors 200 distances plutôt que 10.000, c'est-à-dire un facteur de gain de temps de 50.

La Figure 5.1 présente les profils de température et d'humidité relative des 100 prototypes de première génération (en bleu). On peut noter la grande diversité de ces profils (bien qu'il soient en partie masqués par les autres courbes colorées). Ils recouvrent la variabilité de la température et de la vapeur d'eau dans la base de données complète. Les autres courbes colorées (les vertes et les rouges) sont des profils associés à des prototypes de deuxième génération. Les courbes rouges (les courbes vertes également) sont, ainsi, toutes associées à une courbe bleue donnée, correspondant au prototype de première génération associé à cette famille de prototypes de deuxième génération.

---

1. C'est un choix à faire dans l'échantillonnage de la base de données. On peut chercher à construire un algorithme qui présente des statistiques raisonnables en toute situation (échantillonnage uniforme) ou au contraire chercher à être globalement plus performant (échantillonnage statistique).

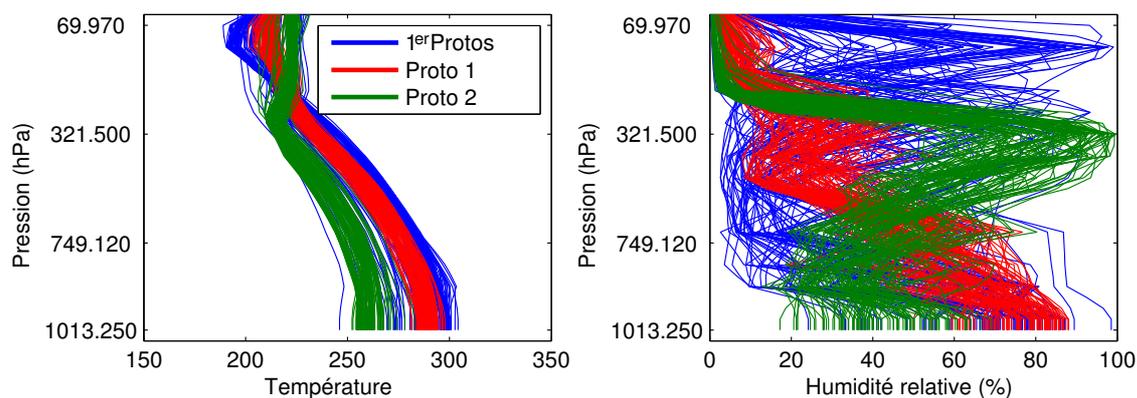


FIGURE 5.1 – Variabilité du profil de température (à gauche) et du profil d’humidité relative (à droite). Les courbes bleues correspondent aux profils associés aux 100 prototypes de première génération. Les courbes vertes et rouges correspondent aux 100 prototypes de deuxième génération associés à deux différents prototypes de première génération (appelés Proto 1 et Proto 2).

Si la variabilité des profils bleus est grande (ce qui est normal, car ils représentent, à eux 100, toute la base de données complète), celle des profils verts et rouges est plus faible. Les profils verts et rouges étant associés chacun à un seul profil bleu (*i.e.*, prototype de première génération), ils sont par conséquent assez semblables entre eux.

La Figure 5.2 présente la répartition géographique des 10.000 situations sélectionnées.

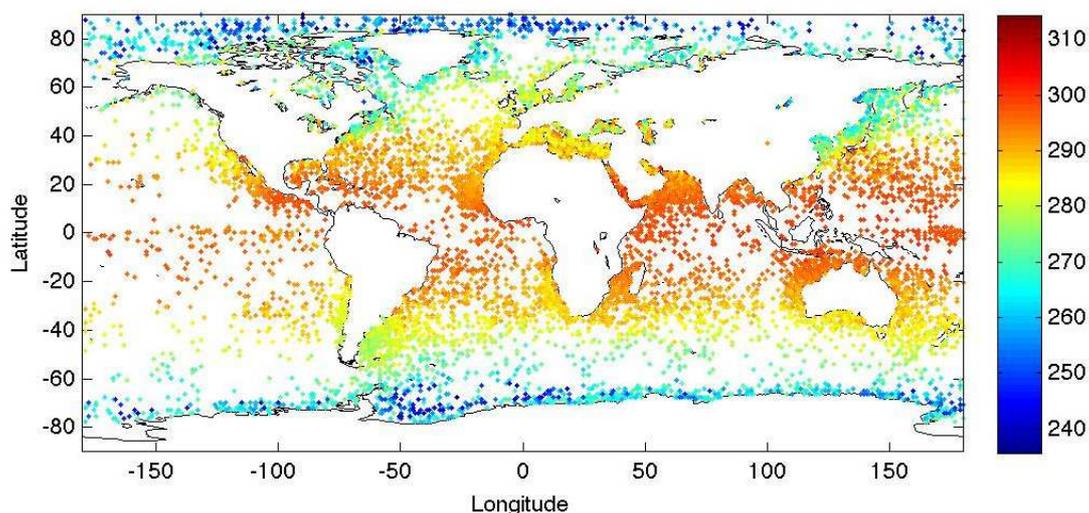


FIGURE 5.2 – Répartition géographique de la base de 10.000 situations échantillonnées par  $k$ -moyennes. La couleur des points dépend de la température de surface.

La couleur des points dépend de la température de surface de chaque situation. L’échantillonnage par  $k$ -moyennes n’a pas pris en compte la répartition géographique de la base

de données originale, mais seulement la variabilité des profils de température et de vapeur d'eau. La base de données échantillonnée est bien répartie sur l'ensemble du globe, tant au niveau des tropiques que des pôles. La surreprésentation des côtes dans la base de données est liée à la grande variabilité des profils proches des côtes. L'algorithme est donc capable d'extraire les profils les plus intéressants de la base de données originale d'un million de situations.

L'algorithme d'échantillonnage a alors réussi à conserver la répartition géographique des profils sans que l'on ait eu à l'intégrer dans l'algorithme. Ceci permet de valider l'échantillonnage. De même, on pourrait montrer que la variabilité temporelle de la base a été respectée. On retrouve les douze mois de l'année de la base d'origine.

### 5.1.3 Simulation de transfert radiatif et bruit instrumental

Une fois que ces 10.000 situations ont été choisies, nous utilisons RTTOV pour associer à chaque profil les températures de brillance correspondant à AMSU-A, MHS et IASI. Les différents profils de températures de brillance sont bruités selon le bruit de mesure de chaque instrument. Cette démarche vise à se rapprocher des données satellites réelles et ainsi donner plus de poids à notre étude qui porte uniquement sur des températures de brillance simulées par RTTOV. À chaque canal d'AMSU-A et MHS est associé un bruit blanc gaussien dont l'écart-type correspond aux informations fournies par le constructeur :

– AMSU-A :

Canal	1	2	3	4	5	6	7	8
Bruit (K)	0,30	0,30	0,40	0,25	0,25	0,25	0,25	0,25
Canal	9	10	11	12	13	14	15	
Bruit (K)	0,25	0,40	0,40	0,60	0,80	1,20	0,50	

– MHS :

Canal	1	2	3	4	5
Bruit (K)	0,22	0,34	0,51	0,40	0,46

Pour l'instrument IASI, le bruit de chaque canal dépend de sa température de brillance. Il faut alors adapter le bruit blanc gaussien à la situation observée (Aires et al. 2002c). L'écart-type du bruit de chaque canal de IASI est présenté sur la Figure 5.3. La courbe bleue correspond à une situation théorique où tous les canaux sont à 280 K. À chaque situation, l'écart-type  $\sigma_i$ , du bruit du canal  $i$  considéré, peut être calculé en fonction de son écart-type  $\sigma_i^{280}$  à 280 K, en fonction de la température de brillance  $Tb_i$  et de la longueur

d'onde  $\lambda_i$  du canal, suivant la formule :

$$\sigma_i = \frac{\frac{\partial B(280 \text{ K}, \lambda_i)}{\partial T}}{\frac{\partial B(Tb_i, \lambda_i)}{\partial T}} \cdot \sigma_i^{280}$$

où  $B$  est la formule de Planck (voir section 1.1.1, page 7). On peut calculer sa dérivée par rapport à la température pour obtenir :

$$\sigma_i = \frac{e^{\frac{hc}{\lambda_i k \cdot 280}}}{e^{\frac{hc}{\lambda_i k \cdot Tb_i}}} \cdot \frac{Tb_i^2}{280^2} \cdot \sigma_i^{280}$$

Cette formule permet alors de calculer l'écart-type du bruit instrumental de chaque canal, à partir duquel on peut simuler un bruit blanc gaussien. Un exemple sur un spectre réel est présenté sur la Figure 5.3 (courbe verte). Plus la température de brillance du canal est faible, plus le bruit est important.

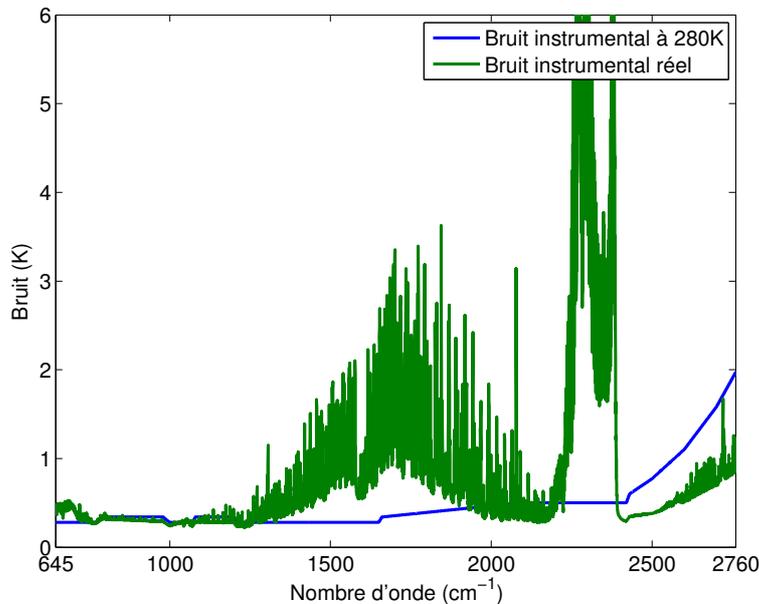


FIGURE 5.3 – Bruit instrumental de l'instrument IASI. La courbe bleue correspond à l'écart-type du bruit sur chaque canal, dans le cas théorique où toutes les températures de brillance valent 280 K. La courbe verte correspond à un exemple d'écart-types du bruit pour un spectre réel.

Les vibrations aléatoires du mécanisme de “bascule”, utilisé au sein de IASI, pour faire varier la largeur de bande de l'interféromètre de Michelson, introduisent des corrélations d'erreur d'un canal à l'autre. N'ayant aucune information sur ces covariances d'erreur, nous ne les prenons pas en compte ici. Il serait toutefois intéressant de pouvoir ajouter cette information dans l'algorithme. Comme présenté à la Section 4.5 (page 117), cette covariance d'erreur des observations permettrait de réduire l'impact du bruit instrumental sur

les inversions et donc de réduire l'erreur de restitution.

Comme expliqué précédemment, la base de données initiale (voir Section 5.1.1, page 123) a dû être échantillonnée afin de permettre un apprentissage plus rapide des méthodes d'inversion. Cet échantillonnage ne donne qu'un nombre restreint de situations (10.000). Le fait d'ajouter, à chaque situation, un bruit instrumental permet d'artificiallement augmenter la taille de la base de données. Ainsi, ce bruit augmente la capacité à généraliser des méthodes de restitution.

La base de données de 10.000 situations est découpée en trois parties distinctes, de la même façon que ce qui a été fait à la Section 2.2.4 (page 56) :

- Une base d'apprentissage de 9.000 situations ;
- Une base de validation de 500 situations ;
- Une base de test de 500 situations.

#### 5.1.4 Réduction de la dimension des données IASI

Le nombre de mesures par sondage de l'instrument IASI (8461) est nettement supérieur à celui des instruments micro-ondes sélectionnés (AMSU-A et MHS ont 20 canaux à eux deux). Ceci est un problème majeur, car la plupart des algorithmes de restitution ne peuvent pas prendre en compte une telle quantité de données. Il est nécessaire d'avoir recours à des méthodes de réduction de dimension des données.

##### 5.1.4.1 Sélection de canaux

Il existe des méthodes qui présentent une approche physique du problème. À partir des mesures de sensibilité de chaque canal à la variable que l'on cherche à restituer, on peut sélectionner quelques canaux particulièrement adaptés à ce que l'on cherche à restituer. Une méthode semblable est mise en place dans Aires et al. (2002a). À partir des fonctions de poids de chaque canal, 442 canaux sont sélectionnés pour représenter l'intégralité du profil vertical de chaque variable atmosphérique. Une telle méthode ne prend pas en compte l'effet de la perte de redondance de l'information sur la qualité de la restitution, on augmente par conséquent l'impact du bruit sur les restitutions.

Une autre grande catégorie de méthodes de compression provient du domaine statistique et de la théorie de l'information. Le nombre de degrés de liberté (Rodgers 2000) ou l'entropie sont utilisés pour sélectionner des canaux ou compresser les données. On peut également citer :

- Les matrices de réduction de données (Menke 1984) ;
- Le minimum d'entropie (Shannon 1949; Huang and Purser 1996) ;
- La décomposition en valeurs singulières (Prunet et al. 1998) ;
- L'approche itérative (Rodgers 2000).

Rabier et al. (2002) comparent quatre méthodes distinctes de sélection de canaux. Certaines méthodes, plus sophistiquées, permettent de compresser les données sans pertes. On

peut citer, par exemple, l'algorithme ZL mis en place par [Ziv and Lempel \(1977\)](#). Ce type d'algorithme utilise la redondance de l'information de manière à la représenter de façon plus concise à l'aide d'un dictionnaire. Ces méthodes sont utilisées au quotidien pour le stockage de données multimédia (*i.e.*, les images gif).

Il est important de savoir si la méthode utilisée perd ou non de l'information. Certaines applications (comme la chimie de l'atmosphère par exemple) ne nécessitent qu'une partie de l'information totale. La perte d'informations n'est alors pas handicapante. D'autres applications (comme le stockage de données d'observations satellites) souffriraient énormément de la perte d'informations.

#### 5.1.4.2 Compression de canaux

La méthode la plus largement utilisée dans la communauté satellite, de nos jours, est l'ACP ( voir Section [2.2.3](#), page [53](#), ou Annexe [B.1](#), page [211](#), pour plus de précisions sur cette méthode), plus souvent appelée Empirical Orthogonal Function (EOF). De nombreuses études ont démontré la capacité de cette méthode à compresser l'information ([Huang and Antonelli 2001](#); [Eriksson et al. 2002](#)) et également à réduire le bruit des mesures satellites ([Aires et al. 2002c](#)). L'ACP permet d'extraire du spectre infrarouge complet les composantes dominantes, afin d'expliquer le maximum de la variance des observations satellites. C'est une transformation linéaire qui est donc facile à intégrer à toutes sortes d'algorithmes de restitution. La représentation du spectre infrarouge par composantes est utile pour transférer une grande quantité de données ou réduire le bruit instrumental. Ces composantes peuvent être directement utilisées dans les algorithmes de restitution ([Aires et al. 2002d,b](#)). Elles peuvent également servir à détecter les nuages ([Smith and Taylor 2004](#)).

Dans le cas de l'ACP sur les spectres IASI, la base de données est constituée de 8461 canaux. On cherche alors une nouvelle base (constituée de vecteurs propres) sur laquelle projeter les données. Comme il a été montré à la Section [2.2.3](#) (page [53](#)), plus le nombre de composantes prises en compte est élevé, plus l'erreur de compression est faible. Cependant, les composantes d'ordre élevé codent des fréquences d'oscillations plus élevées (de la même façon que dans la décomposition en série de Fourier). Le bruit se projette donc sur ces composantes. Les prendre en compte impliquerait plus de représentativité du bruit.

Le bruit de l'instrument IASI est un problème particulièrement important car il atteint un niveau élevé dans certaines zones spectrales (voir Figure [5.3](#)). Utiliser la totalité du spectre IASI permet, par redondance des informations, de diminuer l'impact de ce bruit instrumental.

Par définition, les composantes d'ordre élevé de l'ACP contiennent moins d'informations. Elles ont moins de sens physique car elles sont issues du mélange de la variabilité non expliquée par les premières composantes. Elles ont un rapport signal sur bruit plus faible. L'idée est de séparer les composantes en deux sous-ensembles. Les composantes associées aux premiers vecteurs propres représentent un sous-ensemble "physique" des observations.

Les composantes associées aux vecteurs propres d'ordre plus élevé représentent alors un sous-espace orthogonal correspondant aux variabilités mineures et au bruit. Ne pas prendre en compte ces composantes permettrait alors de s'affranchir du bruit et de ne travailler que sur les composantes contenant de l'information.

En reprenant des notations similaires à celles que l'on avait utilisées à la section 2.2.3 (page 53), on peut écrire l'analyse en composantes principales d'un spectre IASI constitué des températures de brillance  $\{Tb_1, Tb_2, \dots, Tb_{8461}\}$  comme :

$$\{Tb_1, Tb_2, \dots, Tb_{8461}\} = C \cdot VP + \overline{Tb}$$

où  $VP$  est la matrice  $8461 \times 8461$  des vecteurs propres de la matrice de covariance des températures de brillance de IASI.  $C$  correspond à la matrice  $1 \times 8461$  des nouvelles variables dans le nouvel espace de projection.  $\overline{Tb}$  est la moyenne des températures de brillance.

On cherche alors à déterminer  $n$  tel que la base orthornormée  $VP(1 : n, 1 : 8461)$  définisse un nouvel espace de dimension restreinte contenant l'information et  $VP(n + 1 : 8461, 1 : 8461)$  définisse l'espace de projection du bruit instrumental et des informations de variabilité inférieure.

On définit l'erreur de compression de l'ACP avec  $n$  composantes par la racine carrée de la moyenne de l'erreur quadratique du terme  $\{Tb_1, Tb_2, \dots, Tb_{8461}\} - C(1 : n) \cdot VP(1 : n, 1 : 8461)$ . Le critère des moindres carrés, utilisé ici, est particulièrement adapté pour l'ACP (Jolliffe 2002).

En pratique, on calcule les vecteurs propres pour l'ACP sur des spectres IASI non bruités pour obtenir les principales signatures spectrales physiques. Calculer les vecteurs propres sur les spectres bruités fausserait le calcul et certaines composantes seraient constituées essentiellement de bruit. On aurait pu utiliser un filtre de Wien, ou une ACP ajustée pour le bruit (Blackwell 2005), mais l'ACP utilisée ici a déjà fait ses preuves pour la réduction du bruit instrumental sur les observations satellites (Aires et al. 2011a).

#### 5.1.4.3 Statistiques de compression

La Figure 5.4 représente l'erreur instrumentale IASI et diverses erreurs de compression de l'ACP correspondant à divers nombres de composantes utilisées dans la restitution. Chaque spectre est d'abord bruité, puis compressé et décompressé avant d'être comparé au spectre d'origine (sans bruit).

Comme expliqué à la section 5.1.3 (page 127), le bruit instrumental de IASI dépend, pour chaque canal, de la température de brillance observée. Le spectre de bruit, représenté en noir sur cette figure, correspond à l'écart-type du bruit instrumental sur chaque canal pour un spectre IASI moyen en ciel clair. Certaines zones spectrales, particulièrement la bande 3 (2000 à 2760  $\text{cm}^{-1}$ ), sont très bruitées (plusieurs K).

Les courbes colorées correspondent à l'erreur de compression moyenne d'un spectre

bruité par rapport à un spectre sans bruit. La courbe bleue, correspondant à l'erreur de compression en utilisant 5 composantes, montre qu'il est nécessaire de prendre en compte plus de composantes. L'erreur de compression est par endroits supérieure au bruit instrumental (en noir sur la figure), ce qui signifie que de l'information exploitable a été perdue. La bande 3 est par contre déjà bien débruitée avec 5 composantes (1 K d'erreur moyenne contre plus de 2 K de bruit instrumental).

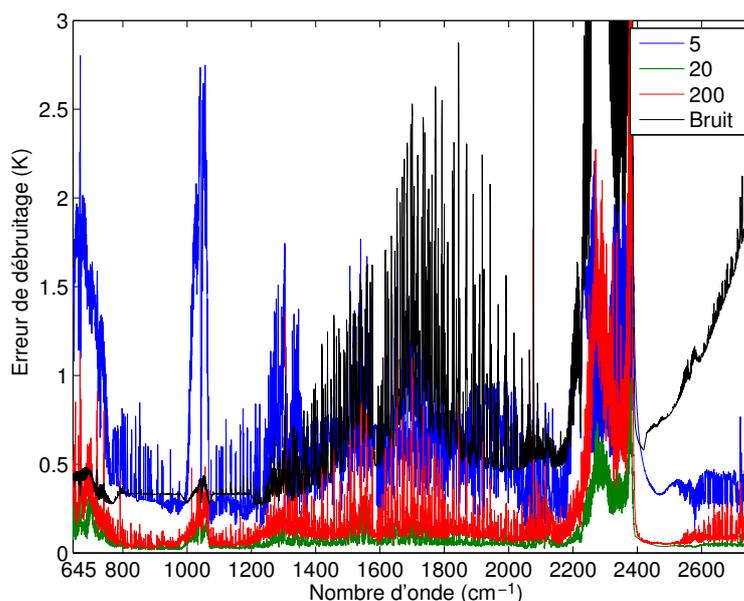


FIGURE 5.4 – Statistiques d'erreur de débruitage de l'ACP. La courbe noire correspond au bruit instrumental. Les courbes colorées correspondent à l'erreur moyenne de compression puis décompression d'un spectre bruité par rapport à un spectre sans bruit pour différents nombres de composantes utilisées (5,20 et 200).

L'utilisation de 20 composantes pour la compression des spectres IASI (courbe verte sur la figure) présente par contre des statistiques intéressantes tout au long du spectre. La racine carrée de l'erreur quadratique moyenne de compression reste proche de 0,1 K pour tous les canaux. L'ACP est alors capable de supprimer le bruit en gardant la majeure partie de l'information contenue dans le spectre. Entre 2200 et 2400  $\text{cm}^{-1}$ , l'erreur de compression reste importante, autour de 0,5 K. Ceci est probablement dû à une contamination atmosphérique par le  $\text{NO}_2$ , qui engendre une variabilité non-expliquée par les premières composantes de l'ACP, mais également dû au bruit qui est important dans cette zone spectrale. Avec 200 composantes (courbe rouge sur la figure), l'erreur de débruitage est plus élevée. Ceci confirme l'idée que les composantes d'ordre plus élevé codent le bruit instrumental et sont donc inutiles voire pénalisantes.

#### 5.1.4.4 Statistiques de débruitage

Nous avons mené une étude plus poussée sur la capacité de l'ACP à débruiter les spectres IASI. Pour ce faire, on compare un spectre IASI non bruité au même spectre compressé puis décompressé, avec différents nombres de composantes. La racine carrée de l'erreur quadratique moyenne alors commise est moyennée spectralement pour obtenir la courbe bleue sur la Figure 5.5. La même méthode est utilisée, en bruitant le spectre avant de le compresser, pour obtenir la courbe rouge (toujours en comparant par rapport au spectre original non bruité).

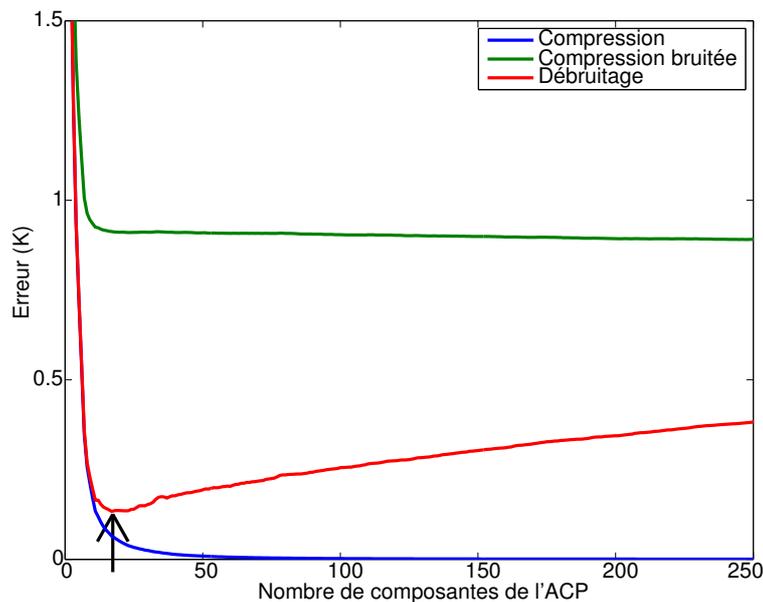


FIGURE 5.5 – Évolution de l'erreur de débruitage de l'ACP. La courbe bleue correspond à l'erreur de compression d'un spectre non-bruité. La courbe rouge correspond à l'erreur de compression d'un spectre bruité par rapport à un spectre non-bruité. La courbe verte correspond à l'erreur de compression d'un spectre bruité par rapport à un spectre bruité.

L'erreur de compression d'un spectre non bruité est, comme attendue, une fonction strictement décroissante du nombre de composantes utilisées (courbe bleue). Pour les premières composantes, la décroissance est plus importante car la variance expliquée est plus élevée. Les composantes suivantes codent ensuite, par définition, une part de plus en plus faible de l'information, ce qui explique le ralentissement de la décroissance. Avec 8461 composantes (non présenté sur la courbe), on retrouverait le spectre complet, soit une erreur de 0 K.

La courbe rouge présente, quant à elle, une forte décroissance suivie d'une lente croissance convergeant vers 1 K (*i.e.*, le bruit moyen du spectre IASI) pour 8461 composantes. La courbe verte correspond à l'erreur de compression d'un spectre bruité mais par rapport à un spectre bruité également. La courbe semble se stabiliser autour de 1 K puis diminue ensuite pour atteindre 0 pour 8461 composantes (non visible sur ce graphique). Ceci correspond à

la partition des composantes que l'on présentait précédemment.

Les  $n$  premières composantes (ici une vingtaine) codent l'information contenue dans le spectre. On peut alors reconstruire le spectre original de façon de plus en plus précise en augmentant le nombre de composantes (la courbe rouge décroît). Les composantes suivantes codent le bruit. Plus on prend en compte un nombre élevé de ces composantes, plus le bruit sera important dans le spectre reconstitué (la courbe rouge croît). L'erreur de reconstruction du spectre sera alors de plus en plus grande. Le fait que la courbe verte ne diminue pas aussi rapidement que les autres vient du fait qu'ici l'erreur est mesurée par rapport à un spectre bruité. Les premières composantes, qui ne permettent pas de coder le bruit, ne permettent alors pas de diminuer l'erreur de compression du spectre. Seule la totalité des composantes permet de coder l'information de bruit manquante. Le meilleur compromis entre compression et débruitage consiste donc à prendre en compte 20 composantes.

Il est important de noter qu'il s'agit d'un point de vue purement statistique, dans l'optique des restitutions à venir. Une utilisation différente des observations IASI nécessiterait un algorithme différent. Il peut paraître surprenant de n'utiliser que 20 composantes pour représenter un spectre de 8461 canaux. Afin de justifier notre point de vue, la restitution par réseaux de neurones présentée ci-après a été recalculée en utilisant 100 composantes (Aires et al. 2011a). La capacité d'un réseau de neurones à ne pas prendre en compte les informations inutiles a permis d'obtenir des résultats en tout point semblables à ceux n'utilisant que 20 composantes. Ces résultats viennent justifier le choix qui a été fait de n'utiliser que 20 composantes<sup>2</sup>.

## 5.2 Restitutions atmosphériques

Actuellement la méthode la plus répandue, à même d'exploiter la synergie entre les différentes mesures atmosphériques, est l'assimilation (voir Section 1.5.2, page 34). Le principe de base consiste à combiner les sorties des modèles de prévision avec toutes les autres données disponibles (satellites, bouées, radiosondages, mesures *in situ*, etc.), de façon inversement proportionnelle à leurs incertitudes respectives. L'utilisation en parallèle de toutes les données permet, par redondance de l'information, de diminuer l'erreur de restitution. Nous présentons ici un cas simple, où plusieurs mesures satellites sont utilisées simultanément, afin de mettre en évidence ce principe. Dans notre situation, nous utilisons les sondages d'AMSU-A, MHS et IASI pour restituer des profils de température et de vapeur d'eau. Dans le cas général, la fusion de toutes les sources de l'information optimise la restitution des différents paramètres atmosphériques.

---

2. Au risque d'introduire ici une légère répétition, cette ACP est adaptée à l'utilisation que l'on en fait et pas forcément à d'autres. La forte influence des profils de température et de vapeur d'eau sur le spectre IASI permet de retrouver les informations sur ces profils dans les premières composantes de l'ACP. Une étude sur des molécules mineures de l'atmosphère perdrait beaucoup à n'utiliser que les composantes calculées ici. Par exemple, dans le chapitre 2 (page 37), on a préféré une sélection de canaux à une ACP.

Comme les variables que l'on cherche à restituer sont multivariées et corrélées les unes aux autres, les matrices de covariance d'erreur (des observations satellites et des variables géophysiques) sont importantes. Ces dernières sont directement impliquées dans les algorithmes de restitution synergique (voir chapitre 4, page 107). La principale limite de l'assimilation réside dans le mélange des données satellites et des sorties de modèles de prévision. Ce qui est pris en compte dans la restitution n'est pas directement l'information contenue dans la mesure satellite, mais son impact sur l'algorithme global. On n'utilise alors uniquement l'information supplémentaire par rapport à ce que contient déjà le système. Avec les progrès des modèles et la quantité toujours croissante de données disponibles, cette information supplémentaire devient de plus en plus négligeable. Pour autant, l'information en elle-même doit être prise en compte. Il faut alors que l'on mette en place un véritable algorithme synergique.

En effectuant ces restitutions atmosphériques, nous cherchons donc à :

- Fournir des restitutions indépendantes des modèles qui permettront de valider les modèles de prévision météorologique *a posteriori* (voire climatologique) ;
- Définir des algorithmes de restitution synergique qui contiennent toute la chaîne d'utilisation des données d'observations satellites. En passant par la compression, le débruitage, la fusion, une calibration éventuelle des données et un modèle de transfert radiatif adapté, il faut mettre en place une méthode multivariée et flexible.

Il y a deux façons de concevoir un algorithme multivarié : hiérarchiquement (*i.e.*, les uns après les autres) ou simultanément. Dans un schéma hiérarchique, il peut s'agir d'utiliser les différentes observations tour à tour ou de restituer les variables les unes à la suite des autres. Une restitution hiérarchique peut également utiliser les variables déjà inversées pour restituer les suivantes. Cette approche repose sur l'interdépendance entre les variables atmosphériques et le schéma de restitution suit cette structure de dépendances. Un tel schéma a été développé dans le ISCCP ("International Satellite Cloud Climatology Project") (Rossow and Schiffer 1999). Dans un premier temps, la présence de nuage est déterminée, en utilisant des mesures infrarouges issues des satellites géostationnaires et polaires. Dans un deuxième temps, la température de surface est estimée pour les situations claires. Le schéma de restitution des L2 IASI d'Eumetsat, qui sera présenté à la section 6.2.1 (page 160), est également hiérarchique. Il utilise les données tour à tour pour converger, de façon semblable à ce qui est fait dans les NWP. Un tel schéma peut être problématique car une erreur dans les étapes préliminaires se répercutera dans toute la suite de la chaîne. Il sera ensuite difficile d'évaluer les incertitudes finales de restitution sur les variables restituées.

L'autre approche consiste à effectuer la restitution multivariée simultanément. On peut restituer les variables simultanément ou utiliser toutes les sources d'informations simultanément. Dans les deux cas, la détermination des incertitudes de restitution est alors plus facile (Aires 2004; Aires et al. 2004a,b). Il est plus simple d'utiliser la synergie entre les différentes sorties et les différentes entrées de l'algorithme en les considérant de façon simultanée. Lors-

qu'on construit une solution à un problème d'optimisation, il est toujours préférable d'avoir une solution globale, plutôt que de la construire par morceaux.

Comme expliqué au chapitre 4 (page 107), il est important de prendre en compte la synergie entre les variables et les mesures utilisées. Ceci peut-être fait par la construction d'une base d'apprentissage adéquate (Aires and Prigent 2007). La base de données construite à l'aide des  $k$ -moyennes contient ces informations d'interdépendances. Le fait d'avoir respecté la variabilité naturelle des variables et gardé les profils de température et de vapeur d'eau associés entre eux, permet de conserver les relations entre ces variables. On s'attend, dès lors, à pouvoir mesurer l'impact de la synergie sur les erreurs de restitution. On met en place différents algorithmes afin de restituer les profils de température ou d'humidité relative à partir des mesures infrarouges et/ou micro-ondes.

### 5.2.1 Corrélations entre les variables et les mesures

Les figures suivantes présentent les différentes matrices de corrélation mises en jeu dans le problème considéré ici. L'algorithme de restitution, que l'on souhaite mettre en place, a comme entrées les observations satellites (les composantes d'ACP pour l'infrarouge et les températures de brillance d'AMSU-A et MHS pour le micro-onde), et les profils de température et d'humidité spécifique en sortie.

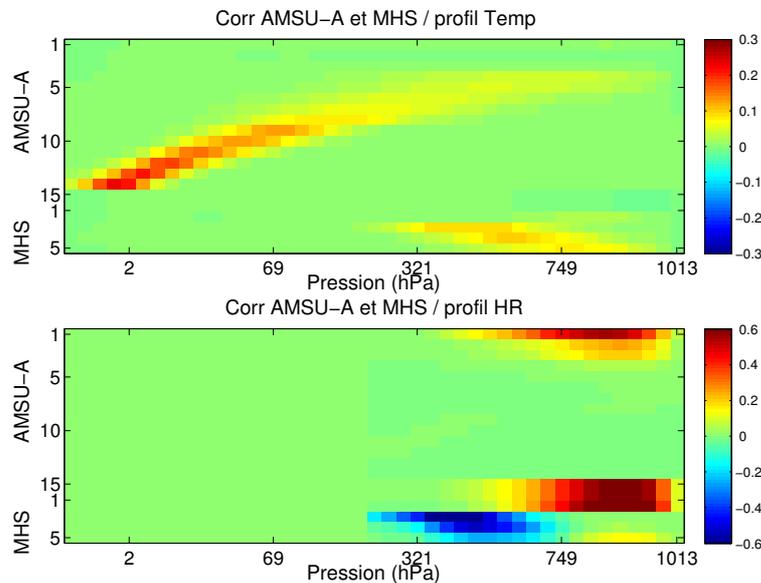


FIGURE 5.6 – Corrélation entre les mesures dans le micro-onde et le profil de température (en haut) et le profil d'humidité relative (en bas).

Les matrices de corrélation de la Figure 5.6 présentent la sensibilité des mesures micro-ondes aux variables atmosphériques considérées (ce qui correspond à la synergie directe). On peut noter une forte corrélation entre la température et les mesures micro-ondes tout le long du spectre. AMSU-A est plus particulièrement sensible dans le haut de l'atmosphère

(partie gauche du graphique), et MHS est plutôt sensible à la température dans le bas de l’atmosphère. La sensibilité des mesures micro-ondes à l’humidité relative semble être nulle dans la partie haute de l’atmosphère, mais compte tenu de la faible quantité de vapeur d’eau qui y est présente, cette faible sensibilité est logique. On remarque néanmoins la très forte sensibilité de MHS à la vapeur d’eau jusque dans le milieu de l’atmosphère.

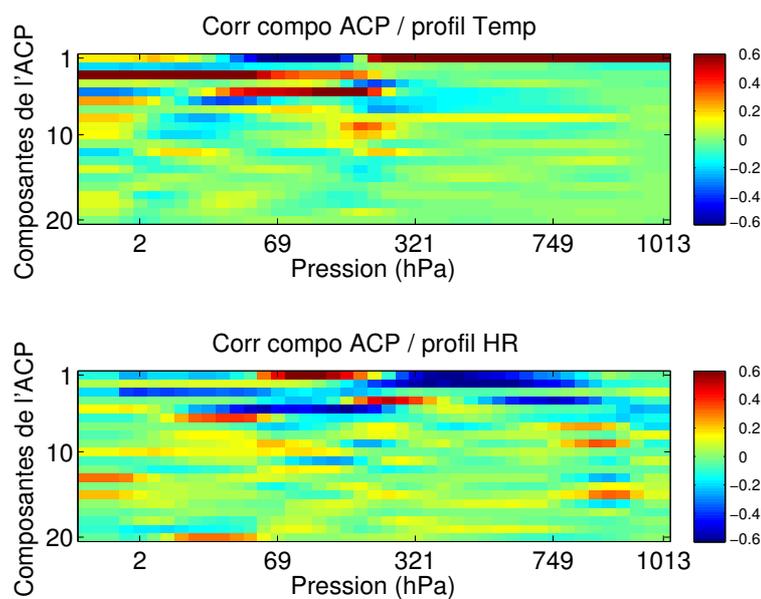


FIGURE 5.7 – Corrélation entre les composantes de l’ACP sur le spectre infrarouge de IASI et le profil de température (en haut) et le profil d’humidité relative (en bas).

Les matrices de corrélation de la Figure 5.7 présentent la sensibilité des composantes de l’ACP sur le spectre IASI aux variables atmosphériques considérées (synergie directe). L’interprétation de ces matrices est moins aisée que précédemment. Du fait de l’utilisation de l’ACP, les différents canaux de IASI sont “mélangés” et on distingue mal la sensibilité de chaque canal à différentes couches atmosphériques. On peut toutefois remarquer, que ce soit pour la température ou l’humidité relative, une forte sensibilité des composantes IASI à tout le profil, depuis la surface jusqu’au sommet de l’atmosphère.

Ces matrices justifient le choix du schéma de restitution. En effet, les entrées présentent une sensibilité importante aux différentes sorties. L’algorithme que l’on souhaite mettre en place doit donc exploiter ces sensibilités pour restituer les variables atmosphériques.

La matrice de corrélation de la Figure 5.8 présente, quant à elle, les corrélations entre les différentes entrées de l’algorithme de restitution. Cela correspond à la synergie de débruitage présentée au Chapitre 4 (page 107), mais également à une partie de la synergie directe car les corrélations que l’on montre ici ne sont pas uniquement des corrélations entre les erreurs des différentes mesures.

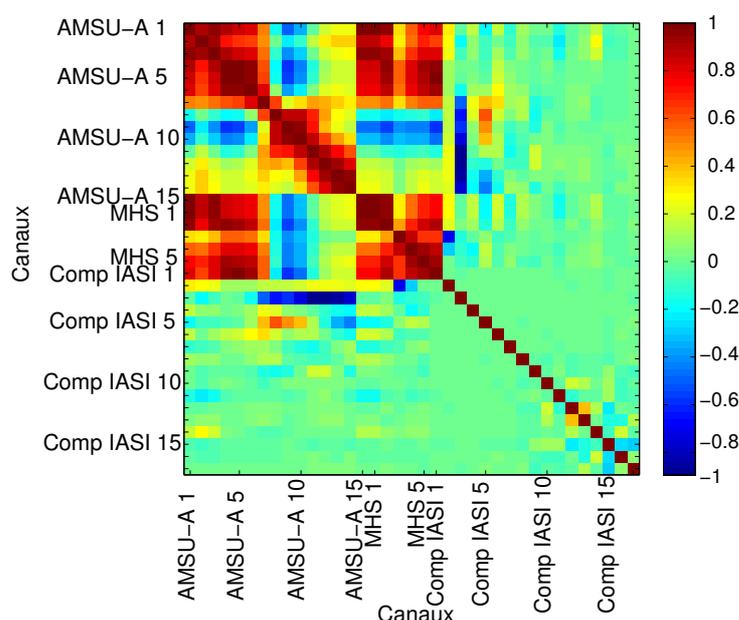


FIGURE 5.8 – Corrélation entre les composantes de l’ACP sur le spectre infrarouge et les mesures micro-ondes.

Les corrélations entre les composantes de IASI et les mesures micro-ondes sont représentées. Par construction, les différentes composantes sont indépendantes entre elles. La partie en bas à droite de la matrice est donc quasi-diagonale. Si les corrélations entre les différents canaux de AMSU-A et MHS sont fortes, on note quand même des corrélations entre les premières composantes de IASI et certains canaux d’AMSU-A. Les composantes d’ordre plus élevé contiennent une information de moindre importance par rapport aux premières, ce qui explique que les corrélations soient plus faibles. Il y a cependant une certaine redondance de l’information du fait de ces corrélations (synergie directe). C’est cette redondance que l’algorithme de restitution mis en place doit savoir exploiter pour utiliser la synergie.

La matrice de corrélation de la Figure 5.9 représente les corrélations entre les différentes sorties de l’algorithme de restitution. Les profils de température et de vapeur d’eau sont très corrélés. On remarque, par exemple, un lien très fort entre les basses couches de température et les couches moyennes de vapeur d’eau ; ou encore une forte anti-corrélation dans les hautes couches. Ces corrélations semblent montrer qu’il est possible d’utiliser la synergie indirecte pour la restitution de la température et de la vapeur d’eau.

Il reste désormais à mettre en place l’algorithme de restitution. Le problème étant complexe, non-linéaire, avec de nombreux degrés de liberté et une profusion de mesures utilisées conjointement, on est amené à penser qu’un réseau de neurones serait le candidat idéal pour effectuer ces restitutions. Afin de vérifier notre intuition, plusieurs schémas sont étudiés. On varie la structure du schéma avec plus ou moins d’entrées et de sorties, afin de vérifier l’impact de la synergie sur la restitution. On utilise également plusieurs méthodes

différentes pour comparer leur capacité à utiliser la synergie. Deux méthodes statistiques, les  $k$ -plus proches voisins et la régression linéaire, sont mises en concurrence avec un réseau de neurones.

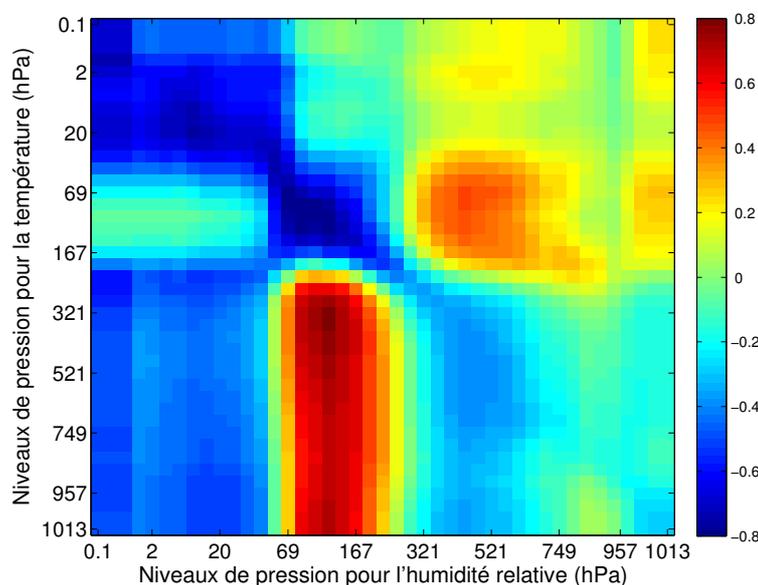


FIGURE 5.9 – Corrélation entre la température et la vapeur d’eau atmosphérique.

## 5.2.2 Restitutions par $k$ -plus proches voisins

### 5.2.2.1 Méthode

La méthode des  $k$ -plus proches voisins utilise une base de données de référence, constituée des variables atmosphériques considérées (*i.e.*, celles que l’on cherche à restituer) associées aux observations satellites correspondantes. La base de référence utilisée ici est la base d’apprentissage construite à la Section 5.1.2 (page 124). On compare chaque nouvelle observation satellite (issue de la base de test) à celles présentes dans la base de référence. Les  $k$ -plus proches observations permettent alors de fournir le profil atmosphérique restitué. Dans le cas où on ne prend en compte que le plus proche voisin (*i.e.*,  $k = 1$ ), l’algorithme est alors une simple LUT (look-up-table). On associe à chaque observation satellite le profil atmosphérique associé à l’observation la plus proche issue de la base de données de référence. Lorsque  $k > 1$ , la solution est donnée par la somme pondérée des profils atmosphériques associés aux  $k$  observations satellites les plus proches dans la base d’apprentissage. La pondération choisie est l’inverse de la distance à l’observation satellite. Plus l’observation de la base de référence est proche de l’observation réelle, plus le profil atmosphérique associé aura un poids important. On a alors le profil restitué  $P_r$  associé à l’observation satellite  $Tb$  en fonction des  $k$  profils atmosphériques de la base d’apprentissage  $\{P_{app}^i\}$  associés aux

observations  $\{Tb_{app}^i\}$  :

$$P_r = \frac{1}{\sum_{i=1}^k \frac{1}{d(Tb, Tb_{app}^i)}} \cdot \sum_{i=1}^k \left( \frac{P_{app}^i}{d(Tb, Tb_{app}^i)} \right)$$

où  $d(Tb, Tb_{app}^i)$  correspond à la distance entre l'observation  $Tb$  et une des  $k$ -plus proches observations de la base d'apprentissage  $Tb_{app}^i$ .

Plus  $k$  est élevé, plus le comportement de la restitution est continu. En effet, si une observation est semblable à deux observations de la base d'apprentissage, en ne prenant en compte que le plus proche voisin, suivant quelle situation est la plus proche, la restitution sera changée brusquement. Si on prend en compte plus de voisins, le changement entre les deux restitutions ne sera pas brutal, la restitution sera une combinaison linéaire des profils associés aux deux (ou plus) observations. Ce paramètre  $k$  contrôle la régularisation de l'algorithme. Le dilemme biais-variance (Geman et al. 1992) nous dit également que plus  $k$  est élevé, plus le biais de la restitution diminue, mais plus l'écart-type de l'erreur de restitution augmente (et inversement). Il faut alors sélectionner  $k$  pour faire un compromis judicieux entre le biais et l'écart-type de l'erreur.

La méthode des  $k$ -plus proches voisins est non-linéaire. Son comportement et ses résultats statistiques convergeront vers ceux obtenus avec un réseau de neurones si la base d'apprentissage est constituée d'assez de situations. Cette méthode n'est pas capable de généraliser la restitution à des situations qui ne lui sont pas données dans la base de référence. Il faut alors que l'espace des solutions soit échantillonné avec un niveau de précision limité uniquement par le bruit instrumental (Rydberg et al. 2009; Jiménez et al. 2007; Evans et al. 2005). Malheureusement, plus la population de la base de référence augmente, plus le temps de restitution par  $k$ -plus proches voisins augmente. Le temps de calcul nécessaire à chaque restitution est une fonction croissante du nombre d'échantillons dans la base d'apprentissage. Avoir une base sur-représentée rendrait donc les restitutions fastidieuses, voire même prohibitives.

Cette méthode est multivariée, même si elle est entièrement basée sur une distance dans l'espace des observations satellites et ne considère pas les relations entre les variables. La distance considérée ne prend pas en compte l'information sur les variables recherchées contenue dans chaque canal. Afin d'éviter que les informations moins pertinentes n'influencent de trop les résultats de la restitution, la distance de Mahalanobis est préférée à la distance euclidienne classique. Si la distance euclidienne accorde le même poids à chaque paramètre, la distance de Mahalanobis pondère les distances par la variabilité du paramètre au sein de la base (Mahalanobis 1936). La distance de Mahalanobis entre une observation satellite  $Tb = \{Tb_1, Tb_2, \dots, Tb_n\}$  quelconque, et une mesure issue de la base de référence  $Tb_{app}^i = \{Tb_{app,1}, Tb_{app,2}, \dots, Tb_{app,n}\}$  peut s'écrire en fonction de la matrice de covariance

des mesures dans la base de référence  $S_{Tb_{app}}$  comme :

$$d(Tb, Tb_{app}) = \sqrt{(Tb - Tb_{app})^T \cdot S_{Tb_{app}}^{-1} \cdot (Tb - Tb_{app})}$$

L'utilisation de cette distance nous permet de prendre en compte les corrélations entre les observations et donc de focaliser la reconnaissance de forme sur les observations de variabilité plus importante au sein de la base. Les températures de brillance non-informatives polluent alors moins la distance et donc la restitution. La méthode idéale serait d'être capable de pondérer chaque distance par sa signification réelle dans l'espace des variables à restituer.

Le lecteur notera ici que les observations infrarouges sont dans l'espace des composantes et non dans l'espace des températures de brillance comme les observations micro-ondes. Cependant, du fait de la linéarité du passage des températures de brillance aux composantes, ce que nous venons de présenter reste valable et nous pouvons nous servir de la distance de Mahalanobis directement sur les composantes. La présentation de la distance de Mahalanobis sur l'espace des composantes aurait nui à son intelligibilité. D'autant que l'apport de la distance de Mahalanobis est faible dans l'espace des composantes seules, celles-ci étant indépendantes entre elles par construction. C'est lorsqu'elles sont mélangées avec les données micro-ondes que cette distance apporte une plus value.

L'inconvénient de cette méthode est le temps nécessaire pour calculer les distances à chaque représentant dans la base d'apprentissage. La base d'apprentissage étant ici constituée de 9.000 situations, il faut calculer 9.000 distances à chaque restitution. Le fait d'avoir échantillonné la base en deux fois à l'aide des  $k$ -moyennes (voir Section 5.1.2, page 124) permet d'accélérer les calculs. Il suffit de calculer les distance pour les 100 prototypes de première génération puis de calculer les 100 distances des prototypes de deuxième génération associés au prototype de première génération le plus proche. Cela permet de ne calculer que 200 distances au lieu de 9.000 et d'accélérer d'autant la restitution.

### 5.2.2.2 Résultats

Cette méthode est donc appliquée en utilisant comme base de référence la base d'apprentissage construite à la Section 5.1.2 (page 124). La base de test est utilisée pour calculer l'erreur de restitution. Une première étude est menée pour déterminer le nombre optimal  $k$  de voisins à prendre en compte. Pour cela, on a calculé la racine carrée de l'erreur de restitution quadratique moyenne sur toutes les variables pour différents nombres de voisins. Cette évolution de l'erreur en fonction de  $k$  est représentée sur la Figure 5.10.

L'erreur diminue jusqu'à un optimum pour  $k = 7$ , puis augmente. Comme énoncé précédemment, une valeur de  $k$  trop faible ne permettra pas à l'algorithme de restituer des structures fines (inversions de profils multiples...). Le principe biais-variance présenté ci-dessus explique l'augmentation de l'erreur pour un nombre de voisins élevé. L'erreur quadratique moyenne est égale à la somme du biais au carré et de la variance de l'erreur. Plus

$k$  augmente, plus le biais diminue et plus la variance augmente. Dans un premier temps, la diminution du biais est telle que l'erreur quadratique augmente. Au bout d'un moment, la forte croissance de la variance de l'erreur entraîne une augmentation de l'erreur quadratique malgré une diminution du biais. Au final, nous décidons de considérer les 7 plus proches voisins, correspondant au nombre optimal.

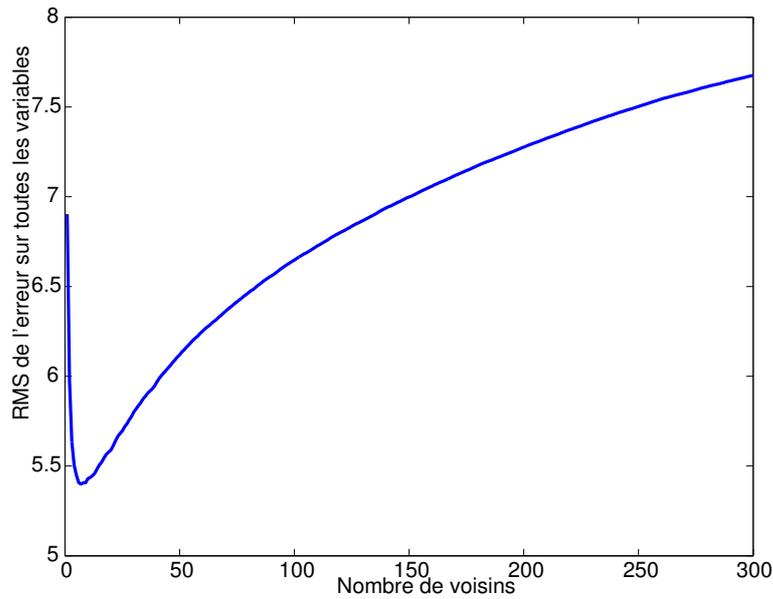


FIGURE 5.10 – Variation de la racine de l'erreur quadratique de restitution, par  $k$ -plus proches voisins, moyennée sur toutes les variables en fonction du nombre de voisins considérés.

On peut alors procéder aux restitutions en elles-mêmes suivant le schéma des entrées et des sorties de la Figure 5.11.

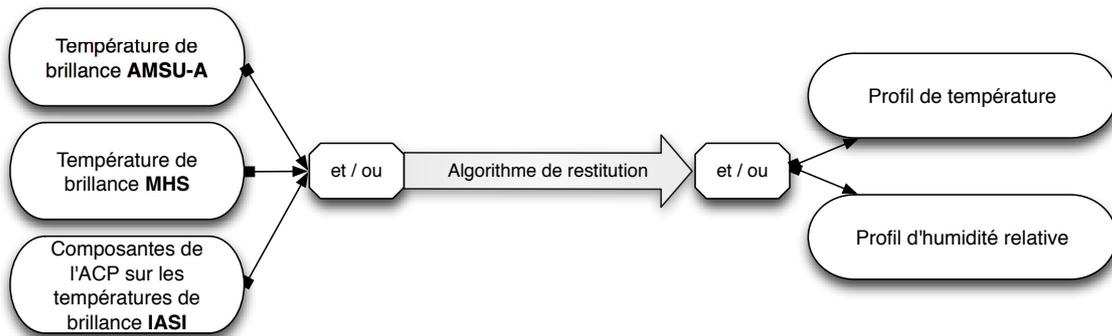


FIGURE 5.11 – Schéma des configurations de restitution.

La Figure 5.12 représente la racine carrée de l'erreur quadratique de restitution par  $k$ -plus proches voisins du profil de température. Pour plus de clarté dans les discussions à

venir, nous appellerons cette erreur par son acronyme anglais : la RMS de l'erreur (Root Mean Square error). Comme nous l'avons dit au préalable, la RMS de l'erreur représente à la fois le biais  $b$  et l'écart-type  $\sigma$  de l'erreur, suivant la relation :  $RMS^2 = b^2 + \sigma^2$ . Nous n'utiliserons ici uniquement cette mesure de l'erreur car nous cherchons à restituer un produit le plus proche possible des variables atmosphériques réelles (minimiser le biais  $b$ ) mais également qui respecte les variations (*i.e.*, cycle diurne...) de ces variables (minimiser l'écart-type  $\sigma$ ).

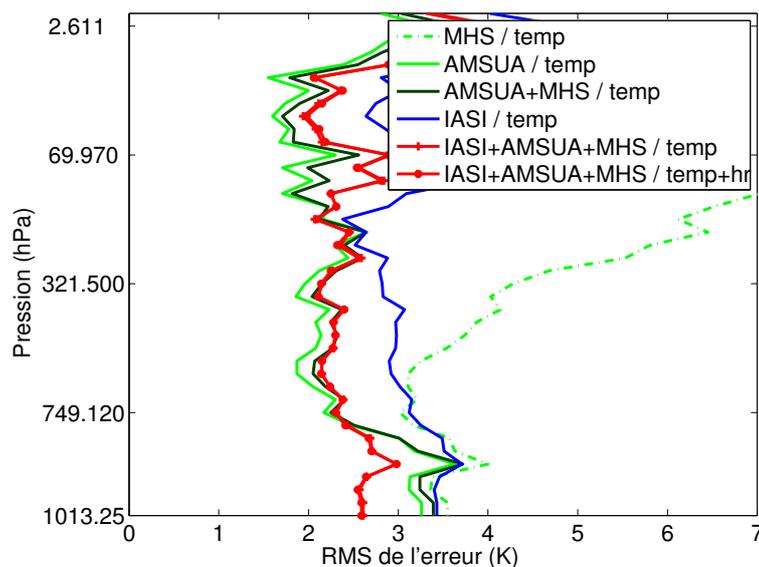


FIGURE 5.12 – RMS de l'erreur de restitution de la température par  $k$ -plus proches voisins en utilisant en entrée : MHS (courbe verte claire en pointillés), AMSU-A (courbe verte claire), AMSU-A et MHS simultanément (courbe verte foncée), IASI (courbe bleue) et les trois instruments simultanément (courbe rouge avec tirets). La courbe rouge avec des points correspond à la configuration où tous les instruments sont utilisés pour restituer la température et la vapeur d'eau simultanément.

Ces RMS d'erreur représentent les différences statistiques entre les profils restitués et les profils atmosphériques réels. Ces restitutions sont effectuées en prenant en compte les observations AMSU-A, MHS, AMSU-A + MHS, IASI et les trois capteurs simultanément. Cette dernière configuration (utilisant les trois capteurs simultanément) est utilisée pour restituer la température uniquement puis la température et la vapeur d'eau simultanément. Les résultats pour ces deux configurations sont identiques. Pour les  $k$ -plus proches voisins, restituer seulement la température, ou la température et la vapeur d'eau simultanément, ne fait pas de différence. Les profils restitués sont ceux associés aux différents plus proches voisins, sans aucune interaction entre les deux variables atmosphériques. Les  $k$ -plus proches voisins ne peuvent utiliser la synergie indirecte car il n'y a pas de contraintes entre les profils atmosphériques.

Si la restitution en utilisant AMSU-A présente des meilleures statistiques que celle uti-

lisant MHS, la combinaison des informations micro-ondes dégrade le résultat. L'ajout des informations de MHS, qui est peu sensible à la température, perturbe la reconnaissance des plus proches voisins. Dans les couches autour de 500 hPa, la RMS de l'erreur en utilisant AMSU-A est inférieure à celle utilisant le capteur hyperspectral IASI. La combinaison des trois capteurs donne de meilleurs résultats proche de la surface, avec une erreur inférieure aux restitutions indépendantes. On peut cependant noter une légère dégradation de la restitution dans le milieu de l'atmosphère, qui s'accroît dans les plus hautes couches. La restitution utilisant IASI, qui présente des statistiques moins bonnes dans le haut de l'atmosphère, vient ici perturber la restitution. Ne pas prendre en compte le contenu en information des différentes mesures, sur lesquelles on effectue la recherche des plus proches voisins, rend compliqué l'utilisation de la synergie.

Les statistiques similaires correspondant au profil de vapeur d'eau sont présentées sur la Figure 5.13. Ici encore, le gain par fusion des informations sur l'erreur de restitution est présent dans les plus basses couches de l'atmosphère. Cependant, ce gain est marginal et MHS seul est meilleur pour certaines couches.

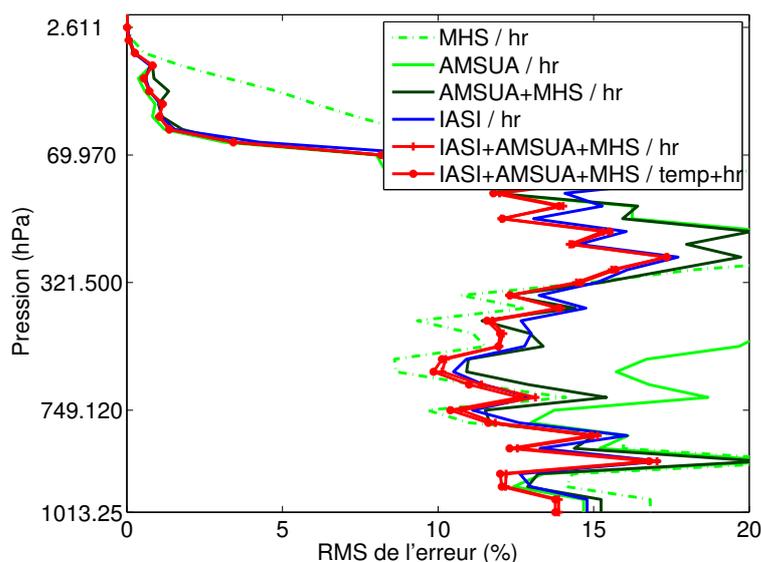


FIGURE 5.13 – Figure identique à la précédente mais pour l'humidité relative.

Pour comprendre ce comportement, il faut garder à l'esprit que, dans la méthode des  $k$ -plus proches voisins, la reconnaissance de forme est effectuée uniquement sur les entrées (*i.e.*, les observations satellites), indépendamment des sorties (*i.e.*, les variables atmosphériques). Lorsque tous les capteurs sont utilisés, le même poids est donné aux canaux de MHS (particulièrement sensibles à la vapeur d'eau), qu'aux observations IASI ou AMSU-A. Les informations de MHS peuvent alors perturber la restitution du profil de température et les autres capteurs feront de même pour la restitution du profil d'humidité relative. Afin de mieux utiliser la synergie, cet algorithme devrait être modifié afin de prendre en compte

les différents capteurs de façon comparable à leur influence sur la variable à restituer. La complexité de la mise en place d'un algorithme des  $k$ -plus proches voisins pondéré par les différentes sensibilités nous a poussés à tester d'autres algorithmes plutôt que de le développer.

### 5.2.3 Restitutions par régression linéaire

#### 5.2.3.1 Méthode

Dans la régression linéaire, chaque sortie de l'algorithme est une combinaison linéaire des différentes entrées plus un biais éventuel. Les paramètres (*i.e.*, les différents poids de chaque entrée et le biais) de la régression sont ajustés de façon à minimiser l'écart entre les variables restituées et les variables réelles (au sens des moindres carrés), grâce à la base d'apprentissage. Cette méthode ne sera pas plus détaillée ici, car elle est relativement simple et classique. C'est une méthode réellement multivariée. Contrairement aux  $k$ -plus proches voisins, seules les informations pertinentes sont prises en compte. Les entrées dénuées d'informations ne pollueront pas la restitution. C'est par contre une méthode linéaire, comme son nom l'indique, qui peut alors être inadaptée pour simuler des relations non-linéaires entre les entrées et les sorties. Il est fréquent d'utiliser cette méthode comme un premier test pour la méthode neuronale utilisée par la suite.

#### 5.2.3.2 Résultats

Les différentes restitutions mises en place ici sont semblables aux schémas utilisés pour les  $k$ -plus proches voisins (voir Schéma 5.11). La base d'apprentissage construite à la Section 5.1.2 (page 124) est utilisée pour calculer les différents paramètres de la régression. Les statistiques de la RMS de l'erreur sont calculées sur la base de test, en comparant les profils restitués aux profils réels. Ici encore, nous n'utilisons pas la base de validation.

Les statistiques de la RMS de l'erreur de restitution du profil de température sont présentées sur la Figure 5.14. Ces courbes correspondent (conformément à ce qui a été fait précédemment) aux restitutions utilisant AMSU-A, MHS, AMSU-A + MHS, IASI et les trois capteurs simultanément. Nous testons la dernière configuration (avec les trois capteurs) pour restituer le profil de température seul, ou les profils de température et d'humidité relative simultanément. Comme on pouvait le prévoir, ces deux configurations donnent des résultats identiques. La restitution linéaire de chaque variable (température ou humidité relative d'une couche donnée) est indépendante des autres. Il ne peut donc pas y avoir de gain ou de perte, à restituer plusieurs variables simultanément. La régression linéaire ne permet pas de mettre en valeur la synergie indirecte.

Sans surprises, MHS est moins informatif que AMSU-A pour le profil de température. IASI permet une amélioration significative des statistiques de restitution dans la basse atmosphère (plus de 700 hPa). AMSU-A est, quant à lui, meilleur pour les hautes couches.

Les informations fournies en entrée d'un modèle linéaire sont additives. La combinaison des différents capteurs met bien en valeur la synergie additive, avec un gain de près de 0,5 K près de la surface. Le reste du profil est systématiquement équivalent voire meilleur que les restitutions indépendantes.

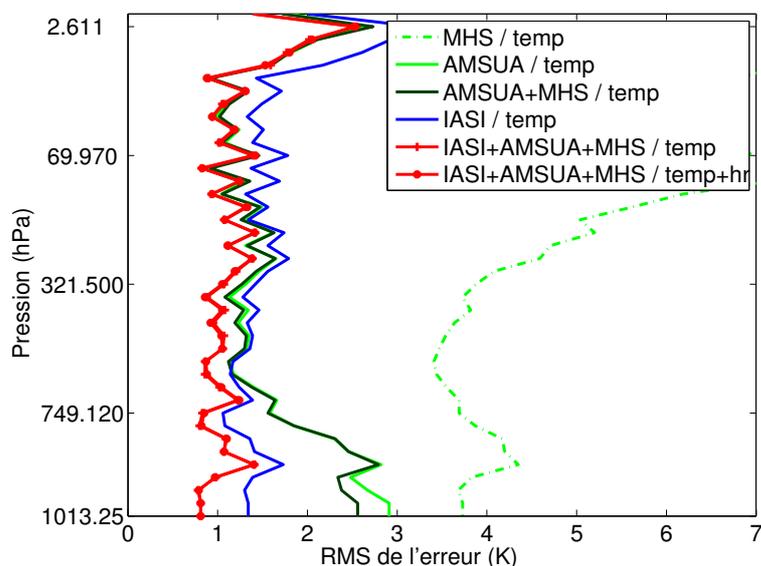


FIGURE 5.14 – RMS de l'erreur de restitution de la température par régression linéaire en utilisant les différentes configurations (les mêmes que précédemment).

La Figure 5.15 présente les statistiques semblables pour la RMS de l'erreur de restitution du profil d'humidité relative. On remarque à nouveau l'impact positif de la synergie. La restitution, en utilisant les capteurs simultanément, est systématiquement meilleure que la meilleure des restitutions indépendantes. L'apport peut être très important, notamment au niveau des couches proches de la surface, où l'erreur passe de 12,5 % à 10 %, soit un gain de 20 % (en erreur relative). Contrairement à ce que l'on obtenait avec les  $k$ -plus proches voisins, la synergie des instruments micro-ondes peut également être mise en valeur. L'erreur de restitution de l'humidité relative, en utilisant AMSU-A et MHS simultanément, est inférieure aux erreurs de restitution en utilisant ces instruments séparément. Le peu d'information que contient MHS sur la température ne permettait pas de mettre en valeur cette synergie sur la restitution de la température. On peut cependant noter que l'apport d'entrées non informatives n'a pas dégradé la restitution. La régression linéaire est donc une méthode plus robuste, car elle permet la prise en compte d'informations multiples sans pertes d'informations.

La synergie indirecte n'est pas présente ici, conformément à ce que nous avons expliqué plus haut sur l'incapacité de la régression linéaire à exploiter les interactions entre ses sorties.

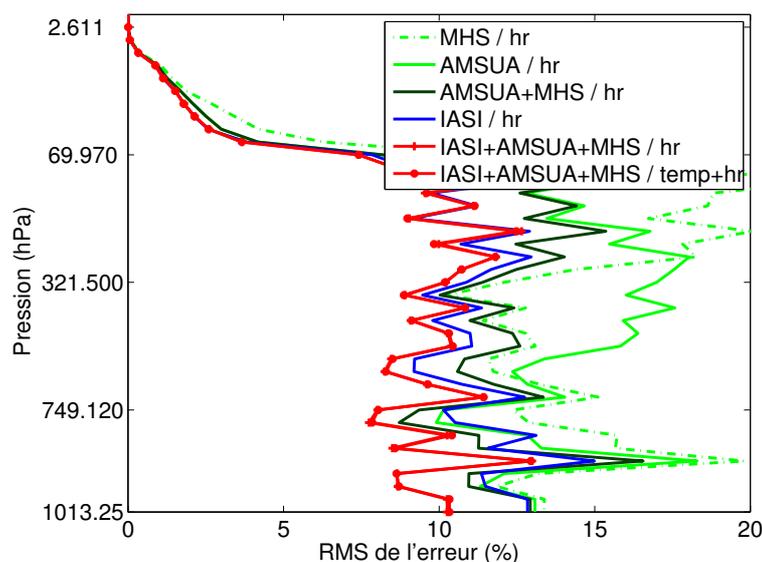


FIGURE 5.15 – RMS de l’erreur de restitution de l’humidité relative par régression linéaire en utilisant les différentes configurations.

## 5.2.4 Restitutions grâce à un réseau de neurones

### 5.2.4.1 Méthode

Les méthodes neuronales sont des algorithmes particulièrement efficaces pour la télédétection. Le perceptron multicouche (Rumelhart et al. 1986) est utilisé ici (semblable à celui utilisé pour l’interpolateur d’émissivité infrarouge à la Section 2.2.1, page 50). Il s’agit d’un modèle non-linéaire. Nous utilisons ici un réseau de neurones à une seule couche cachée. L’utilisation d’une seconde couche cachée introduirait plus de non-linéarité dans le modèle. Nous avons préféré ne pas considérer de deuxième couche cachée pour deux raisons :

- La plupart du temps, le gain n’est pas significatif. Ce que l’on peut gagner en représentation de la complexité d’un modèle peut être perdu en surapprentissage ;
- Le processus d’apprentissage devient bien plus long, il est alors chronophage de tester de multiples combinaisons.

Ceci illustre le principe du rasoir d’Ockham (Pearl 2000) qui énonce, de façon résumée et simplifiée, que les hypothèses valables les plus simples sont les bonnes. Il est inutile de complexifier un modèle qui fonctionne.

Un perceptron multicouche est défini par son nombre d’entrées, de sorties et le nombre de neurones sur la couche cachée qui contrôle la complexité du modèle. Trop de paramètres libres dans le modèle (*i.e.*, trop de neurones cachés) peut entraîner une sur-paramétrisation et donc une dégradation de la capacité à généraliser du réseau (augmentation de la variance de l’erreur de restitution). Au contraire, trop peu de paramètres entraînera une sous-paramétrisation, et donc des erreurs de restitution du modèle plus élevées (augmentation du biais de la restitution).

Le réseau de neurones est étudié pour reproduire le comportement décrit par une base d'échantillons, constituée d'entrées (des observations satellites) et de sorties correspondantes (les variables atmosphériques). Si assez d'échantillons sont fournis, toutes les relations continues, aussi complexes soient-elles, entre les entrées et les sorties peuvent être représentées par un perceptron multicouche (Hornik et al. 1989; Cybenko 1989). La base d'apprentissage présentée à la Section 5.1.2 (page 124) est utilisée pour paramétrer le réseau de neurones. L'algorithme d'apprentissage utilisé est une rétropropagation du gradient, suivant la méthode de Levenberg-Marquardt (voir Section B.2.2, page 215). Cette méthode d'optimisation a déjà prouvé son efficacité pour de tels problèmes. Le critère de qualité à minimiser doit être choisi précautionneusement, surtout lorsque plusieurs types de variables sont restituées simultanément (*i.e.*, la température et la vapeur d'eau).

Les réseaux de neurones sont paramétrés sur la base de données atmosphériques de 9.000 situations. Leur capacité à généraliser les restitutions est mesurée sur la base de validation de 500 situations. L'apprentissage de chaque réseau est stoppé lorsque l'erreur de généralisation cesse de diminuer. Afin d'éviter que les réseaux "apprennent" la base de validation, une troisième base est utilisée pour mesurer les erreurs finales de restitution, de façon plus réaliste. Ces trois estimations d'erreur restent assez proches pour justifier qu'il n'y a pas de surapprentissage ou de surparamétrisation. Les erreurs présentées par la suite sont celles calculées sur la base de test.

#### 5.2.4.2 Résultats

Les statistiques de la RMS d'erreur de restitution du profil de température par des réseaux de neurones sont présentées sur la Figure 5.16, pour les mêmes configurations que précédemment. MHS est ici encore moins informatif sur le profil de température. AMSU-A est meilleur pour les hautes couches de l'atmosphère tandis que IASI est meilleur dans les basses couches. La combinaison des trois instruments améliore la restitution. Près de la surface, la meilleure restitution indépendante (IASI) donne une erreur de 1,2 K alors que la restitution utilisant les trois capteurs donne une erreur de seulement 0,5 K. De la même façon que pour la régression linéaire, un réseau de neurones est capable de combiner plusieurs informations. Plus la quantité d'informations en entrée augmente, plus l'erreur de restitution diminue. Un réseau de neurones est capable de prendre en compte toutes les informations pertinentes disponibles en entrée, même si MHS pourrait "contaminer" la restitution de la température, ce n'est pas le cas ici.

Il est intéressant de noter que la restitution est dégradée lorsque la vapeur d'eau et la température sont restituées simultanément. Le réseau de neurones utilisé ici n'est donc pas capable d'exploiter la synergie indirecte. Ceci peut être dû à plusieurs raisons :

- La complexification du problème, liée à la restitution simultanée de la température et de la vapeur d'eau, a rendu le réseau moins précis. Il faudrait reconfigurer le réseau afin de lui permettre de s'adapter à un problème plus complexe (*i.e.*, augmenter le

- nombre de neurones sur la couche cachée) ;
- Le critère de convergence (*i.e.*, les moindres carrés) augmente l’influence relative du profil de vapeur d’eau par rapport au profil de température, car l’erreur commise est plus importante. Le réseau se focalise alors sur la restitution de la vapeur d’eau ;
  - La restitution du profil de température seul présente déjà une erreur faible, il est alors difficile, voire impossible de diminuer cette erreur.

D’autres applications ont montré la capacité d’un réseau de neurones à exploiter la synergie indirecte. Cela dépend de l’interdépendance des sorties à restituer.

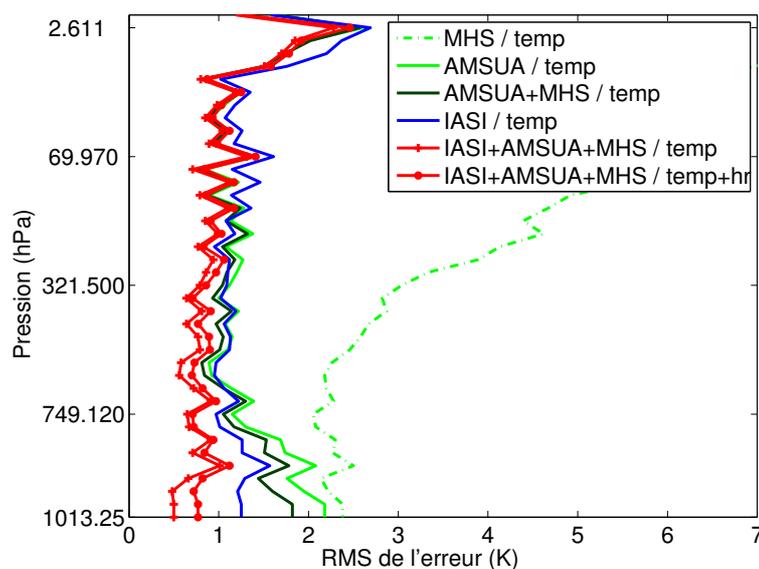


FIGURE 5.16 – RMS de l’erreur de restitution de la température par réseau de neurones en utilisant les différentes configurations.

La Figure 5.17 présente les statistiques de la RMS de l’erreur de restitution du profil d’humidité relative pour les différentes configurations. MHS permet une erreur plus faible que AMSU-A entre 300 et 700 hPa, c’est l’inverse dans le reste de l’atmosphère. Les restitutions utilisant IASI sont systématiquement meilleures que celles utilisant les données micro-ondes. La synergie de ces trois instruments est toujours positive, notamment dans le bas de l’atmosphère.

La synergie indirecte n’est encore une fois pas exploitée par le réseau de neurones. Si la restitution simultanée de la température et de la vapeur d’eau n’améliore pas la restitution de la vapeur d’eau, elle ne la dégrade pas non plus. Ceci tend à prouver que le problème de dégradation de la restitution de la température, vu précédemment, est lié à une mauvaise caractérisation du réseau, puisque l’erreur devrait, au moins, être équivalente. À architecture constante, il faut faire un compromis pour restituer en même temps la température et la vapeur d’eau. S’il n’y a pas assez de neurones, il faut “partager” les différents effets du réseau de neurones. En changeant l’architecture du réseau, il suffit de rassembler les réseaux

restituant la température et la vapeur d'eau pour obtenir un réseau capable de restituer simultanément les deux profils de façon équivalente.

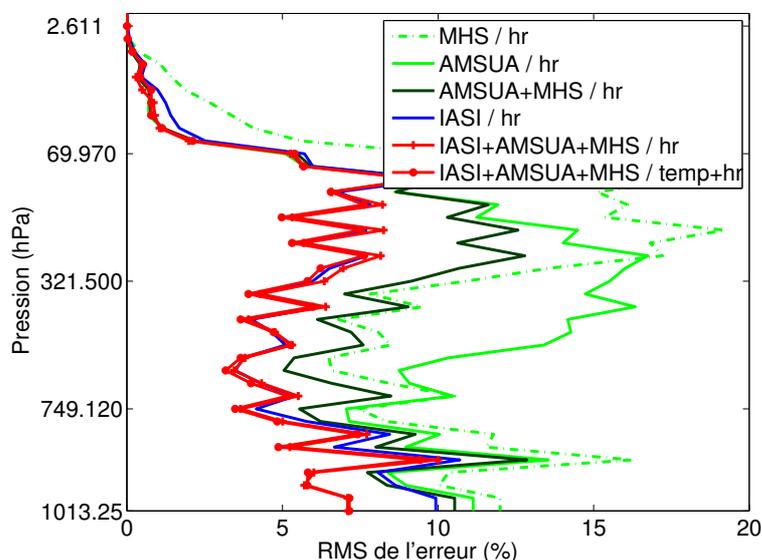


FIGURE 5.17 – RMS de l’erreur de restitution de l’humidité relative par réseau de neurones en utilisant les différentes configurations.

La restitution de la température ou de la vapeur d’eau en utilisant IASI est complexe. L’utilisation de composantes de l’ACP en lieu et place des températures de brillance mélange les différentes informations. La combinaison de ces composantes avec les informations micro-ondes permet de “démêler” les différentes informations et de restituer de façon plus précise la température ou la vapeur d’eau. La fonction de coût incluse dans l’apprentissage est mieux caractérisée et l’apprentissage en est facilité. L’utilisation des trois capteurs permet de régulariser le problème.

### 5.2.5 Comparaison des trois méthodes

Nous comparons les trois méthodes de restitution dans la configuration où tous les capteurs (AMSU-A, MHS et IASI) sont utilisés simultanément pour restituer la température ou la vapeur d’eau. Cette configuration correspond, pour chaque méthode, à celle présentant les RMS d’erreurs de restitution les plus faibles. Les rares fois où ce n’est pas le cas, la différence est très faible (dans le cas des  $k$ -plus proches voisins notamment).

La Figure 5.18 représente la RMS de l’erreur de restitution de la température pour les trois différentes méthodes. On représente également sur la même figure l’écart-type du profil de température. Il correspond à la variabilité naturelle de la température au sein de la base d’apprentissage. On voit ici que même la restitution par  $k$ -plus proches voisins (qui est la moins précise) présente une erreur nettement inférieure à cette variabilité naturelle. Les restitutions par régression linéaire et réseau de neurones sont nettement meilleures. Le

réseau de neurones présente même une erreur systématiquement inférieure à celle avec la régression linéaire. Ce qui est normal puisque les réseaux de neurones sont une extension non-linéaire de la régression linéaire.

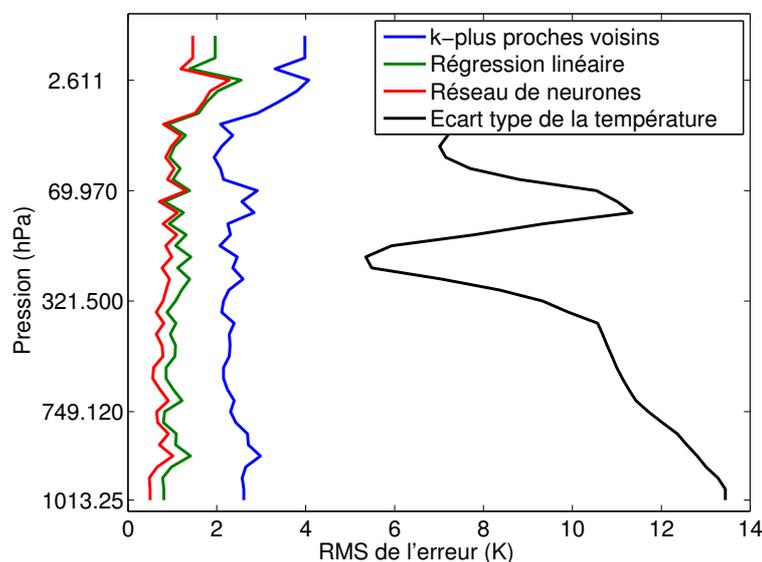


FIGURE 5.18 – RMS de l’erreur de restitution de la température en utilisant AMSU-A, MHS et IASI simultanément par  $k$ -plus proches voisins (courbe bleue), régression linéaire (courbe verte) et réseau de neurones (courbe rouge). L’écart-type de la température est représenté en noir à des fins de comparaisons.

Un comportement similaire peut être remarqué sur les profils de RMS d’erreurs de restitution du profil d’humidité relative présentés sur la Figure 5.19. Les  $k$ -plus proches voisins présentent les résultats les plus mauvais avec plus de 10 % d’erreur. L’erreur de restitution par régression linéaire reste autour de 10 %. Le réseau de neurones se comporte bien mieux, avec des erreurs avoisinant les 5 % tout au long du profil. La relation entre les observations satellites et la vapeur d’eau est plus complexe et moins linéaire que celle avec la température. C’est pourquoi le réseau de neurones est meilleur que les autres méthodes. L’impact de la non-linéarité du réseau de neurones est présent également mais de façon moins significative pour la restitution du profil de température.

Divakarla et al. (2006) ont fait des comparaisons entre la température et la vapeur d’eau restituées par AIRS (Atmospheric InfraRed Sounder, un instrument similaire à IASI), issues de radiosondages et des données issues de modèles de prévision météorologiques. De cette étude, on voit que l’erreur sur le profil de température avoisine le 1 à 1,5 K et celle sur la vapeur d’eau avoisine 20 %. Ces valeurs sont supérieures aux résultats que nous obtenons ici, mais la coïncidence difficile entre les mesures satellites et les radiosondages a tendance à surestimer les erreurs de restitution. De plus, l’utilisation de données satellites simulées que nous faisons ici sous-estime les erreurs de restitution.

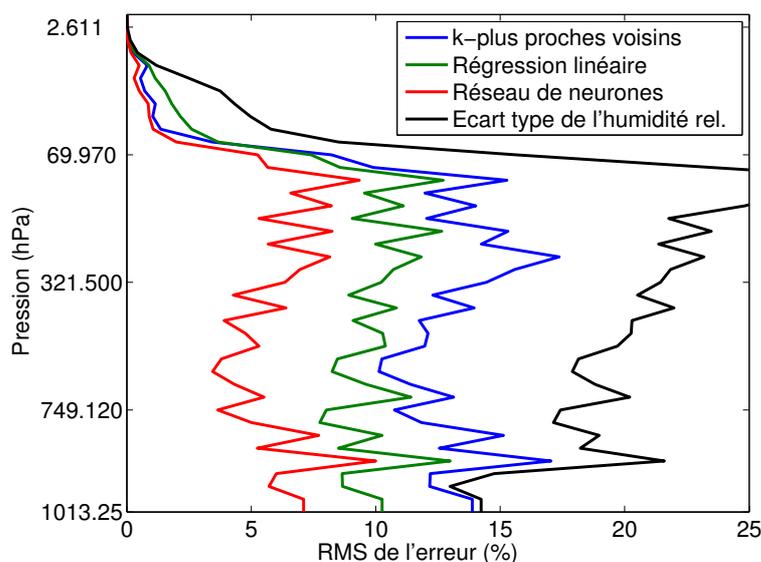


FIGURE 5.19 – RMS de l’erreur de restitution de l’humidité relative en utilisant AMSU-A, MHS et IASI simultanément par  $k$ -plus proches voisins (courbe bleue), régression linéaire (courbe verte) et réseau de neurones (courbe rouge). L’écart-type de l’humidité relative est représenté en noir à des fins de comparaisons.

### 5.2.6 Évaluation de la synergie

Comme expliqué au Chapitre 4 (page 107), le facteur de synergie est défini par le quotient entre la meilleure restitution individuelle et la restitution combinée. Un facteur de 100% (attention il s’agit ici d’un pourcentage d’erreur et non plus des pourcentages d’humidité relative) indique qu’il n’y a pas de synergie. Considérer les différents capteurs simultanément ou indépendamment ne change pas le résultat de la restitution. Un facteur inférieur correspond à un cas où la prise en compte de plus d’informations a dégradé la restitution. Si le facteur est supérieur à 100%, alors il y a de la synergie.

Le facteur de synergie pour la restitution du profil de température, en utilisant les trois méthodes présentées ci-dessus, est représenté sur la Figure 5.20. L’algorithme des  $k$ -plus proches voisins a un facteur de synergie proche de 100%, avec un léger gain en dessous de 800 hPa, mais une dégradation dans la haute atmosphère. Comme nous l’avons indiqué, cet algorithme nécessiterait une amélioration par la pondération appropriée des entrées, afin de mieux exploiter la synergie infrarouge/micro-onde. La régression linéaire présente de la synergie tout au long du profil et plus particulièrement proche de la surface. L’impact peut être vraiment important avec un gain de plus de 50% dans les basses couches de l’atmosphère et autour de 10% au milieu de la troposphère. Le facteur de synergie pour les réseaux de neurones est encore supérieur, les erreurs de restitutions peuvent être réduites d’un facteur de 2,5 pour les couches proches de la surface. En effet, un facteur de synergie de 250% signifie que le rapport entre l’erreur de restitution indépendante et celle simultanée est de

2,5. Tout le long du profil, le facteur de synergie reste très élevé.

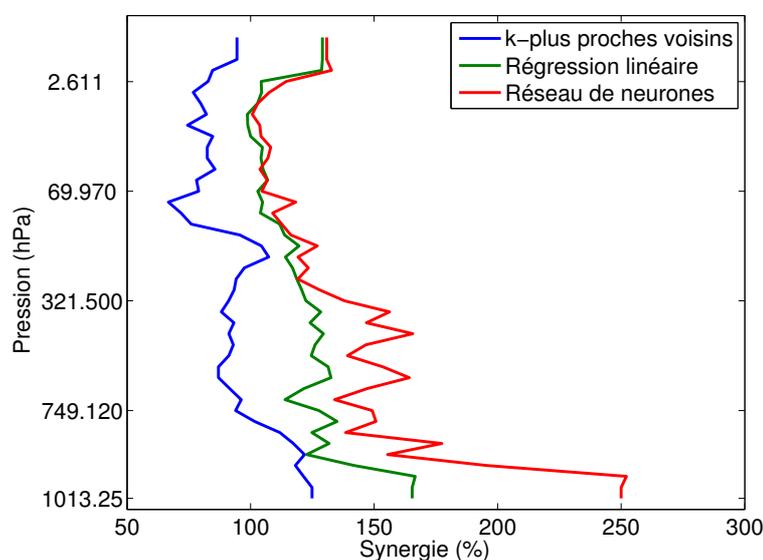


FIGURE 5.20 – Facteur de synergie pour la restitution de la température par  $k$ -plus proches voisins (courbe bleue), régression linéaire (courbe verte) et réseau de neurones (courbe rouge).

La Figure 5.21 montre les statistiques du facteur de synergie pour la restitution du profil d'humidité relative en utilisant les  $k$ -plus proches voisins, la régression linéaire et un réseau de neurones. Ici encore, la méthode des  $k$ -plus proches voisins ne permet pas une fusion optimale des informations. L'impact de l'utilisation des trois capteurs simultanément est systématiquement négatif. Le pic fortement positif dans les hautes couches de l'atmosphère n'a que peu de signification physique. Du fait de la faible variabilité de l'humidité relative dans ces couches, un tel gain peut être un simple artifice statistique. La régression linéaire présente un facteur de synergie proche de 100% (*i.e.*, il n'y a presque pas de synergie), sauf pour les niveaux en dessous de 700 hPa où l'impact est positif. Les réseaux de neurones profitent également de la synergie. Si celle-ci est peu visible au milieu de l'atmosphère, le gain proche de la surface est important.

Le facteur de synergie pour la restitution de la vapeur d'eau est supérieur pour la régression linéaire par rapport à celui des réseaux de neurones dans le milieu de l'atmosphère. Il faut cependant bien garder à l'esprit les statistiques de RMS de l'erreur de restitution présentée précédemment. Les restitutions à l'aide de réseaux de neurones présentent des erreurs systématiquement inférieures à celles obtenues avec des régressions linéaires. L'erreur plus faible des réseaux de neurones explique le gain plus faible lié à l'utilisation de la synergie. De plus, si le gain par synergie est plus faible, l'objectif reste tout de même de restituer les profils de vapeur d'eau et de température. Une méthode présentant des erreurs systématiquement inférieures aux autres est donc préférée.

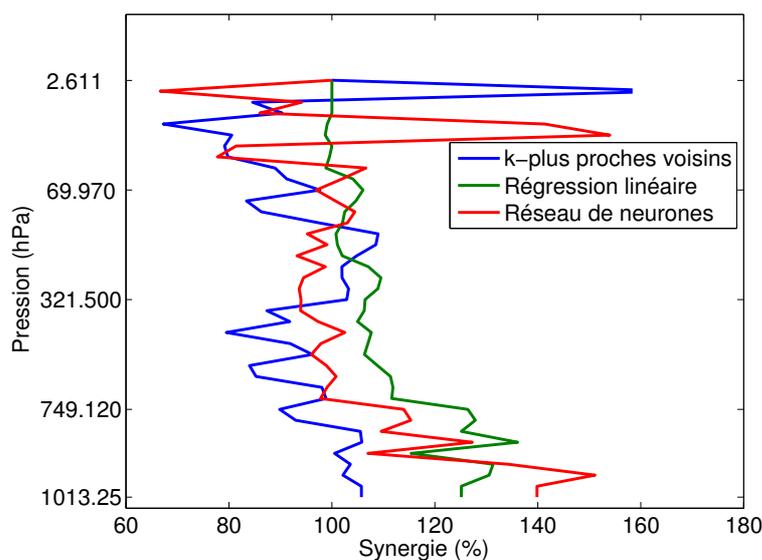


FIGURE 5.21 – Facteur de synergie pour la restitution de l’humidité relative par  $k$ -plus proches voisins (courbe bleue), régression linéaire (courbe verte) et réseau de neurones (courbe rouge).

### 5.3 Conclusion

Nous avons réussi, dans ce chapitre, à mettre en application directe les principes de synergie présentés au Chapitre 4 (page 107). Nous avons mis au point une méthode statistique à même d’exploiter les synergies potentielles entre différents instruments d’un satellite. L’importante synergie entre l’infrarouge et le micro-onde permet d’améliorer la restitution des profils de température et de vapeur d’eau, même en ciel clair. Nous avons également démontré la capacité des réseaux de neurones à exploiter cette synergie de façon multivariée et non-linéaire.

La synergie indirecte n’a pas pu être mise en valeur dans cette étude. Elle reste cependant primordiale. Au niveau des NWP, elle est prise en compte dans la mesure où les matrices de corrélations entre les variables atmosphériques servent à contraindre les solutions des modèles. La méthode développée ici est très générale et peut être appliquée à beaucoup de situations, avec d’autres capteurs, d’autres variables atmosphériques...

Après avoir mis au point une méthode de restitution des paramètres de surface et une méthode de restitution des profils atmosphériques, il nous reste à combiner les deux pour restituer des profils atmosphériques au-dessus des continents.

# ÉCHANTILLONNAGE PAR ENTROPIE



## Sommaire

<b>6.1</b>	<b>L'entropie au sens de Shannon</b>	<b>157</b>
6.1.1	Historique et définition	158
6.1.2	Exemple de calcul d'une entropie	159
<b>6.2</b>	<b>Base de données atmosphériques initiale</b>	<b>160</b>
6.2.1	Les produits L2 IASI d'EUMETSAT	160
6.2.2	Niveaux de pression	162
<b>6.3</b>	<b>Échantillonnage</b>	<b>164</b>
6.3.1	Méthode d'échantillonnage : l'entropie	165
6.3.1.1	Discrétisation	166
6.3.1.2	Variabilité de la vapeur d'eau	166
6.3.1.3	Pondération de chaque variable	167
6.3.1.4	Optimisation du calcul	169
6.3.1.5	Convergence de l'algorithme d'échantillonnage	170
6.3.2	Apport de l'entropie pour l'échantillonnage multivarié	171
6.3.2.1	Exemple d'échantillonnage sur une seule variable	171
6.3.2.2	Prise en compte de plusieurs variables grâce à l'entropie	174
<b>6.4</b>	<b>Résultats</b>	<b>176</b>
6.4.1	Représentation des différentes variables	176
6.4.1.1	Base de données complète	176
6.4.1.2	Base de données échantillonnée	177
6.4.2	Variabilité spatiale	179
<b>6.5</b>	<b>Conclusion</b>	<b>180</b>

La télédétection satellite pour l'observation de la Terre génère par essence une grande quantité de données. De nombreux algorithmes de traitement de ces observations satellites requièrent l'utilisation d'une base de données représentative de ce qui peut exister dans l'atmosphère. Certains algorithmes peuvent utiliser la totalité des données à disposition. Pour

certaines applications, comme l'apprentissage de méthodes statistiques ou l'étude de certains phénomènes locaux et la mise en valeur de régimes spéciaux, seule une faible quantité de données peut être utilisée.

De plus, certains instruments comme IASI (voir Section 1.3.1, page 29) génèrent une telle quantité de données qu'il est nécessaire de réduire au maximum la taille des bases de données pour minimiser les temps de calcul. Il faut alors mettre en place des méthodes d'échantillonnage des données.

La contrainte réside dans la grande quantité de variables à échantillonner et la nécessité de garder certaines de leurs caractéristiques attenantes (extrema, variabilité...). Il faut mettre en place des méthodes multivariées qui la plupart du temps dépendent de l'utilisation qui doit être faite de la base de données (Aires and Prigent 2007). La méthode doit permettre de faire un compromis entre le nombre restreint d'observations qui peuvent être utilisées dans la pratique et une bonne description des variables atmosphériques et de surface.

Nous allons introduire dans ce chapitre une méthode d'échantillonnage originale, multivariée et capable d'agir sur des variables inhomogènes. Il peut être complexe de comparer au cours d'un échantillonnage des températures à des valeurs d'humidité absolue. Même en ne considérant que l'humidité, la différence de variabilité entre les basses couches de l'atmosphère et de la stratosphère complique grandement l'échantillonnage. La méthode mise en place doit pouvoir s'affranchir de ce problème.

La suite de l'étude porte sur la synergie des capteurs infrarouges (IASI) et micro-ondes (AMSU-A et MHS) de la plate-forme MetOp en prenant en compte les caractéristiques de surface au-dessus des continents. Elle nécessite la construction d'une base d'apprentissage robuste, complète et adaptée aux réseaux de neurones qui seront utilisés.

Cette problématique est relativement récurrente dans le cadre de la restitution de variables atmosphériques ou de surface par des méthodes statistiques. De telles méthodes reposent en grande partie sur la qualité de la base d'apprentissage qui est utilisée. Ces bases de données peuvent être utiles, tant pour l'apprentissage de régressions linéaires ou de réseaux de neurones (Aires et al. 2001), que pour des régressions bayésiennes (Kummerow et al. 2001).

Une telle base de données doit être constituée de températures de brillance correspondant aux sondeurs considérés, associées aux profils atmosphériques correspondants. Il existe deux grandes méthodes pour construire une base d'apprentissage :

- **Base "empirique"** : Il est possible de colocaliser des mesures atmosphériques *in situ* et des observations satellites. Cette méthode est complexe car il faut prendre en considération les possibles décalages temporels et spatiaux entre les deux mesures, ainsi que les différentes erreurs mesures ;
- **Base "simulée"** : Une autre méthode possible est de constituer une base de données atmosphériques et d'utiliser ensuite un modèle de transfert radiatif pour y associer

les observations satellites dites “simulées”. Cette méthode paraît plus simple, mais les algorithmes issus de telles bases peuvent ne pas avoir des résultats similaires lorsqu’ils sont appliqués à des données satellites réelles. Il peut y avoir des erreurs systématiques liées au code de transfert radiatif ou à la calibration instrumentale (Aires et al. 2010).

Le développement des codes de transfert radiatif permet d’obtenir des températures de brillance assez proches de la réalité pour préférer opter pour cette deuxième solution. Il nous faut alors construire des bases de données atmosphériques bien représentatives sur lesquelles nous calculerons ensuite le transfert radiatif.

Des bases de données ont déjà été créées. On peut citer par exemple la base de données TIGR (Thermodynamical Initial Guess Retrieval) (Chedin et al. 1985; Escobar 1993; Chevallier et al. 1998). Il s’agit d’une base de données de profils atmosphériques issus de radio-sondages. Cette base est beaucoup utilisée et certains paramètres ont été améliorés pour des utilisations spécifiques comme la variabilité de la température de surface (Aires et al. 2002a) ou avec des re-analyses de l’ECMWF (Chevallier et al. 2000). Dans la continuité de ce travail, une autre base a été développée par l’ECMWF (Chevallier et al. 2006). Cette base a été conçue pour échantillonner, à partir d’une grande base de données atmosphériques, une base plus limitée de situations représentatives, uniformisant la distribution de chaque variable indépendamment. Suivant la variable que l’on souhaite considérer (température, vapeur d’eau, ozone, eau liquide, glace), une base de données différente doit être utilisée.

L’échantillonnage multivarié de variables inhomogènes a toujours posé de grands problèmes. La solution choisie par l’ECMWF (construire différentes bases de données suivant la variable échantillonnée) permet théoriquement de retrouver des distributions quasi-uniformes mais sur une seule variable à la fois. Ceci ne nous semble pas optimal pour représenter la variabilité des situations atmosphériques qui est, par nature, multivariée. L’échantillonnage par variable donne des bases trop ciblées et complexes à réutiliser. Les anomalies statistiques présentes dans la base de données se répercuteront ensuite sur l’algorithme de restitution et donc sur les profils restitués. Nous avons donc décidé de mettre en place notre propre algorithme d’échantillonnage adapté à nos besoins. Nous cherchons ici à mettre au point une méthode pour extraire une base de données multivariée qui échantillonne sur toutes les variables simultanément. Ce chapitre a fait l’objet d’une publication : Paul and Aires (2013a).

## 6.1 L’entropie au sens de Shannon

L’objectif ici est de réduire le nombre d’échantillons d’une base de données. Nous avons déjà présenté au préalable la méthode des k-moyennes qui permet d’échantillonner une base de données (voir Section 5.1.2, page 124, et Aires and Prigent (2007)). Cette approche effectue un échantillonnage dit statistique car les distributions dans la base échantillonnée seront

les mêmes que dans la base d'origine, représentatives de ce que l'on retrouve dans la nature. Cette approche est parfaite pour certaines applications. Pour l'apprentissage de méthodes d'inversion, il est plus optimal d'effectuer un échantillonnage dit uniforme. L'apprentissage d'un réseau de neurones nécessite une base bien construite afin de pouvoir être efficace (voir Annexe B.2.2, page 215). Notre but est de réussir à créer une base multivariée dont chaque variable sera bien représentée, notamment dans ses valeurs extrêmes. Si ces dernières ne sont pas présentes dans la base, l'apprentissage ne pourra les prendre en compte et les résultats seront erronés pour ce type de situation. Dans cette optique, un travail conséquent a été effectué afin de mettre au point une méthode innovante d'échantillonnage capable de travailler sur des variables inhomogènes.

### 6.1.1 Historique et définition

Nous avons décidé de nous servir ici du concept d'entropie. Ce concept a été introduit en théorie de l'information par Claude Shannon (1916-2001) (Shannon et al. 1949; Shannon 1951). Pour lui, l'entropie d'une variable correspond à la quantité d'informations qu'il faut connaître pour pouvoir la déterminer de façon certaine. Par exemple, une variable fixe aura une entropie nulle, car il suffit de connaître la valeur que prend la variable pour la déterminer. À l'inverse, une variable aléatoire aura une entropie supérieure car connaître sa moyenne ne suffit pas à la déterminer complètement. Claude Shannon n'a pas utilisé ce terme d'entropie. C'est John Von Neumann qui l'a nommée ainsi. L'histoire veut que ce dernier ait choisi ce terme pour sa ressemblance à la notion d'entropie en physique statistique mais également car ce terme était suffisamment mal compris pour pouvoir triompher dans tout débat. Le lien avec l'entropie macroscopique et microscopique fut démontré plus tard. L'entropie de Shannon est devenue un concept mathématique de base, en particulier pour la théorie de l'information. Elle est utilisée aujourd'hui pour mesurer, par exemple, le nombre de bit minimal en lequel on peut compresser un fichier sans pertes.

Soit  $X$  une variable aléatoire pouvant prendre les valeurs  $(x_1, x_2, \dots, x_n)$ <sup>1</sup> avec respectivement une probabilité  $(p_1, p_2, \dots, p_n)$ . L'entropie de la variable  $X$  au sens de Shannon est alors :

$$E(X) = \sum_{i=1}^n \left( p_i \times \log \left( \frac{1}{p_i} \right) \right) \quad (6.1)$$

L'entropie de cette variable est maximale si  $\forall i \in \{1 \dots n\}, p_i = \frac{1}{n}$ . Dans ce cas,  $X$  prend chacune des valeurs avec la même probabilité.

Pour une variable à plusieurs dimensions  $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ , du fait que l'on manipule des probabilités, on peut comparer les entropies de  $X_1$  et  $X_2$ . Le calcul de l'entropie de  $X$  néces-

1. On note ici que l'on considère une variable  $X$  discrète. Une définition de l'entropie d'une variable continue est possible en  $\int p \times \log \left( \frac{1}{p} \right) dp$ . Dans la pratique, on utilise l'entropie discrète. Il est aussi important de noter que l'entropie est dépendante de la discrétisation choisie.

site alors de discrétiser un espace multidimensionnel. Le nombre de situations à considérer est alors une fonction exponentielle de la dimension de  $X$ . Dans un espace de dimension  $m$ , si chaque variable peut prendre  $p$  valeurs, on doit alors considérer un espace de dimension  $m^p$ . Techniquement cela supposerait de calculer  $m^p$  probabilités, ce qui n'est pas faisable si  $m > 10$ . L'espace que l'on va considérer dans ce chapitre (les variables géophysiques) est constitué de 68 variables (43 niveaux de pression en température, 20 en vapeur d'eau, 4 couches d'ozone et la température de surface). Afin de remédier à ce problème, nous allons considérer que l'entropie d'une variable à plusieurs dimensions correspond à la somme des entropies des sous-variables :  $E(X) = E(X_1) + E(X_2)$ . Ceci revient à sous estimer l'entropie de la variable  $X$ . Ainsi, calculer l'entropie d'une variable multivariée comme la somme des entropies de ses variables conduit à une sous estimation globale de l'entropie de la variable. L'utilisation de cette approximation n'est pas optimale, car les corrélations entre les variables ne sont pas prises en compte. Cependant, la maximisation des entropies indépendantes permet de maximiser l'entropie globale. La suite de cette étude nous prouvera que cette approximation est satisfaisante et conduit à de bons résultats.

L'entropie est donc une mesure dans l'espace des probabilités, ce qui permet de surmonter aisément l'inhomogénéité des variables. On peut facilement comparer ou additionner les entropies de deux variables différentes.

### 6.1.2 Exemple de calcul d'une entropie

Nous allons utiliser l'entropie en remplaçant la probabilité des variables aléatoires par des probabilités d'appartenance à un intervalle donné. En télédétection spatiale, les variables sont continues et non discrètes (température, vapeur d'eau...). En considérant le minimum et le maximum d'une variable, on peut découper son intervalle d'existence en intervalles plus restreints et ainsi les discrétiser. Pour une base de données, on peut calculer l'entropie d'une variable de la même façon que l'entropie de Shannon, en utilisant la probabilité de la variable d'appartenir à chaque sous-intervalle prédéfini.

#### Exemple

Considérons trois bases de températures de surface composées de :

- **Cas 1** : {279; 287; 299; 294; 300; 299; 282; 277; 282; 275} °K.
- **Cas 2** : {275; 277; 276; 279; 278; 297; 295; 299; 296; 300} °K.
- **Cas 3** : {282; 287; 296; 282; 276; 289; 297; 275; 294; 292} °K.

La variable varie entre 275 et 300° K. On découpe cet intervalle de variation en 5 sous-intervalles : [275; 280[, [280; 285[, [285; 290[, [290; 295[ et [295; 300]. Par rapport aux trois bases de données considérées, on peut associer à chaque intervalle une probabilité qui correspond au prorata du nombre de situations dans cet intervalle sur le nombre total de situations. On a alors :

Cas	[275; 280[	[280; 285[	[285; 290[	[290; 295[	[295; 300]	Entropie
1	0,3	0,2	0,1	0,1	0,3	0,65
2	0,5	0	0	0	0,5	0,35
3	0,2	0,2	0,2	0,2	0,2	0,7

Pour calculer l'entropie de la température de surface dans cette base de données, il suffit de sommer les  $p \times \log\left(\frac{1}{p}\right)$  de chaque intervalle. Ceci nous donne une entropie de 0,65 (cas numéro 1 présenté dans le tableau). Si on change la base de données (cas numéro 2 présenté dans le tableau), on a uniquement des cas extrêmes, équitablement répartis entre le maximum et le minimum. On obtient une entropie de 0,35. Enfin, dans le cas 3 du tableau, on a une base de données où les variables sont présentes en quantités égales dans chaque intervalle, ce qui donne une entropie de 0,7. Plus l'histogramme de répartition de la variable est lisse, plus la variable est équirépartie dans les différents sous-intervalles considérés, et plus l'entropie est élevée.

Il est important de noter ici que l'entropie est dépendante de la discrétisation choisie. Plus celle-ci est fine (*i.e.*, les sous-intervalles sont fins et nombreux), moins son influence est importante. Néanmoins, il faut être vigilant et choisir une discrétisation adaptée à la variable considérée.

L'entropie est donc bien une mesure de la variabilité d'une variable au sein d'une base de données. Plus l'entropie est élevée, mieux la variable est répartie uniformément dans sa gamme de variabilité. L'échantillonnage par entropie consiste à extraire un certain nombre de situations d'une base de données, en cherchant à maximiser l'entropie de la nouvelle base. Cette mesure de l'entropie ne dépend plus de la valeur des variables en elles-mêmes. On peut alors facilement additionner et comparer les entropies de plusieurs variables sans soucis d'homogénéité.

## 6.2 Base de données atmosphériques initiale

On décrit ici la base de données atmosphériques que l'on utilise. C'est une base de variables atmosphériques issues de restitutions. Le nombre de profils disponibles est trop élevé, c'est donc cette base que nous allons échantillonner par la suite.

### 6.2.1 Les produits L2 IASI d'EUMETSAT

Dans cette partie, nous utiliserons les données IASI L2 d'EUMETSAT. Chaque situation comprend :

- Le profil de température ;
- Le profil de vapeur d'eau ;
- La quantité totale d'ozone par couches ;
- La température de surface ;

- L'émissivité de surface ;
- La couverture nuageuse ;
- La pression au sommet du nuage ;
- La température au sommet du nuage ;
- La phase nuageuse ;
- La quantité totale de protoxyde d'azote (N<sub>2</sub>O) ;
- La quantité totale de monoxyde de carbone (CO) ;
- La quantité totale de méthane (CH<sub>4</sub>) ;
- La quantité totale de dioxyde de carbone (CO<sub>2</sub>) ;
- La matrice de covariance d'erreur ;
- Différents flags.

La résolution spatiale de chaque situation est de 12 km (restitutions à partir de IASI, donc à la même résolution, voir Section 1.3, page 25). Ces produits L2 sont générés par le EPS (EUMETSAT Polar System) Core Ground Segment (*i.e.*, le segment sol d'EUMETSAT pour les satellites polaires) à partir des observations de IASI. Ils sont restitués en utilisant les radiances L1C de IASI mais également les instruments ATOVS (AMSU-A, AVHRR, MHS), les données L2 de ces instruments et des prédictions météorologiques (EUMETSAT Document 2012).

La restitution se fait en trois étapes. La première est une restitution des variables des nuages, plus particulièrement de la couverture nuageuse et de la hauteur du sommet du nuage. Les situations sont ensuite séparées en deux catégories en fonction de leur nébulosité. La deuxième étape est une restitution basée soit sur une régression à partir de composantes d'ACP soit sur un réseau de neurones. Si le résultat obtenu n'appartient pas à une gamme de résultats considérés comme raisonnables (des limites supérieures et inférieures sont fixées), il est remplacé par une valeur issue d'une climatologie. La troisième étape est une amélioration itérative de la restitution. À l'aide de FRTM (Fast Radiative Transfer Model, un modèle de transfert radiatif semblable à RTTOV, voir Section 1.4, page 33), une simulation du transfert radiatif est effectuée afin de calculer les températures de brillance de IASI ainsi que les Jacobiens associés aux variables précédemment restituées. Les nouvelles variables atmosphériques sont ensuite calculées de façon itérative à l'aide d'une fonction de coût sur la différence entre les températures de brillance simulées et réelles. Cela permet d'obtenir des profils atmosphériques indépendant des modèles de prévision météorologique, qui peuvent donc être comparés à ces derniers.

Dans cette étude, nous nous intéressons uniquement aux situations restituées en ciel clair (soit environ 12% de toutes les situations). Nous nous intéressons uniquement aux variables suivantes :

- Le profil de température ;
- Le profil de vapeur d'eau ;
- La quantité totale d'ozone par couches ;

– La température de surface.

La quantité d’ozone est disponible en quantités intégrées par couches. Quatre valeurs sont disponibles, elles représentent la quantité totale d’ozone intégrée entre la surface et les niveaux de pression : 0,005 ; 132,49 ; 222,94 et 478,54 hPa. Ces quatre valeurs sont utilisées tel quel dans la première partie de l’étude (pour l’échantillonnage), mais elles seront réparties sur les 82 à 89 niveaux de pressions, par la suite, pour servir aux calculs de transfert radiatif avec RTTOV. Cette répartition est faite en utilisant une interpolation verticale, mais de manière à respecter la forme globale du profil d’ozone atmosphérique avec un maximum dans la partie basse de la stratosphère.

Une année de données IASI correspond à plus d’une dizaine de millions de situations. Effectuer un calcul de transfert radiatif pour simuler les 8461 canaux de IASI pour toutes ces situations serait beaucoup trop long. De plus, un apprentissage d’un algorithme de restitution ne peut pas utiliser autant de situations. Il faut donc échantillonner cette base.

### 6.2.2 Niveaux de pression

La base de données considérée est constituée de produits L2 IASI d’EUMETSAT en ciel clair pour toute l’année 2010. Suivant l’altitude de la surface et les conditions météorologiques, la pression de surface varie. Les profils n’ont donc pas toujours le même nombre de niveaux à pression fixe. Les niveaux à pression plus élevée peuvent ne pas exister dans certaines situations. On distingue les différents cas de figures en huit grandes catégories suivant leur nombre de niveaux de pression : de 82 à 89 niveaux. Les niveaux qui ne sont pas représentés sont les plus élevés en pression. Le tableau ci-dessous présente les 89 niveaux de pression.

Pression (hPa)	0,0050	0,0082	0,0135	0,0223	0,0368	0,0607	0,0920	0,1305
Niveaux	1	2	3	4	5	6	7	8
Pression (hPa)	0,1703	0,2222	0,2900	0,3602	0,4473	0,5588	0,6981	0,8722
Niveaux	9	10	11	12	13	14	15	16
Pression (hPa)	1,0896	1,3612	1,6600	2,0613	2,5125	3,1095	3,7839	4,6652
Niveaux	17	18	19	20	21	22	23	24
Pression (hPa)	5,6618	6,9500	8,4895	10,3700	12,3927	14,8100	17,3817	20,4000
Niveaux	25	26	27	28	29	30	31	32
Pression (hPa)	23,5819	27,2600	31,1127	35,5100	40,1030	45,2900	50,6883	56,7300
Niveaux	33	34	35	36	37	38	39	40
Pression (hPa)	63,0031	69,9700	77,2013	85,1800	93,2342	102,0500	111,5983	122,0400
Niveaux	41	42	43	44	45	46	47	48

## 6.2. BASE DE DONNÉES ATMOSPHÉRIQUES INITIALE

Pression (hPa)	132,4924	143,8400	155,4281	167,9500	180,6731	194,3600	208,1601	222,9400
Niveaux	49	50	51	52	53	54	55	56
Pression (hPa)	237,8279	253,7100	269,6541	286,6001	303,5489	321,4999	339,3921	358,2800
Niveaux	57	58	59	60	61	62	63	64
Pression (hPa)	377,0533	396,8100	416,3966	436,9500	457,2725	478,5399	499,5392	521,4601
Niveaux	65	66	67	68	69	70	71	72
Pression (hPa)	543,0530	565,5400	587,6382	610,6000	638,6005	667,7082	696,9703	727,4356
Niveaux	73	74	75	76	77	78	79	80
Pression (hPa)	759,1557	792,1840	826,5760	862,3900	899,6864	938,5284	978,9818	1007,1151
Niveaux	81	82	83 (opt)	84 (opt)	85 (opt)	86 (opt)	87 (opt)	88 (opt)
Pression (hPa)	1021,115	1050						
Niveau	89 (opt)	90 (non)						

Les 82 premier niveaux (haut de l'atmosphère) sont présents dans presque tous les profils (exceptés les plus hauts sommets). Les niveaux 83 à 89 sont les niveaux les plus proches de la surface. Ces derniers peuvent ne pas exister sur un profil si la pression de surface est inférieure à leur pression de référence, ils sont notés (opt) pour optionnels. Le nombre de profils dont la pression de surface est supérieure à 1050 hPa est très faible, c'est pourquoi on ne le prend pas en considération (noté non sur le tableau). On ne les prend donc pas en compte. On dispose alors de 8 bases de données (pour les 8 nombres de niveaux de pression) complètes de situations de l'année 2010<sup>2</sup>.

Ces situations sont représentées sur la Figure 6.1. Sur chacune de ces cartes, la température de surface de chaque situation est représentée en couleur. Chaque carte correspond à un nombre de niveaux de pression dans les profils et donc à une gamme de pressions de surface donnée. De gauche à droite et de haut en bas, il y a de 82 à 89 niveaux de pression. On remarque que les profils constitués de moins de niveaux (pression de surface plus élevée) correspondent à des situations en altitude (autour des hauts sommets comme l'Himalaya). Ces différentes bases de données couvrent donc bien l'ensemble du globe. On peut noter que les hauts sommets ne sont pas contenus dans ces bases de données car la pression de surface est trop faible. Il y aurait donc dans ce cas moins de 82 niveaux de pression. Par souci de clarté et également par souci d'exactitude dans les résultats (RTTOV manque de précision sur les hautes altitudes), ces situations ne sont pas prises en compte. L'altitude élevée de l'Antarctique explique qu'il n'y ait pas de situations dans cette zone. Il faudra après l'échantillonnage vérifier la variabilité spatiale des différentes bases échantillonnées.

2. Cette séparation en 8 bases distinctes a été choisie afin de disposer de niveaux de pression fixes. Dans la pratique, une restitution préalable de la pression de surface est nécessaire pour déterminer le nombre de niveaux de pression fixe.

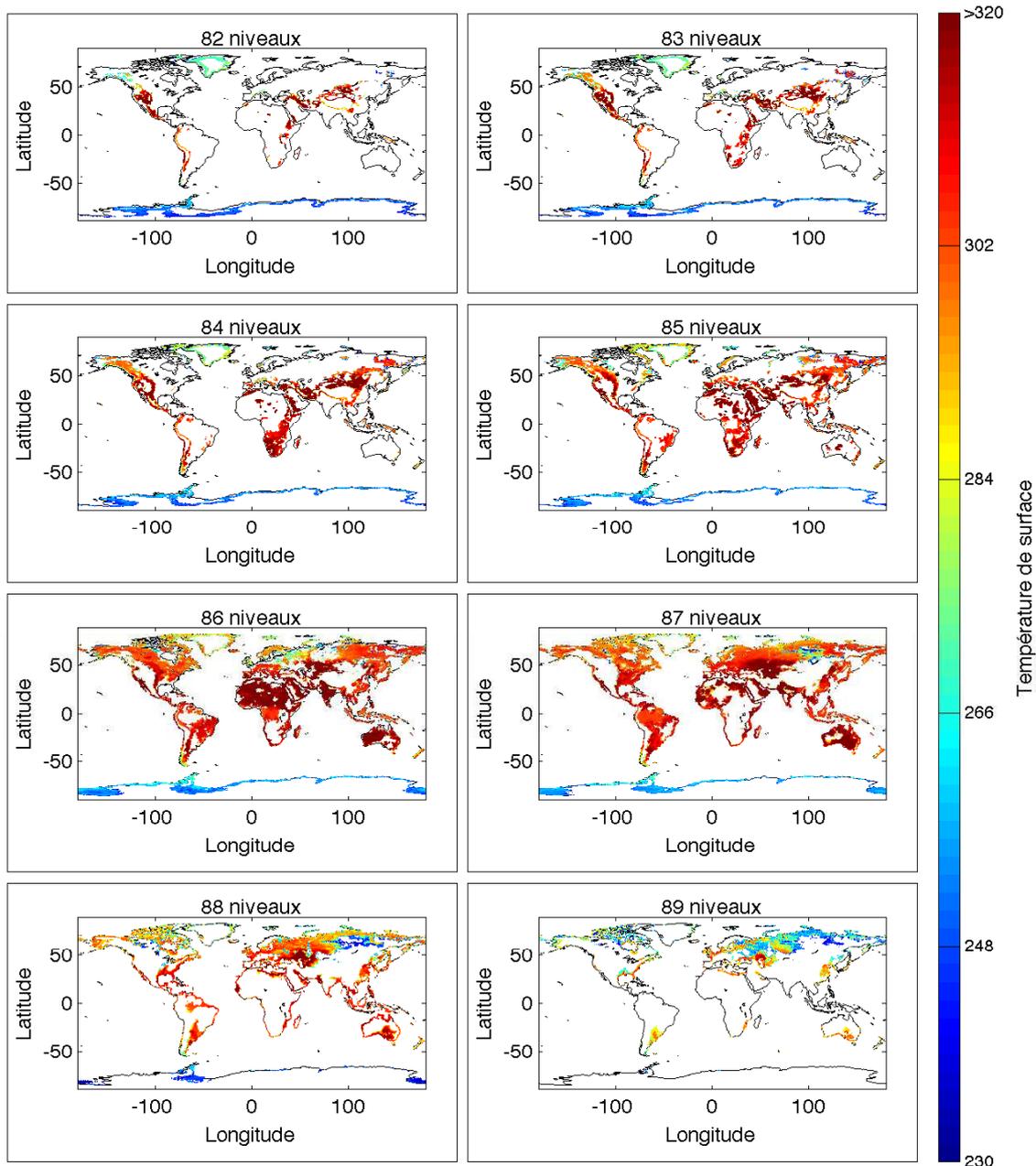


FIGURE 6.1 – De haut en bas et de gauche à droite, répartition spatiale des situations contenues dans les 8 bases de données atmosphériques. La couleur correspond à la température de surface pour chaque point.

### 6.3 Échantillonnage

L'objectif est d'extraire des bases présentées ci-dessus (à la Section 6.2.2, page 162) des bases d'apprentissage qui seront utilisées par la suite pour paramétrer les algorithmes d'inversion. Ces bases d'apprentissage seront utilisées uniquement pour l'apprentissage des

réseaux de neurones, il faut donc les construire spécifiquement pour cette technique. Il faut représenter l'ensemble des situations atmosphériques possibles pour que le réseau, une fois paramétré, puisse les restituer correctement.

La méthode des k-moyennes dont on s'est servie précédemment pour échantillonner de telles bases (voir Section 5.1.2, page 124) n'est pas utilisée ici. Cette méthode conserve la variabilité de la base, alors que nous voulons la lisser afin d'avoir autant de représentations des situations moyennes que des situations extrêmes. Nous avons fait le choix d'un algorithme statistique, c'est-à-dire qui conserve les statistiques de variabilité de la base de données. Si l'on cherche à obtenir une erreur globale plus faible, il est préférable de sur-échantillonner les situations les plus courantes et donc d'opter pour un échantillonnage statistique. C'est le choix que nous avons fait précédemment. Si au contraire on cherche à éviter les pics d'erreurs de restitution, il faut représenter les situations extrêmes dans la base d'apprentissage et donc s'orienter vers un échantillonnage uniforme.

Nous présentons dans cette partie l'algorithme d'échantillonnage mis en place pour extraire 8 bases d'apprentissage à partir des 8 bases présentées ci-dessus (Section 6.2.2, page 162). La même technique est utilisée à chaque fois. Rappelons que ces huit bases correspondent à des situations présentant 8 différents nombres de niveaux de pression. La base contenant 86 niveaux de pression contient le plus de profils, les résultats intermédiaires seront présentés uniquement sur cette base. Les résultats sur les autres bases étant similaires, ils ne sont pas présentés ici afin d'éviter d'alourdir le manuscrit.

### 6.3.1 Méthode d'échantillonnage : l'entropie

Nous allons utiliser l'échantillonnage par entropie présenté en Section 6.1.1 (page 158). Cela consiste à maximiser l'entropie de Shannon d'une variable au sein de la base de données échantillonnée. Nous avons une base de données multivariée, où chaque profil est constitué de températures et de quantités de vapeur d'eau sur 82 à 89 niveaux de pression, de quantités d'ozone sur 4 couches et de la température de surface. Afin de clarifier les explications, nous nous plaçons ici dans le cas où 86 niveaux de pression sont disponibles pour les profils. Le même raisonnement est mené pour les 7 autres bases de données.

Afin de réduire la quantité de variables à échantillonner et de réduire les temps de calcul et parce que cela n'a pas d'impact sur les résultats obtenus, nous ne considérons (uniquement dans le cadre de l'échantillonnage) qu'un niveau de pression sur deux en température et en vapeur d'eau. Il ne reste ainsi que 91 variables différentes sur lesquelles effectuer l'échantillonnage (43 niveaux en température et en vapeur d'eau, 4 niveaux en ozone et la température de surface). Il faut toutefois noter que le nombre total de variables sur lesquelles est effectué l'échantillonnage est variable d'une base à l'autre, suivant le nombre total de niveaux de pression.

### 6.3.1.1 Discrétisation

Comme expliqué à la Section 6.1.1 (page 158), l'entropie est applicable à une variable discrète. Nous considérons des températures et des quantités de vapeur d'eau ou d'ozone. Ces variables étant continues, il a fallu les discrétiser. Pour cela, nous considérons l'intervalle de variation  $[min; max]$  de chacune des variables considérées. Ces intervalles sont découpés en 20 sous-intervalles. Nous allons alors considérer pour chaque variable (chaque valeur de température ou de vapeur d'eau aux différents niveaux de pression, les 4 valeurs d'ozone et la température de surface) leur probabilité d'appartenance à ces 20 sous-intervalles.

Les sous-intervalles ici définis sont propres à chaque variable. Les probabilités d'appartenance à chaque intervalle sont calculées statistiquement sur la base de données. C'est à partir de ces probabilités d'appartenance aux différents intervalles d'existence que nous pourrions calculer l'entropie de chaque variable. En sommant ces différentes entropies, nous obtiendrions l'entropie globale de la base.

### 6.3.1.2 Variabilité de la vapeur d'eau

Après diverses analyses, la faible variabilité de la vapeur d'eau (en humidité spécifique) dans les hautes couches de l'atmosphère rendait inutile le calcul de l'entropie sur ces couches. Chaque niveau de pression supplémentaire considéré implique un nouveau calcul d'entropie. Calculer l'entropie pour des niveaux de pression n'ayant pas de variabilité augmente inutilement la quantité de calculs à effectuer et donc le temps d'échantillonnage de la base. Afin de déterminer le nombre de niveaux en vapeur d'eau qui sera pris en compte, une courte étude de la variabilité de la vapeur d'eau dans la base de données globale est nécessaire.

La Figure 6.2 représente la quantité de vapeur d'eau dans chacun des 20 sous-intervalles pour chaque niveau de pression. La couleur de chaque intervalle correspond à la probabilité que la quantité de vapeur d'eau appartienne à cet intervalle. Les intervalles laissés en blanc représentent des valeurs de vapeur d'eau qui n'existent pas.

La quantité de vapeur d'eau au sommet de l'atmosphère n'est pas comparable à celle à la surface. Afin de pouvoir visualiser la répartition de vapeur d'eau sur chaque couche, nous avons choisi de représenter en abscisse les intervalles d'existence. Les valeurs minimales et maximales de vapeur d'eau de chaque niveau de pression sont normalisées, ce qui explique que le rendu final soit carré (contrairement aux densités de probabilité présentées à la Figure 6.4, par exemple).

Ces statistiques ont été calculées sur la base complète de profils sur 86 niveaux (voir Section 6.2.2, page 162). On remarque sur cette figure qu'à partir de 102.05 hPa (cela correspond au 46<sup>ème</sup> niveau), seuls quelques intervalles sont peuplés. Ainsi, presque tous les profils de vapeur d'eau présentent les mêmes valeurs dans les hautes couches de l'atmosphère. L'échantillonnage ne pourra donc pas être en mesure de sélectionner des profils différents susceptibles de maximiser l'entropie.

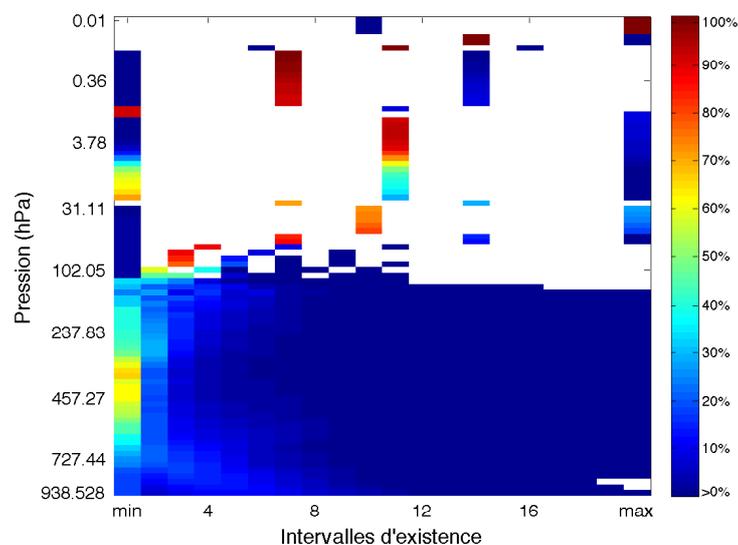


FIGURE 6.2 – Densité de probabilité de la vapeur d’eau par niveau de pression. Pour chaque niveau de pression, 20 intervalles sont créés entre le minimum et le maximum du niveau. La couleur représente le pourcentage de situations appartenant à cet intervalle.

Le choix est fait de ne pas prendre en compte les niveaux de vapeur d’eau au-dessus de cette valeur de pression. Seuls sont pris en compte les niveaux 47 à la surface (*i.e.*, le niveau 86). Il ne reste alors que 40 niveaux en vapeur d’eau. Nous avons choisi de ne considérer qu’un niveau sur deux pour le calcul de l’entropie, il ne reste alors que 20 niveaux de vapeur d’eau à prendre en compte dans le calcul.

### 6.3.1.3 Pondération de chaque variable

Pour passer de l’entropie d’une variable à l’entropie générale multivariée, nous sommes les entropies de chaque variable. Cependant, du fait du nombre de niveaux de pression sur lesquels sont disponibles la température (86 niveaux, on n’en garde qu’un sur deux donc 43 valeurs) et la vapeur d’eau (40 niveaux, avec un sur deux sélectionné, donc 20 valeurs), il a fallu attribuer un poids à ces différentes entropies. Ces différents poids permettent d’éviter que la température et la vapeur d’eau n’influencent trop l’entropie totale rendant ainsi sous-représentée la température de surface et l’ozone. L’entropie d’un profil complet de température serait 43 fois supérieure à l’entropie de la température de surface.

Après divers tests sur l’influence du poids de l’entropie totale du profil de température et de vapeur d’eau, nous avons constaté qu’un poids trop élevé pour la vapeur d’eau entraînait une mauvaise représentation du profil de température. Ceci est dû au fait que les situations avec une humidité élevée sont généralement des situations tropicales, présentant un profil de température toujours semblable. Ainsi, plus nous cherchons à uniformiser la répartition des différents profils d’humidité, plus les profils de température associés seront semblables.

Il a donc fallu diminuer le poids de l'entropie totale de la vapeur d'eau par rapport à celui de la température. Au final, il a été décidé de ramener à 15 l'entropie sur les niveaux de température et à 5 l'entropie sur les niveaux de vapeur d'eau. Ainsi, la somme de toutes les entropies des différents niveaux de température sera 15 fois plus importante que l'entropie de la température de surface ou que l'entropie d'un niveau d'ozone, celle sur les niveaux de vapeur d'eau sera seulement 5 fois plus importante. Plusieurs tests ont été menés afin de déterminer cette pondération. Il a fallu faire un choix entre la bonne représentation des profils de température et ceux de vapeur d'eau.

Pour calculer l'entropie multivariée de notre base de données, on somme donc l'entropie de la température de surface, celle de chacun des quatre niveaux d'ozone, celle des différents niveaux de température pondérés par  $\frac{15}{\text{nombre total de niveaux considérés pour la température}}$  et celle des différents niveaux de vapeur d'eau pondérés par  $\frac{5}{\text{nombre total de niveaux considérés pour la vapeur d'eau}}$ .

Chaque variable  $x$  de la base échantillonnée de 10.000 situations présente une probabilité  $p_n^x$  d'appartenir à l'intervalle  $n$ . Son entropie vaut donc, par définition :

$$\sum_{inter=1}^{20} \left( p_n^x \cdot \ln\left(\frac{1}{p_n^x}\right) \right)$$

Il ne reste alors plus qu'à sommer les entropies des différentes variables considérées : les 43 niveaux de température, les 20 niveaux de vapeur d'eau, les 4 couches d'ozone et la température de surface. Il s'agit de trouver 10.000 situations qui maximisent la somme :

$$\sum_{niveaux=1}^{43} \left( \sum_{inter=1}^{20} p_n^{temp} \cdot \ln\left(\frac{1}{p_n^{temp}}\right) \right) \cdot \frac{15}{43} + \sum_{niveaux=1}^{20} \left( \sum_{inter=1}^{20} p_n^{vap} \cdot \ln\left(\frac{1}{p_n^{vap}}\right) \right) \cdot \frac{5}{20} + \sum_{niveaux=1}^4 \left( \sum_{inter=1}^{20} p_n^{ozo} \cdot \ln\left(\frac{1}{p_n^{ozo}}\right) \right) + \sum_{inter=1}^{20} p_n^{TS} \cdot \ln\left(\frac{1}{p_n^{TS}}\right) \quad (6.2)$$

Le premier terme correspond à la somme sur les 43 niveaux considérés de l'entropie de la température de l'atmosphère, calculée sur les 20 sous-intervalles. Cette somme est pondérée par  $\frac{15}{43}$  afin d'éviter que la somme des 43 entropies ne soit prépondérante par rapport aux autres termes, comme expliqué précédemment. Le deuxième terme correspond à la même somme pour la vapeur d'eau, la pondération est ici de  $\frac{5}{20}$  car seuls les niveaux 47 à 86 sont pris en compte et on n'en utilise que un sur deux. Il ne reste donc que 20 niveaux. Le troisième terme correspond à la somme de l'entropie de chacune des quatre valeurs d'ozone sur les couches considérées et le dernier terme correspond à l'entropie de la température de surface.

### 6.3.1.4 Optimisation du calcul

On cherche donc à trouver une base de données de 10.000 échantillons, parmi 15 millions de situations, qui maximise le terme (6.2). Notre méthodologie consiste à sélectionner dans un premier temps 10.000 situations aléatoirement parmi la base complète. Ensuite, en relisant la totalité de la base, situation par situation, on regarde si la situation augmente l'entropie de la base échantillonnée en l'interchangeant avec chacune des situations sélectionnées. Cela conduira au final à ne garder que les situations qui maximisent l'entropie totale de la base.

Cette méthode présente de nombreuses limitations techniques. Le fait que la base complète soit constituée de millions de profils est une première limitation. Une telle base de données ne peut être stockée dans la mémoire vive d'un ordinateur actuel. S'il faut accéder au disque dur à chaque fois que l'on cherche si un profil peut s'échanger avec un autre de la base de 10.000 échantillons, le temps de calcul devient trop élevé. Afin de réduire les accès au disque dur, trop chronophages, nous décidons de fonctionner par "batchs" (*i.e.*, paquets). Il s'agit de procéder par itérations : à chaque itération, 100.000 profils sont lus (de façon aléatoire parmi la base complète) et stockés dans la mémoire vive. Les 10.000 situations maximisant l'entropie sont mises à jour et ensuite une autre base de 100.000 profils aléatoires est lue. On reproduit cette itération jusqu'à convergence de l'échantillonnage. Cette technique permet, d'une part, de concentrer les accès au disque dur, mais aussi de considérer potentiellement plusieurs fois chaque profil de la base complète et ainsi d'envisager le plus de combinaisons possibles pour déterminer la base de 10.000 profils optimale.

Le temps de calcul reste cependant encore élevé. À chaque fois que l'on teste un profil, il faut recalculer l'entropie de la base des 10.000 situations en interchangeant ce profil avec chacun des 10.000 profils un à un. On calcule alors 10.000 entropies pour chaque profil que l'on teste. Afin de pallier ce problème, un algorithme d'optimisation est mis en place. Ce dernier rend le code nettement moins accessible et compréhensible, mais le gain en rapidité justifie ce choix.

Il s'agit, à chaque fois que la base de 10.000 profils est modifiée, de calculer le changement dans les densités de probabilités par intervalles  $p_n$  (plus exactement le changement de la somme des  $p_n \cdot \ln(p_n)$ ), définies ci-dessus, associé au retrait de chaque profil de la base. On calcule également le changement de ces probabilités lié à l'ajout du profil courant. Il suffit alors de chercher le maximum parmi la somme de ce changement et la différence pour chacun des 10.000 échantillons. Si le maximum de ces 10.000 sommes est positif, alors on échange le spectre courant avec le spectre associé à ce maximum, sinon, ce spectre ne peut pas augmenter l'entropie de la base échantillonnée.

À chaque fois qu'un profil est modifié dans la base échantillonnée le calcul du changement des densités de probabilités en cas de retrait de l'un d'eux est mis à jour. Cette méthode permet de diminuer le nombre de calculs dans le cas où le profil courant n'améliore pas l'entropie. Au lieu de recalculer les densités de probabilités et l'entropie totale, dans chacun

des 10.000 cas où le profil courant remplace un des profils présélectionnés, il suffit de calculer directement son apport entropique hypothétique. Dans le cas où le profil considéré améliore l'entropie de la base échantillonnée, le nombre de calculs reste le même.

Il s'agit ici d'un simple artifice technique visant à accélérer le temps d'échantillonnage. Le principe d'échantillonnage reste toutefois le même : pour chaque nouvelle situation, on calcule l'entropie de la base échantillonnée en échangeant un profil avec cette situation. On échange éventuellement cette situation avec un profil de la base échantillonnée si cela augmente l'entropie totale.

### 6.3.1.5 Convergence de l'algorithme d'échantillonnage

Il s'agit désormais de déterminer un critère de convergence afin de stopper les itérations l'algorithme d'échantillonnage lorsque les résultats sont satisfaisants. En effet, le nombre de sélections aléatoires par batchs de 100.000 profils peut être conséquent. D'autant plus qu'un groupe de profils assez proches peuvent être interchangeable et ne modifier que très peu l'entropie. Le critère de convergence choisi dépend du nombre de changements de profils parmi la base échantillonnée de 10.000 profils effectués au cours d'un batch. Si pour deux batchs consécutifs ce nombre est inférieur à 30 alors on considère que l'algorithme a convergé, car la base finale échantillonnée n'est presque plus modifiée d'une itération à l'autre.

La Figure 6.3 illustre la convergence de l'algorithme d'échantillonnage de la base avec 86 niveaux de pression. La courbe (A), en haut, représente l'entropie de la base échantillonnée (celle de 10.000 profils, correspondant au terme (6.2)) en fonction du nombre de batchs effectués. On constate que l'entropie de la base a bel et bien augmenté et qu'elle converge pour se stabiliser à une valeur avoisinant 28,1. La sélection par batchs implique que les profils peuvent être pris en compte plusieurs fois chacun. Le nombre de batchs effectués est donc une meilleure indication, pour la convergence de l'algorithme, que le nombre de situations examinées (qui vaut 100.000 fois plus).

Nous ne représentons pas le nombre de profils interchangeés pour chaque batch. Ce nombre est trompeur puisque après plusieurs interchanges, la base échantillonnée peut n'avoir été modifiée que d'une situation. Un même profil de la base échantillonnée peut être modifié plusieurs fois. C'est pourquoi la courbe (B) au-dessous présente le nombre de profils différents dans la base échantillonnée à l'issue de chaque batch (cette courbe correspond au critère de convergence). On constate qu'au cours des premiers batchs, un grand nombre de profils est modifié. Ce nombre de modifications diminue très rapidement.

Ces courbes viennent valider le choix du critère de convergence. L'évolution typique des différentes courbes, avec une dérivé qui tend à s'annuler, prouve la convergence de notre algorithme. L'échantillonnage des bases de données correspondant à d'autres nombres de niveaux de pression présentent des statistiques similaires.

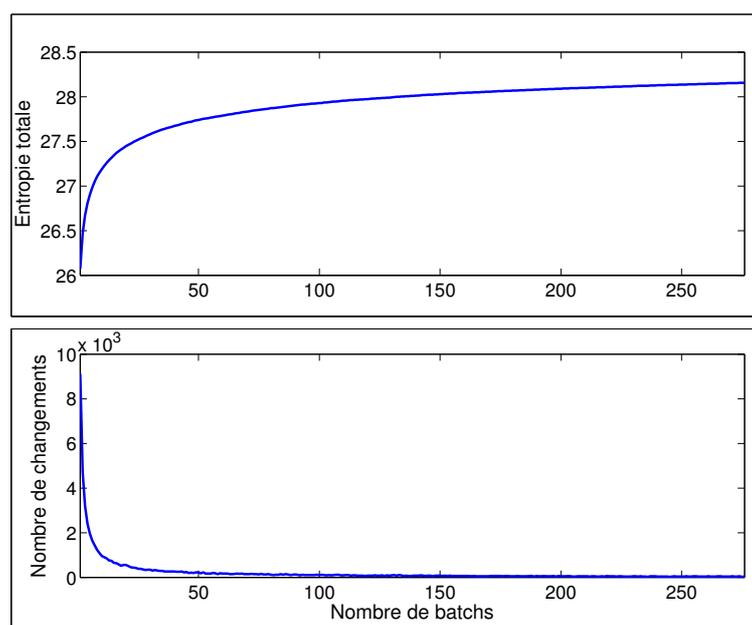


FIGURE 6.3 – De haut en bas : (A) évolution de l’entropie de la base échantillonnée en fonction du nombre de batches de 100.000 échantillons pris en compte ; (B) Nombre de profils différents dans la nouvelle base échantillonnée après chaque batch.

Il reste désormais à justifier que la maximisation de l’entropie de la base échantillonnée nous a permis d’extraire une base de profils représentant les extrema de la base complète autant que les situations moyennes et ce sur toutes les variables prises en compte.

### 6.3.2 Apport de l’entropie pour l’échantillonnage multivarié

Dans un premier temps, une étude simple de la capacité de l’entropie à échantillonner une base de données multivariée a été menée. Il s’agit de la principale limitation des algorithmes d’échantillonnage. Il faut que ces derniers puissent prendre en compte des variables très différentes tant en termes de dimensions que de gammes de variabilité. Cette étude a été menée en marge de l’échantillonnage des bases de données, afin de démontrer le caractère multivarié de l’entropie.

#### 6.3.2.1 Exemple d’échantillonnage sur une seule variable

Nous utilisons ici, pour comparaison, les bases de données développées par l’ECMWF (Chevallier et al. 2006). Ces bases sont constituées de profils de température, d’humidité spécifique, d’ozone, d’eau liquide, d’eau sous forme de glace, de pluie et de neige, le tout sur 91 niveaux. Pour chaque situation, la température de surface est également fournie. Un algorithme d’échantillonnage a été mis en place, mais uniquement pour une variable donnée, ou un profil (température, vapeur d’eau, ozone, eau liquide ou glace). Il existe donc plusieurs

bases de données correspondant aux échantillonnages sur la température, la vapeur d'eau, l'ozone, eau liquide, glace...

L'algorithme d'échantillonnage consiste en une maximisation de la distance inter-profil au sein de la base de données. Cette méthode permet d'obtenir des bases de données presque uniformes pour la variable considérée (à l'image de ce que permet le calcul de l'entropie). Il serait possible, mais difficile, d'appliquer une telle méthode sur plusieurs variables, du fait de la différence d'unité et de gamme de variation de chaque variable. Aucune contrainte n'est donc imposée sur les variables non prises en compte dans l'échantillonnage.

Afin de mettre en avant l'utilité d'un échantillonnage multivarié, nous nous intéresserons uniquement dans cette partie à deux variables : le profil de température et le profil de vapeur d'eau. Nous considérerons le profil de vapeur d'eau en humidité spécifique et relative car, même si l'échantillonnage est effectué sur l'humidité spécifique, il est plus simple d'analyser la variabilité de l'humidité relative.

La Figure 6.4 représente la variabilité de ces variables au sein des bases de l'ECMWF échantillonnées en température (en haut) et en humidité spécifique (en bas). Nous discrétisons les profils en découpant leurs intervalles d'existence en 20 sous-intervalles, comme nous l'avons fait précédemment. Cette discrétisation nous permet de mieux caractériser la variabilité des profils en représentant leur probabilité d'appartenance à chacun des sous-intervalles de chaque niveau de pression.

La couleur de ces carrés correspond à la probabilité d'appartenance du profil à ce sous-intervalle. Plus le carré sera rouge, plus le pourcentage de profils passant par ce sous-intervalle est élevé. Inversement, plus le pourcentage de profils y passant est faible, plus la couleur sera bleue.

Contrairement à ce que l'on avait représenté tout à l'heure, nous utilisons ici en abscisses les valeurs de température et de vapeur d'eau. Ainsi, la figure obtenue n'est plus carrée, puisque la position de chaque sous-intervalle dépend de ses bornes. On voit clairement que la gamme de variabilité de la vapeur d'eau n'est pas la même dans les hautes couches de l'atmosphère que proche de la surface. L'humidité spécifique est nettement plus faible dans les hautes couches de l'atmosphère. Cette différence est moins visible en termes d'humidité relative, mais il faut garder à l'esprit que plus on monte le long de l'atmosphère, plus la densité de molécules diminue. Nous avons choisi une représentation en fonction des valeurs de température et de vapeur d'eau, afin de mieux visualiser les profils atmosphériques en eux-mêmes. En effet, la coloration des sous-intervalles laisse apparaître certains pics rouges, correspondant aux profils les plus représentés dans la base. Pour chaque niveau de pression, le maximum et le minimum de la variable considérée sont représentés en noir. Ils encadrent donc les 20 sous-intervalles de chaque niveau de pression.

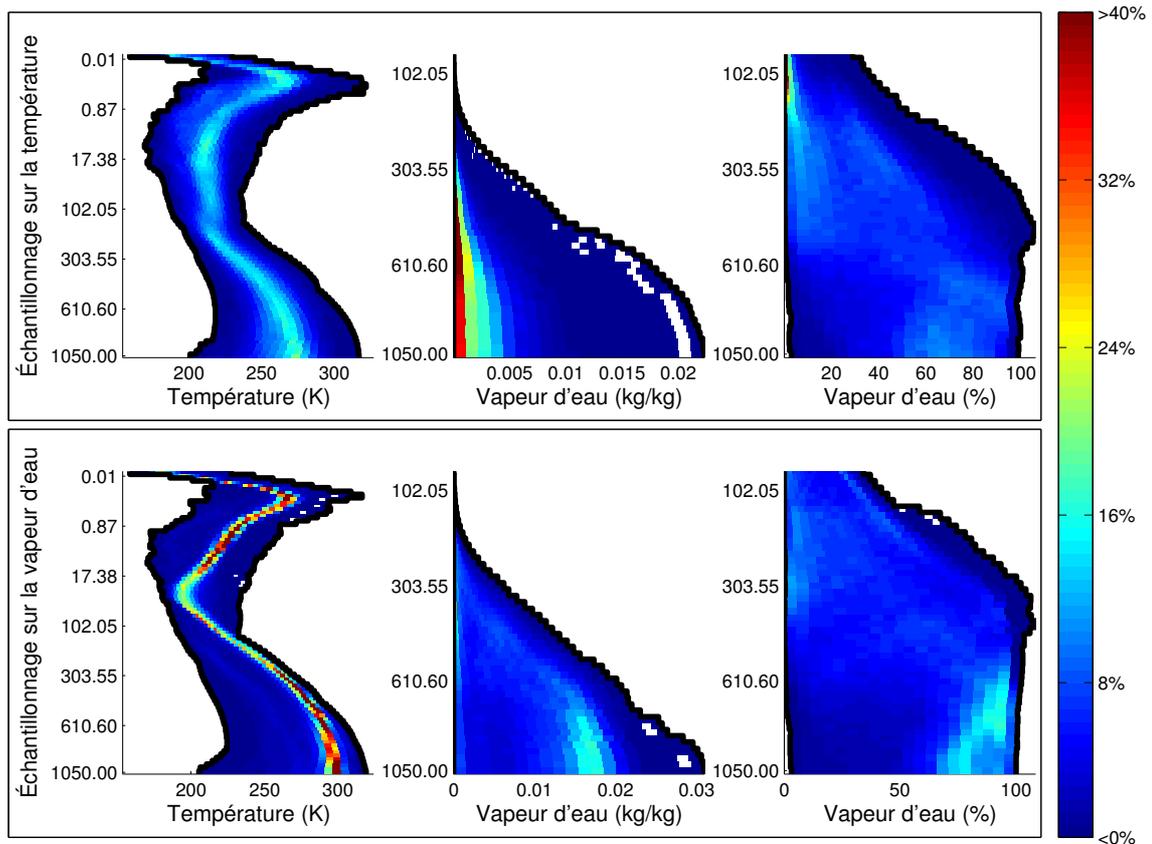


FIGURE 6.4 – Densité de probabilité d’existence par intervalles des profils de température et d’humidité spécifique et relative de deux bases de l’ECMWF. Au-dessus, la base échantillonnée en température, au-dessous, celle échantillonnée en humidité spécifique.

Comme expliqué à la Section 6.3.1.2 (page 166), les plus hautes couches de l’atmosphère présentent peu de variabilité en vapeur d’eau. Nous n’avons donc pas montré ici les plus hautes couches de l’atmosphère sur les profils de vapeur d’eau. Ceci explique la différence d’ordonnées entre les profils de température et de vapeur d’eau.

On constate sur cette figure que la base de données échantillonnée en température présente, dans des proportions similaires, tout type de profils de température (la couleur des rectangles correspondant à chaque intervalle reste très bleutée). Les profils de vapeur d’eau de cette base sont nettement moins variables. Si les profils d’humidité spécifique présentent souvent une faible quantité proche de la surface, on constate également que la représentation de l’humidité relative est faible pour des taux élevés d’humidité.

Les courbes dans la partie basse de la figure correspondent à la base de données échantillonnée en humidité spécifique. Ici, la variabilité des profils d’humidité (spécifique ou relative) est importante, même pour des quantités de vapeur d’eau élevées. Par contre, la représentation des profils de température présente un maximum (en rouge), ce qui montre que certains profils (ceux passant par les rectangles les plus rouges) sont sur-représentés au

sein de la base. Ceci s'explique par la même raison qui nous a poussé à diminuer le poids relatif de l'entropie sur la vapeur d'eau par rapport à celle sur la température : les profils présentant des taux d'humidité élevés sont pour la plupart situés au niveau des tropiques et présentent des profils de température similaires.

On voit, ici, la difficulté : échantillonner sur une seule variable conduit à un mauvais échantillonnage des autres variables. Il est nécessaire de considérer les variables de façon simultanée afin de parvenir à échantillonner correctement une base de données. Rassembler les deux bases de données correspondant aux échantillonnages distincts ne masquera pas la sur-représentativité du profil moyen de température et la sous-représentativité des profils humides.

### 6.3.2.2 Prise en compte de plusieurs variables grâce à l'entropie

La Figure 6.5 présente le résultat de l'utilisation de l'entropie comme critère d'échantillonnage multivarié. Les figures sont construites de la même façon que pour la base de l'ECMWF.

Trois bases de données de 10.000 profils sont extraites de la base complète de profils sur 86 niveaux. La première base de données extraite ne prend en compte que l'entropie sur les profils de température (représenté dans la partie haute de la figure). On y retrouve les mêmes résultats que sur la base de l'ECMWF échantillonnée en température, à savoir une bonne représentativité de la température mais des profils de vapeur d'eau très similaires.

La base échantillonnée en vapeur d'eau (au milieu sur la figure) est également équivalente à celle construite par l'ECMWF. On retrouve encore ici le même problème : la sur-représentation d'un profil de température moyen. Les sous-intervalles laissés en blanc sont ceux par lesquels aucun profil ne passe. Certaines valeurs de température ne sont donc jamais atteintes. En échantillonnant sur une seule variable, nous avons perdu beaucoup de représentativité sur la deuxième variable. Nous pouvons faire le même constat sur la partie supérieure de la figure où certains profils de vapeur d'eau ne sont plus présents dans la base de données échantillonnée.

La série de figures présentée au bas de la Figure 6.5 montre, quant à elle, qu'en prenant en compte à la fois la température et la vapeur d'eau dans le calcul de l'entropie totale, on peut obtenir une bonne représentativité des deux profils. L'humidité spécifique (au milieu) présente certes un grand nombre de profils avec une faible valeur, mais en se reportant à l'humidité relative, on a des profils plus divers.

On peut également remarquer que dans le cas d'échantillonnages sur une seule variable (partie supérieure de la figure), les extrema de l'autre variable ne sont plus représentés. Les zones blanches (que l'on aperçoit sur les profils de vapeur d'eau sur la partie haute de la figure et sur les profils de température sur la partie du milieu) correspondent à des intervalles qui ne sont plus représentés dans la base échantillonnée. L'échantillonnage multivarié permet une meilleure représentation des extrema des deux variables simultanément.

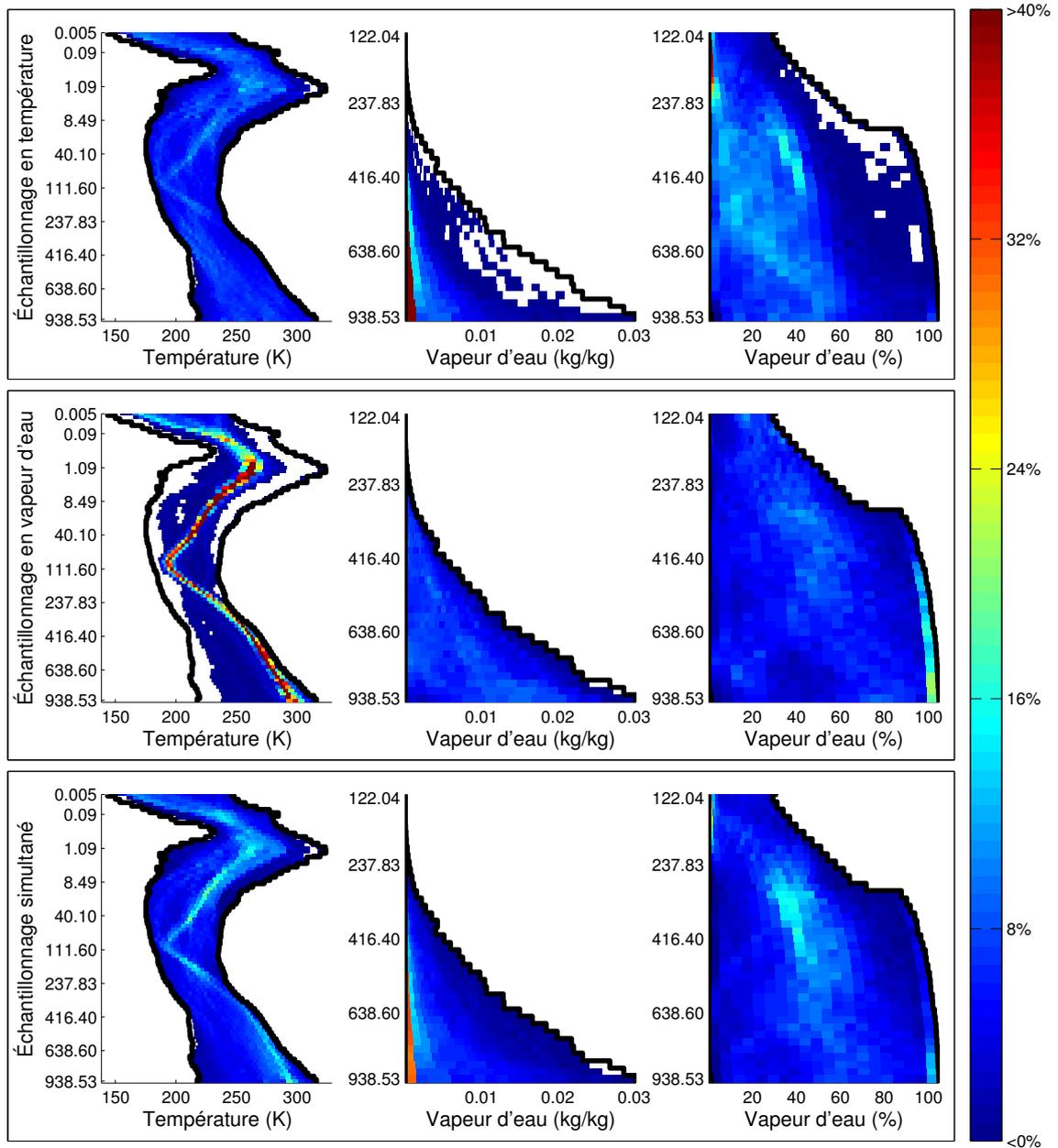


FIGURE 6.5 – Densité de probabilité d’existence par intervalle des profils de température et d’humidité spécifique des bases de données échantillonnées par entropie en considérant uniquement la température (en haut), uniquement la vapeur d’eau (au milieu) et les deux simultanément (en bas).

L’entropie permet donc d’obtenir un critère d’échantillonnage facilement comparable d’une variable à l’autre et peut donc aisément prendre en compte ces dernières simultanément. Il est nécessaire de mettre au point une telle méthodologie d’échantillonnage multivarié car, comme il a été démontré, ne pas prendre en considération certaines variables peut les rendre très peu variables dans la base échantillonnée.

## 6.4 Résultats

### 6.4.1 Représentation des différentes variables

Afin de représenter au mieux la base de données multivariée et les résultats de l'échantillonnage, nous présentons les densités de probabilités des différentes variables pour la base complète et la base échantillonnée.

#### 6.4.1.1 Base de données complète

Les figures représentées ici contiennent les 86 niveaux de pression des profils, même si seul un niveau sur deux a été utilisé au cours de l'échantillonnage. Ici encore, nous ne présentons les résultats que pour l'échantillonnage de la base de données sur 86 niveaux (voir Section 6.2.2, page 162), les résultats de l'échantillonnage des autres bases sont similaires.

La Figure 6.6 représente ces densités de probabilité pour la base de données complète. Pour chaque niveau de température, de vapeur d'eau ou pour chaque couche d'ozone, on représente les 20 sous-intervalles entre le minimum et le maximum. Les couleurs représentent la densité de probabilité, c'est-à-dire le pourcentage de profils passant par ce sous-intervalle. En sommant ces densités de probabilité sur un niveau donné, on retrouve 100%, car chaque profil passe par un seul sous-intervalle de chaque niveau.

On trace la densité de probabilité en humidité relative pour mieux voir la répartition de l'humidité dans l'atmosphère. La grande variabilité de l'humidité spécifique, même sur les 40 plus basses couches considérées ici, ne permettrait pas de visualiser l'échantillonnage dans les plus hautes couches. Le pourcentage de profils dans chacun des sous-intervalles de la température de surface est représentée sous la forme d'une courbe.

On peut remarquer sur cette figure la forte représentation de profils "moyens". Tant en température atmosphérique, qu'en température de surface, contenu d'ozone ou qu'en vapeur d'eau, certains profils sont fréquents dans la base. ils sont symbolisés par des zones rouges ou par un pic de pourcentage pour la température de surface. La base de données complète comporte donc une grande quantité de profils similaires. L'objectif de l'échantillonnage est justement de ne pas conserver autant de profils similaires et de s'appliquer à conserver les extrema et les profils rares. Cette forte représentation de certains profils nous a poussés à développer un algorithme d'échantillonnage complexe. En effet, une sélection aléatoire de profils parmi cette base serait constituée en grande partie de ces profils "moyens".

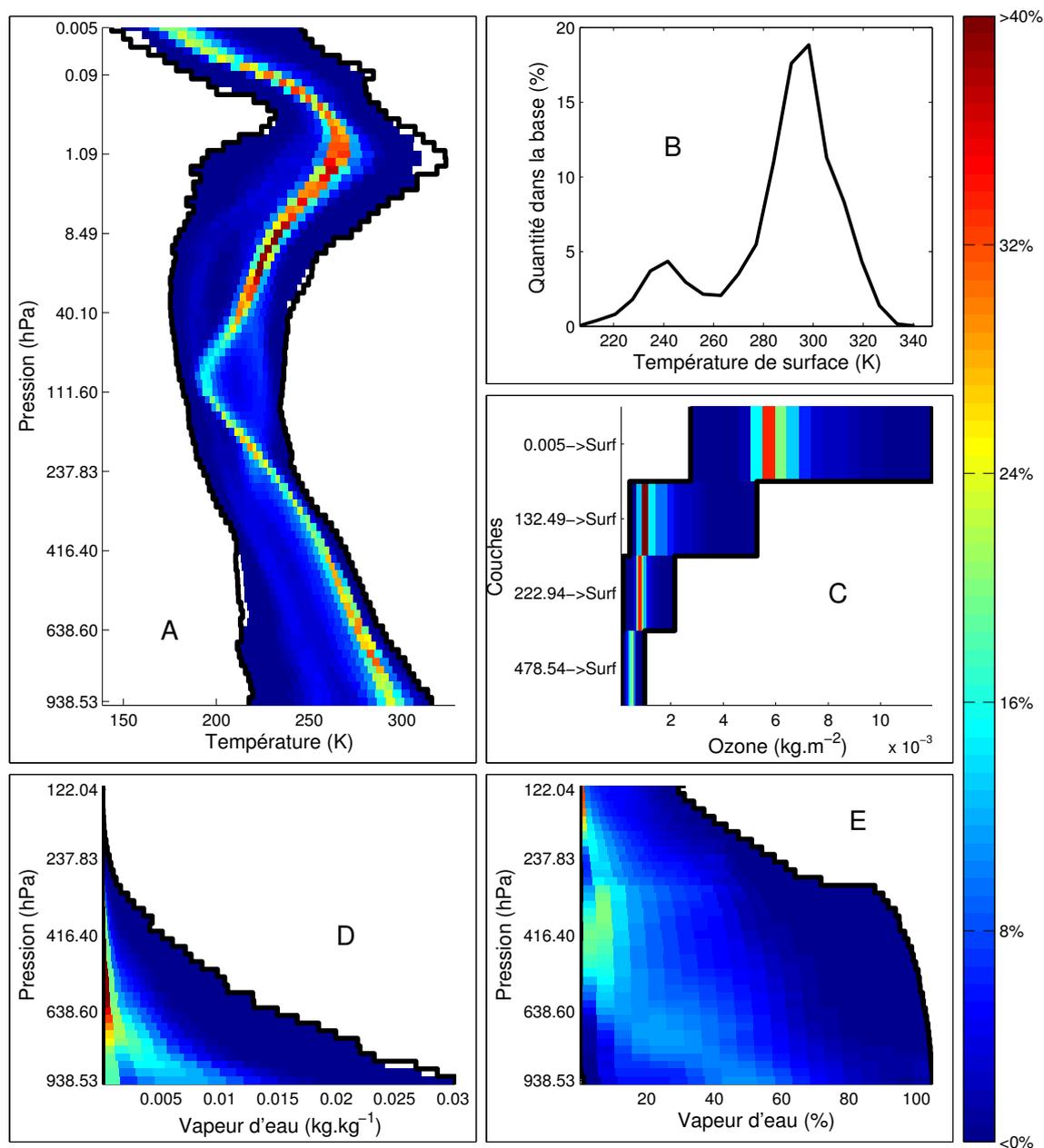


FIGURE 6.6 – De gauche à droite et de haut en bas : représentation de la densité de probabilité : (A) des profils de température ; (B) de la température de surface ; (C) des différentes couches d’ozone ; (D) des profils de vapeur d’eau en humidité spécifique ; (E) des profils de vapeur d’eau en humidité relative. Toutes ces densités de probabilité sont calculées sur la base de données atmosphériques complète (voir Section 6.2.2, page 162).

#### 6.4.1.2 Base de données échantillonnée

La Figure 6.7 correspond à la même figure mais pour la base de données échantillonnée par notre technique basée sur l’entropie. On remarque que le pic de densité de probabilité

de la température de surface autour de 300K a disparu et que la courbe est maintenant plus uniforme. De plus, les extrema ( $< 220$  ou  $> 320$ ) sont maintenant mieux représentés. L'échantillonnage a donc bien fonctionné en température de surface.

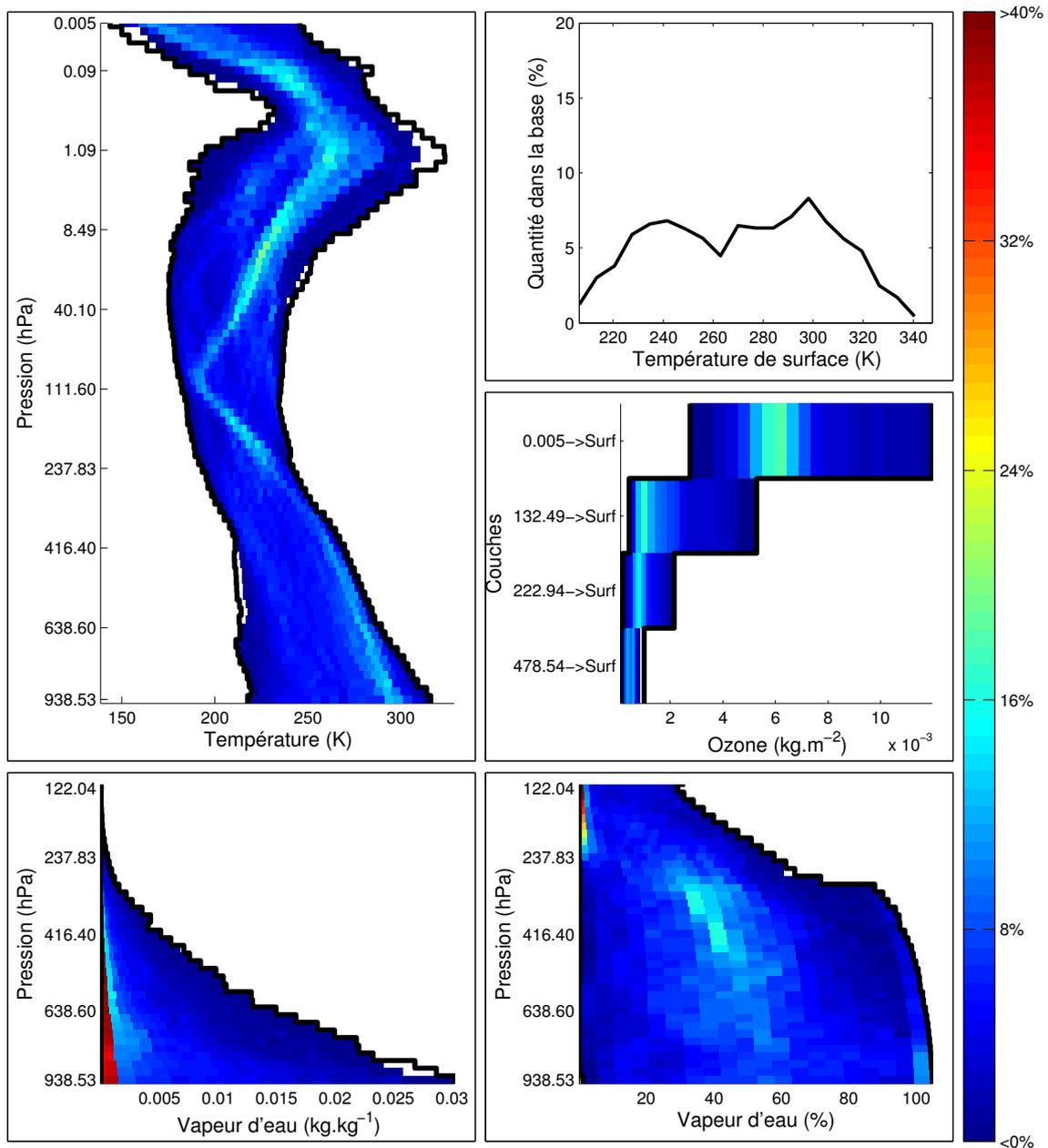


FIGURE 6.7 – Figure similaire à la précédente, mais pour la base de données échantillonnée.

Un même constat peut être fait pour les quatre couches d'ozone. Les zones très rouges ont été estompées, et le reste des sous-intervalles est toujours bleu. Il n'y a donc plus de valeurs sur-représentées, ou plutôt nettement moins, tout en ayant conservé les valeurs extrêmes.

Le profil de température a également perdu la forte densité de probabilité de certains

profils, même si certaines traces subsistent car le profil moyen reste plus représenté que les autres. Notre échantillonnage fait un compromis entre les différentes variables ; s’il cherchait à focaliser l’échantillonnage sur une seule variable, les autres en pâtiraient. Ici, le profil moyen correspond à environ 15% des cas (couleur bleu-vert), contre plus de 30% dans la base complète initiale. Nous sommes donc satisfaits de la répartition des profils.

Les profils d’humidité spécifique sont plus complexes à examiner. En effet, toutes les situations extrêmement humides sont des situations correspondant à un profil de température “moyen”. Ainsi, en augmentant l’entropie pour l’humidité spécifique, on diminue celle du profil de température. Cependant, l’examen des densités de probabilité en humidité relative est plus parlant. On remarque sur cette figure (en bas à droite de la Figure 6.7) une meilleure répartition des densités de probabilité. En regardant, par exemple, les niveaux de pression autour de 416,40 hPa, le pic de densité, présent sur la Figure 6.6 correspondant à la base de données complète, a ici disparu.

L’entropie a donc permis d’échantillonner la base de données complète initiale, constituée de plusieurs millions de situations, pour obtenir une base constituée de seulement 10.000 situations uniformisant la dispersion des différentes variables.

### 6.4.2 Variabilité spatiale

L’échantillonnage par entropie nous a permis de mener à bien un échantillonnage multivarié complexe. Nous pouvons maintenant analyser la répartition géographique des différentes bases échantillonnées (correspondant aux différents nombres de niveaux de pression) et la comparer à la répartition spatiale des bases complètes présentée à la Figure 6.1.

La répartition spatiale des différentes bases de données échantillonnées est présentée sur la Figure 6.8. On retrouve sur cette figure les mêmes répartitions que précédemment. Les différentes bases de données correspondent à des altitudes données et donc à des localisations géographiques restreintes. Cependant, on retrouve une répartition spatiale équivalente aux bases de données complètes, ce qui est un signe de la bonne qualité de l’échantillonnage.

On peut regarder plus précisément la base de données sur 86 niveaux qui est à l’origine bien répartie sur le globe. La base de données échantillonnée constituée de 100 fois moins de représentants couvre tout de même une surface semblable à la base totale. Attention, afin de mieux voir la localisation des différentes bases de données échantillonnées, les points correspondant aux profils sont plus gros que pour les bases complètes. L’objectif étant de montrer que l’échantillonnage n’a pas favorisé une zone du globe en particulier.

Sans aucune indication sur la localisation géographique des situations, l’algorithme a réussi à représenter la variabilité spatiale des bases de données, en plus de respecter la variabilité des différents profils. L’entropie est donc un critère adapté à un échantillonnage, quoique complexe à mettre en place.

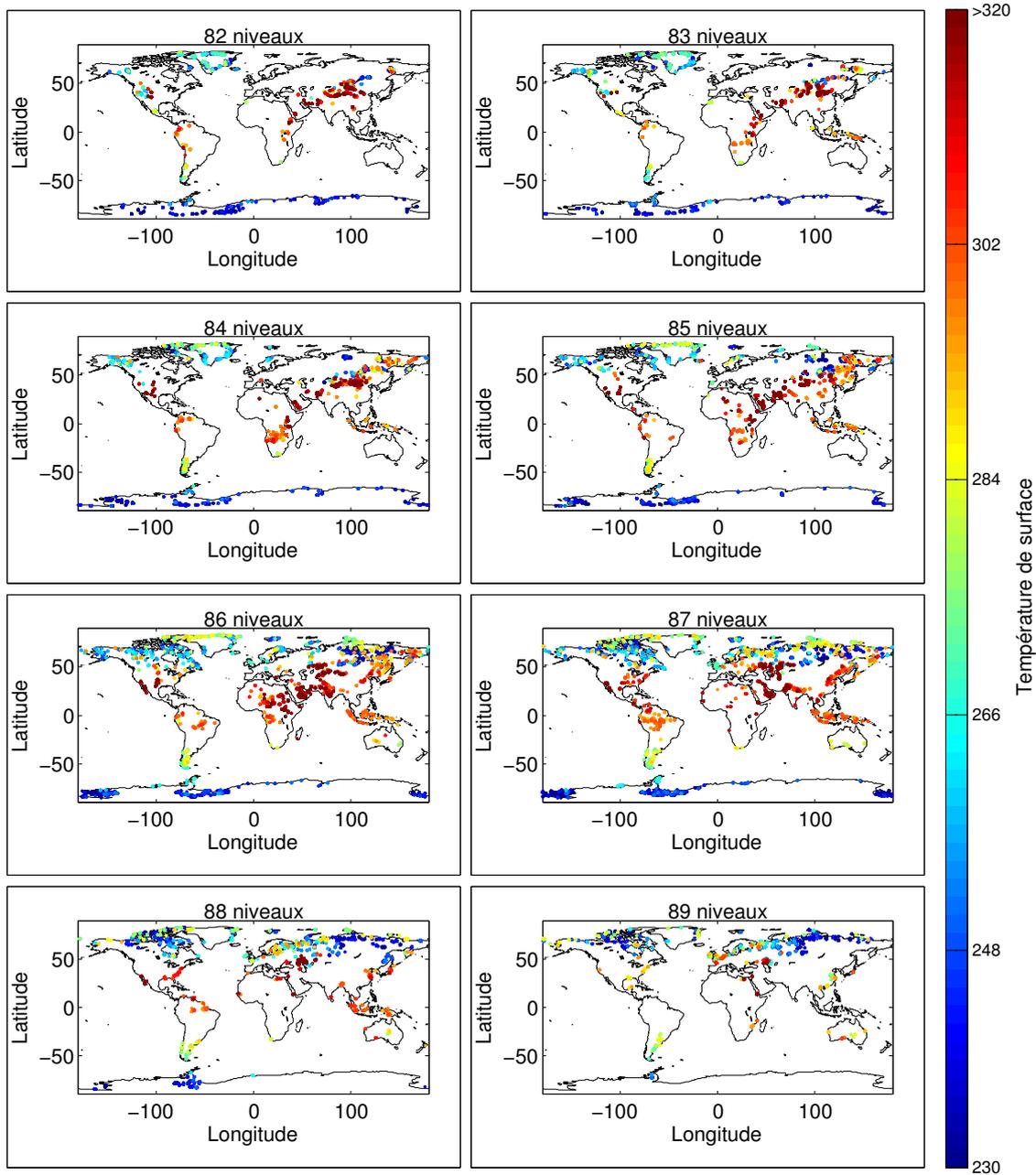


FIGURE 6.8 – De haut en bas et de gauche à droite, répartition spatiale des situations contenues dans chaque base de données échantillonnée. La couleur correspond à la température de surface pour chaque point.

## 6.5 Conclusion

Nous avons mis au point une méthode d'échantillonnage originale. Nous avons prouvé son efficacité pour effectuer un échantillonnage sur un espace de grande dimension avec des variables inhomogènes. Cette méthode est générale et peut être utilisée dans de nombreux

contextes différents.

La base qui a été créée ici peut être fournie à la communauté. L'approche que nous avons utilisée sera présentée aux centres opérationnels, comme l'EUMETSAT et l'ECMWF, qui pourraient avoir besoin de telles bases. Cela peut leur permettre de :

- Retrouver des premières ébauches climatiques pour les algorithmes itératifs, avec une base indépendante des modèles ;
- Faire des études de contenu en informations sur des bases atmosphériques limitées mais représentatives de la nature ;
- Calibrer des méthodes d'inversion.

Cette base de données échantillonnée nous permet, dans la suite, de mettre au point des algorithmes de restitution de profils atmosphériques, plus facilement et plus rapidement, tout en prenant en compte toute la variabilité des profils de température, de vapeur d'eau, d'ozone et également de la température de surface, en particulier leurs extrêmes.



---

# SYNERGIE INFRAROUGE ET MICRO-ONDE AU-DESSUS DES CONTINENTS, EN CIEL CLAIR : PRISE EN COMPTE DES PRO- PRIÉTÉS DE SURFACE



## Sommaire

---

<b>7.1 Configuration des restitutions</b>	<b>185</b>
7.1.1 Base d'apprentissage	186
7.1.2 Méthode de restitution	187
7.1.3 Convergence des différents réseaux de neurones	188
<b>7.2 Apports de l'émissivité et de la température de surface</b>	<b>190</b>
7.2.1 Dans le micro-onde	190
7.2.2 Dans l'infrarouge	192
<b>7.3 Synergie infrarouge et micro-onde</b>	<b>194</b>
7.3.1 Sans informations de surface	195
7.3.2 Synergie complète entre toutes les données disponibles	196
<b>7.4 Synergie et informations de surface</b>	<b>198</b>
<b>7.5 Conclusion</b>	<b>200</b>

---

Nous avons montré, au long du Chapitre 5 (page 121), que la synergie entre les sondeurs infrarouges et micro-ondes de la plateforme MetOp permettait de diminuer l'erreur de restitution des profils atmosphériques au-dessus des océans. Nous allons désormais poursuivre cette étude au-dessus des continents. Ces restitutions impliquent cependant quelques précautions supplémentaires.

L'émissivité micro-onde du sol dépend de son humidité, de la présence de végétation ou de neige, de la topographie et de la composition du sol (Prigent et al. 2006; Jiménez et al. 2010). Le rayonnement mesuré par un satellite sera donc lié à tous ces paramètres de surface. L'émissivité est un bon moyen de caractériser le rayonnement de la surface. On a aussi

montré que l'émissivité infrarouge dépendait de la composition du sol et de la présence de végétation (voir Chapitre 3, page 83).

Pour pouvoir restituer les profils atmosphériques, notamment dans les basses couches, il faut être capable d'identifier et de démêler les différents rayonnements mesurés. Connaître l'émissivité de la surface permet alors de mieux caractériser son rayonnement et donc de mieux restituer la température et la vapeur d'eau dans les basses couches de l'atmosphère. De plus, la covariance entre les émissivités infrarouges et micro-ondes peut permettre de diminuer encore l'erreur de restitution. Si les deux émissivités ont un comportement différent, elles présentent néanmoins des dépendances statistiques exploitables.

L'impact de la surface sur le rayonnement est défini par sa température et son émissivité. Connaître ces deux variables permet de mieux comprendre le rayonnement global mesuré au sommet de l'atmosphère par les satellites. De nombreuses études ont été menées afin de mettre en évidence leur rôle dans les restitutions de profils atmosphériques. Dès les années 70, [Kornsfield and Susskind \(1977\)](#) ont montré l'impact de la prise en compte de l'émissivité infrarouge de la surface sur l'erreur de restitution d'un profil de température. D'autres études ont été conduites concernant le profil de vapeur d'eau, avec des conclusions similaires ([Seemann et al. 2008](#)).

La température de la surface en elle-même est également une des clefs du rayonnement mesuré par le satellite au sommet de l'atmosphère. Il a également été mis en évidence que la corrélation entre la température de surface et la température de l'air doit être prise en compte dans l'assimilation des données satellites afin de mieux exploiter ces données ([Garand et al. 2004](#)). Du point de vue des centres NWP, la Terre est découpée en mailles. Les données *in situ* mesurées dans chaque maille permettent de fournir une première ébauche aux modèles de prévision. À partir de cette première ébauche, les données satellites vont être utilisées afin d'affiner l'estimation de l'état de l'atmosphère. Les bouées, qui couvrent une grande partie des surfaces océaniques, permettent une bonne assimilation des données satellites au-dessus des océans. Certaines zones sur les continents restent très peu instrumentalisées. L'absence de ces données oblige la prise en compte d'informations de surface précises venant d'autres sources d'information ([English 1999](#)).

La prise en compte de l'émissivité et de la température de surface est importante pour une restitution précise des profils atmosphériques de température et de vapeur d'eau. La variabilité spatiale et temporelle de la température de surface impose une restitution en temps réel afin de la déterminer de façon précise. Des études ont prouvés l'existence d'une variabilité diurne de l'émissivité de surface dans l'infrarouge ([Li et al. 2012](#)), nous avons également étudié sa variabilité spatiale. Il est donc nécessaire de restituer en temps réel l'émissivité de la surface.

Ce chapitre repose sur tout ce que nous avons créé jusqu'à présent. La compréhension du rayonnement mesuré par le satellite nous est fournie par le Chapitre 1 (page 5). Les restitutions de températures et d'émissivités de surface effectuées aux Chapitres 2 (page 37)

et 3 (page 83) nous fournissent l'information nécessaire sur l'état de la surface. À partir du savoir faire développé au-dessus des océans aux Chapitres 4 (page 107) et 5 (page 121), nous pouvons construire des schémas de restitution de profils atmosphériques au-dessus des continents. L'échantillonnage par entropie développé au Chapitre 6 (page 155) nous fournit une base d'apprentissage robuste pour paramétrer nos algorithmes d'inversion. Ce chapitre a fait l'objet d'une publication : [Paul and Aires \(2013b\)](#).

## 7.1 Configuration des restitutions

Nous allons effectuer des restitutions de profils atmosphériques de température et de vapeur d'eau. L'objectif est de montrer à la fois la synergie entre les mesures des instruments satellitaires (comme on l'a fait au Chapitre 5, page 121), mais également de montrer l'importance de la prise en compte des informations de surface dans la restitution. Pour cela, nous avons légèrement complexifié le schéma de restitution utilisé au-dessus des océans, en y rajoutant les informations de surface en entrée (voir Schéma 7.1).

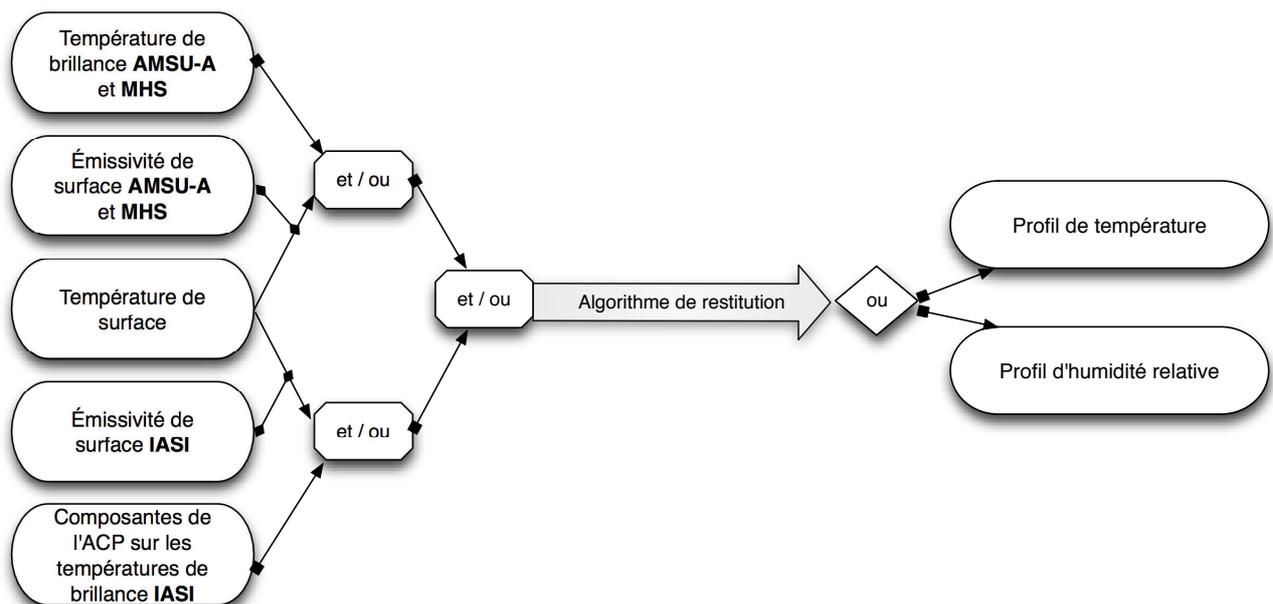


FIGURE 7.1 – Schéma des configurations de restitutions au-dessus des continents.

Cette étude sera menée en deux temps, nous étudierons d'abord l'apport des informations de surface. Pour cela nous considérerons des schémas de restitution utilisant comme entrée les mesures des instruments seuls ou combinés avec les émissivités correspondantes et la température de surface. Puis, nous étudierons l'impact de la synergie entre les rayonnements infrarouges et micro-ondes sur la restitution en combinant les différentes entrées.

### 7.1.1 Base d'apprentissage

Nous allons utiliser ici les bases de données atmosphériques que nous avons créées au chapitre précédent (Chapitre 6, page 155). Il s'agit de huit bases de données différentes (correspondant à différents nombres de niveaux de pression) constituées de :

- Un profil de température ;
- Un profil de vapeur d'eau ;
- Quatre valeurs d'ozone intégrées ;
- La température de surface.

À partir des 4 quantités d'ozone intégrées, des couches 0,005 ; 132,49 ; 222,94 et 478,54 hPa jusqu'à la surface, nous extrapolons un profil complet sur le même nombre de niveaux que les deux autres profils. Nous répartissons les quantités intégrées d'ozone tout au long du profil en respectant le profil typique de l'ozone atmosphérique : assez faible à basse altitude, croissant jusqu'à un maximum dans la basse stratosphère, puis une rapide décroissance.

Les profils de vapeur d'eau sont utilisés en humidité relative. Cette unité, plus facilement intelligible, nous permet de mieux prendre conscience des erreurs commises sur les restitutions effectuées.

Chaque base est donc composée de 10.000 situations avec 3 profils atmosphériques sur les mêmes niveaux verticaux et la température de surface. Ces profils proviennent de restitutions à partir de mesures IASI de l'année 2010. Ils sont donc associés à une date et une géolocalisation précises. À partir de ces informations, nous associons à chaque situation l'émissivité infrarouge hyperspectrale issue des moyennes mensuelles détaillées à la Section 3.3 (page 90). Il s'agit de moyennes mensuelles de l'émissivité restituée à partir de mesures IASI. Elles sont projetées sur une grille "equal-area" à  $0,25^\circ \times 0,25^\circ$  (voir Section 2.1.3, page 49). Nous utilisons l'interpolateur d'émissivité micro-onde TELSEM (Aires et al. 2011b) pour associer les émissivités micro-ondes correspondant aux canaux d'AMSU-A et MHS à chaque situation.

Les spectres d'émissivité infrarouge ne sont pas décompressés et sont utilisés sous la forme de composantes principales de l'ACP (voir Section 2.2.3, page 53). Neuf composantes sont prises en compte dans la restitution car l'algorithme bayésien de restitution de l'émissivité et de la température (voir Chapitre 2, page 37) en donne neuf en sortie.

Nous disposons donc d'une base de données de profils atmosphériques associés à la température de surface et aux émissivités micro-ondes et infrarouges correspondantes. Nous pouvons désormais utiliser le code de transfert radiatif RTTOV pour simuler les mesures des instruments IASI, AMSU-A et MHS pour chaque situation.

À l'image de ce que nous avons fait à la Section 5.1.3 (page 127), chaque fois qu'une situation est prise en compte (dans la phase d'apprentissage, de validation ou de test), les mesures satellites sont bruitées suivant les caractéristiques instrumentales données par le constructeur.

Les données IASI sont compressées à l'aide d'une ACP. Les raisons nous ayant amenés à faire ce choix ont déjà été largement exposées dans la Section 5.1.4.2 (page 130). Nous ne

nous attarderons donc pas plus sur l'ACP ou sur le choix des composantes.

Présenter les résultats obtenus sur les 8 bases de données différentes serait fastidieux et répétitif pour le lecteur. Nous ne présenterons que les résultats concernant la base sur 86 niveaux, car c'est celle qui présente la plus grande variabilité de profils atmosphériques. Les statistiques obtenues sur les autres bases sont très similaires.

Chacune des 10.000 situations de la base de données est donc composée au final de :

- Un profil de température sur 86 niveaux de pression ;
- Un profil de vapeur d'eau sur 86 niveaux de pression ;
- Un profil d'ozone sur 86 niveaux de pression ;
- Une température de surface ;
- L'émissivité infrarouge de la surface à la résolution IASI ;
- L'émissivité micro-onde à la résolution de AMSU-A et MHS ;
- 20 composantes de l'ACP sur le spectre de mesures simulées IASI ;
- Les mesures simulées de AMSU-A et MHS.

Du fait de la linéarité de l'ACP, le bruit de l'instrument IASI peut être calculé directement sur les composantes, pour chaque itération d'apprentissage. Ces bruits instrumentaux permettent d'augmenter artificiellement la base de données de 10.000 situations et de diversifier les différents profils. Cela permet d'augmenter les capacités de généralisation de l'algorithme de restitution.

Nous découpons ici, de la même manière que nous avons procédé à chaque fois qu'il était question de créer des algorithmes statistiques, la base de données en trois bases distinctes :

- **Une base d'apprentissage** de 8.000 situations, qui sera présentée à l'algorithme et qui permettra de le paramétrer ;
- **Une base de validation** de 1.000 situations, qui permet de s'assurer au cours de l'apprentissage que l'algorithme conserve sa capacité à généraliser. C'est sur cette base qu'est calculé le critère de convergence et que sont faits les choix d'architecture des réseaux ;
- **Une base de test** de 1.000 situations, qui servira, une fois l'algorithme paramétré, à calculer l'erreur finale de restitution.

Ces trois bases sont séparées une fois pour toutes et ce sont les trois mêmes bases qui sont utilisées pour chaque configuration de restitution.

### 7.1.2 Méthode de restitution

Les restitutions effectuées ci-après seront toutes effectuées grâce à des réseaux de neurones. Les réseaux de neurones que nous allons utiliser sont des Perceptrons Multi-Couches ([Rumelhart et al. 1986](#)) (semblables à ceux utilisés pour l'interpolateur d'émissivité infrarouge à la Section 2.2.1 (page 50) et pour les restitutions au-dessus des océans à la Section 5.2.4.1, page 147). Nous avons montré dans le Chapitre 5 (page 121) que les réseaux

de neurones étaient les outils les mieux adaptés à ces restitutions. En nous basant sur cet acquis, nous allons donc nous concentrer uniquement sur cette méthode de restitution.

Afin de pouvoir comparer facilement les différents résultats, nous avons choisi de conserver la même architecture neuronale pour toutes les configurations envisagées. Tous les réseaux de neurones que l'on utilise sont constitués d'une seule couche cachée comprenant 60 neurones. De nombreux tests ont été menés afin de déterminer ce nombre, qui semble adapté compte tenu de la configuration minimale et maximale des réseaux que l'on construit. Ces réseaux contiennent tous 86 sorties (nombre de niveaux pour les différents profils) et de 20 à 70 entrées. Ces 60 neurones cachés leur permettront donc d'établir la relation entre ces variables.

Afin de caractériser au mieux les erreurs en sortie des différents réseaux de neurones que nous allons construire, nous utiliserons systématiquement la racine de l'erreur quadratique moyenne (notée RMS). Comme nous l'avons déjà expliqué, cette mesure de l'erreur nous permet de prendre en compte à la fois son biais et son écart-type. Les réseaux de neurones, par construction, génèrent peu de biais, la plus grande partie de la RMS de l'erreur est de la variance. La RMS est alors une mesure proche de l'écart-type de l'erreur.

### 7.1.3 Convergence des différents réseaux de neurones

Nous présentons dans cette partie la stabilité des apprentissages effectués dans les différentes configurations. Nous ne montrons ici que la convergence des réseaux de neurones destinés à la restitution des profils d'humidité relative. Les résultats sur les réseaux restituant le profil de température sont sensiblement identiques, ils ne sont donc pas présentés ici.

La Figure 7.2 présente les convergences des différents réseaux de neurones. Chaque itération correspond à une présentation de la base d'apprentissage au réseau de neurones. À chaque itération, le bruit instrumental des trois capteurs est tiré aléatoirement en suivant la loi gaussienne définie par le constructeur (voir Section 5.1.3, page 127).

Les valeurs indiquées en ordonnées correspondent à la moyenne de la RMS de l'erreur sur tout le profil de l'humidité relative. On retrouve des valeurs relativement faibles car l'erreur est très faible dans les hautes couches de l'atmosphère. Les erreurs sont calculées sur la base de validation au cours de l'apprentissage. Les représentations du profil d'erreur sur les figures suivantes correspondent, quant à elles, à des moyennes de l'erreur sur la base de test avec le réseau paramétré à l'issue de la dernière itération.

Le codage couleur utilisé ici sera conservé par la suite pour plus de clarté pour identifier les différentes configurations. Ici, seuls les réseaux restituant le profil d'humidité relative sont présentés. Le codage couleur est le même pour ceux restituant le profil de température. Les couleurs correspondent aux entrées des réseaux :

- **Bleu foncé** : les mesures de IASI;

- **Bleu clair** : les mesures de IASI et les informations de surface (température et émissivité infrarouge) ;
- **Vert foncé** : les mesures de AMSU-A et MHS ;
- **Vert clair** : les mesures de AMSU-A et MHS et les informations de surface (température et émissivité micro-onde) ;
- **Rouge** : les mesures de IASI, AMSU-A et MHS ;
- **Jaune** : les mesures de IASI, AMSU-A et MHS et les informations de surface (température et émissivité infrarouge et micro-onde).

On peut constater ici que les informations de surface ont amélioré les différentes restitutions. Dans chaque configuration, la prise en compte de l'émissivité et de la température de surface permet à la courbe plus claire d'atteindre des valeurs inférieures. On note aussi que les courbes rouge et jaune (correspondant à l'utilisation des trois capteurs simultanément) sont en dessous des autres courbes. C'est l'effet de la synergie. Il semblerait que l'apport de la synergie soit comparable à l'effet de la prise en compte des caractéristiques de surface.

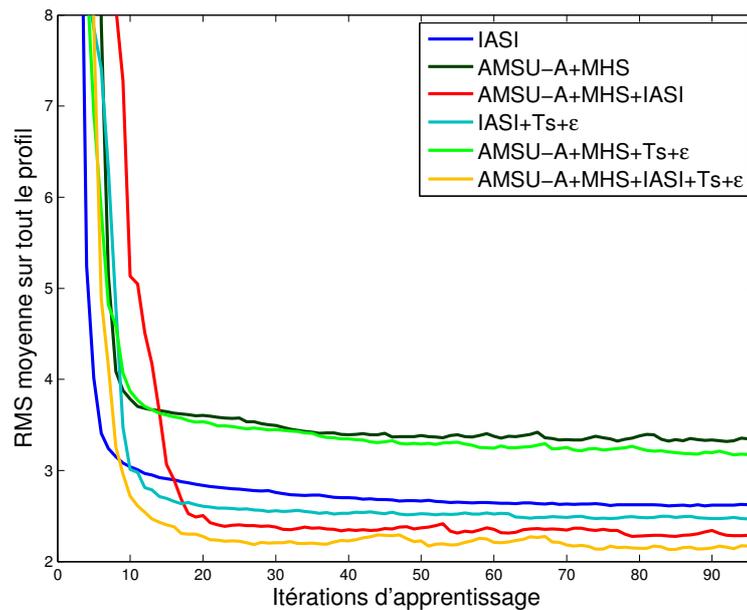


FIGURE 7.2 – Moyenne sur le profil d'humidité relative de la RMS de l'erreur de restitution pour les différentes configurations envisagées, pour chaque itération au cours de l'apprentissage.

Les différentes courbes présentent une structure similaire : une forte décroissance due à un paramétrage grossier des différents poids synaptiques du réseau, suivie d'une lente stabilisation de l'erreur liée à des modifications plus fines des différents paramètres des réseaux. Toutes les courbes se sont bien stabilisées. Aucune ne présente de sur-apprentissage, qui serait visible par une croissance progressive de l'erreur. Ceci est dû à l'utilisation du

bruit instrumental qui est modifié à chaque itération et à un nombre limité de paramètres dans le réseau. Cette légère modification de la base d'apprentissage permet d'éviter que le réseau diminue son erreur de restitution, en codant des relations existant uniquement au sein de la base d'apprentissage. Les réseaux ainsi paramétrés conservent leur capacité de généralisation à des situations qui ne leur ont pas été présentées au cours de l'apprentissage.

Le nombre d'itération d'apprentissage de chaque réseau est, ici, volontairement semblable, afin de pouvoir les comparer sans soucis liés à différents apprentissages. Il est cependant important de noter que le temps d'itération de chaque réseau dépend du nombre de ses entrées (et également de ses sorties, mais ici ce nombre ne change pas). Le calcul de l'erreur et sa rétropropagation au sein du réseau est plus long pour un réseau avec plus d'entrées. Si l'apprentissage d'un réseau de neurones plus complexe est plus long, lors de son utilisation en configuration de restitution, le temps de calcul des sorties est sensiblement équivalent. Il n'est donc pas problématique de construire un réseau relativement complexe.

## 7.2 Apports de l'émissivité et de la température de surface

Dans cette première partie de l'étude, nous nous concentrons sur l'apport des informations de surface. Nous utilisons en parallèle deux configurations d'entrée différentes : une avec uniquement les mesures satellites, et l'autre en ajoutant la température et les émissivités de surface. Cette étude sera menée pour les instruments micro-ondes et infrarouges séparément. Nous construisons à chaque fois deux réseaux de neurones : un pour restituer le profil de température et l'autre pour restituer le profil d'humidité relative.

### 7.2.1 Dans le micro-onde

Nous considérons ici les deux instruments micro-ondes simultanément. Nous avons déjà montré précédemment l'importante synergie entre ces deux capteurs (voir Chapitre 5, page 121). Les réseaux de neurones que l'on construit ont 20 entrées (dans le cas avec seulement les mesures des instruments) ou 41 entrées (en y ajoutant les 20 émissivités et la température de surface).

La partie gauche de la Figure 7.3 présente les RMS des erreurs de restitution du profil de température en utilisant les mesures de AMSU-A et MHS. La courbe verte foncée correspond à la configuration où on ne prend en compte que les mesures des instruments, et la courbe verte claire correspond au cas où on prend également en compte l'émissivité et la température de la surface.

La courbe verte sombre nous indique que la restitution de la température en utilisant AMSU-A et MHS a une erreur moyenne proche de 1,5 K. Cette erreur est plus importante dans les hautes couches de l'atmosphère, 7 à 8 K, ainsi que proche de la surface, 2,5 K. L'erreur de restitution plus élevée proche de la surface est liée aux diverses inversions possibles du profil qui compliquent l'inversion. La plus grande variabilité de la température dans les

basses couches implique une restitution moins précise. Dans les hautes couches de l'atmosphère, la faible densité de molécules diminue l'intensité du rayonnement de ces couches. La détection de ce rayonnement, combiné aux signaux venant d'autres couches de l'atmosphère implique une erreur de restitution plus forte dans les hautes couches de l'atmosphère.

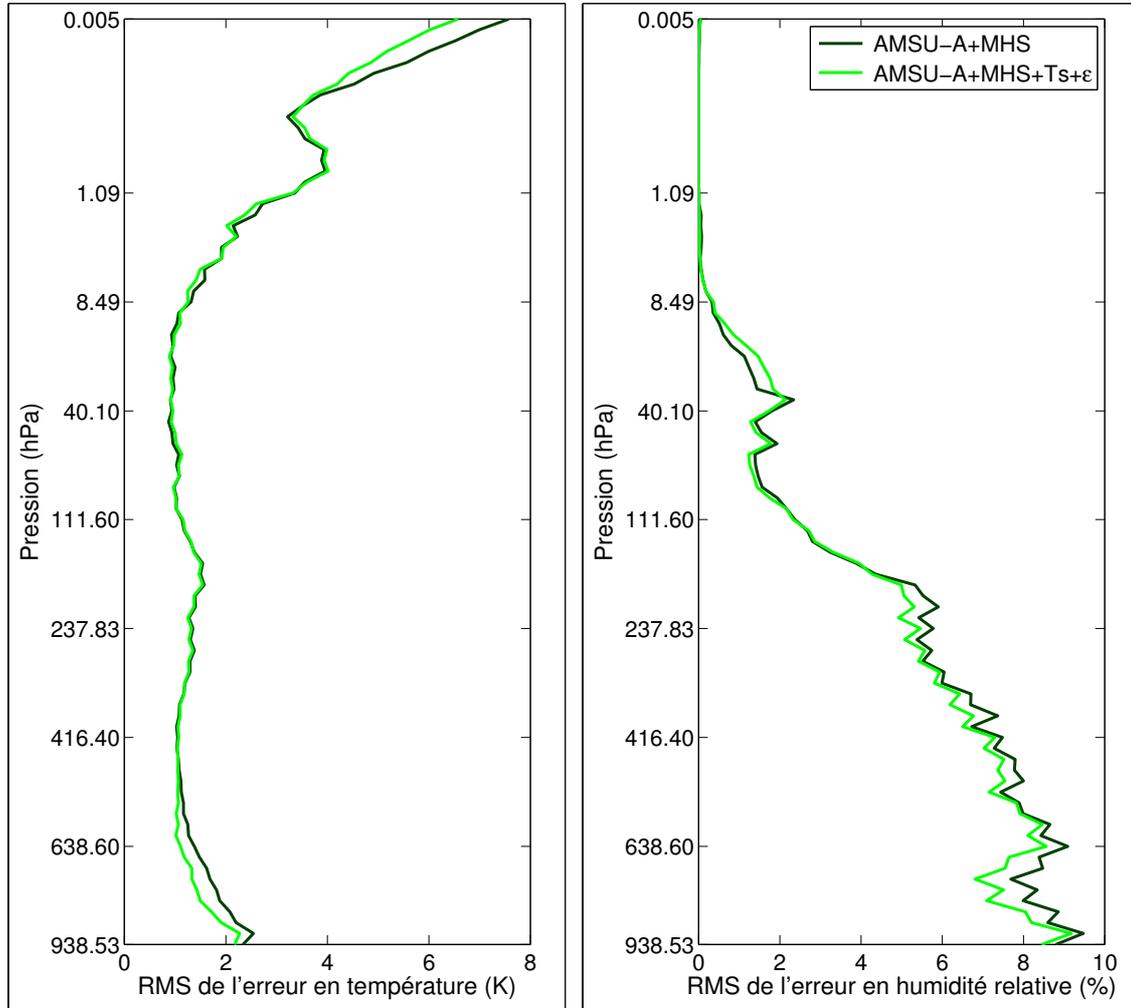


FIGURE 7.3 – Partie gauche : RMS de l'erreur de restitution du profil de température en utilisant AMSU-A et MHS (en vert sombre) et en utilisant également les informations de surface (en vert clair). Partie droite : représentation similaire pour le profil d'humidité relative.

L'ajout de l'émissivité et de la température de surface en entrée de l'algorithme ne modifie pas les statistiques obtenues au milieu de l'atmosphère. On retrouve par contre, comme on s'y attendait, une amélioration de la restitution du profil dans les basses couches de l'atmosphère. L'erreur est désormais plus proche de 2 K.

De façon plus surprenante, la prise en compte des informations de surface a également amélioré la restitution du profil dans les plus hautes couches de l'atmosphère. Par construc-

tion des profils, qui viennent des restitutions IASI, il existe une certaine corrélation entre les plus hautes couches de l’atmosphère et celles proches de la surface. La faible sensibilité de IASI aux hautes couches de l’atmosphère rend compliquée la détermination de la température et de l’humidité dans ces couches. Afin de contraindre la restitution, il est nécessaire d’extrapoler le profil, ce qui introduit ces corrélations entre le haut du profil et le bas.

La partie droite de la Figure 7.3 présente les statistiques de la RMS de l’erreur de restitution du profil d’humidité relative. L’erreur d’estimation de l’humidité relative dans les basses couches de l’atmosphère avoisine les 6 à 9 %. Plus haut dans l’atmosphère, cette erreur décroît jusqu’à atteindre des valeurs inférieures à 1 %. Cette décroissance est liée à un faible taux d’humidité dans la haute atmosphère. L’utilisation de l’humidité relative en lieu et place de l’humidité absolue permet de comparer plus facilement la quantité de vapeur d’eau à des niveaux de pression différents. Cependant, les plus hautes couches de l’atmosphère restent très sèches. Des valeurs de l’humidité relative faibles impliquent une erreur de restitution plus faible. De plus, comme nous l’avons vu à la Section 6.3.1.2 (page 166), les couches les plus hautes de l’atmosphère ont une variabilité de l’humidité faible voire nulle. Cette faible variabilité facilite la restitution et diminue l’erreur sur les hautes couches.

L’ajout en entrée du réseau de neurones des émissivités et de la température de surface (courbe verte claire) permet de réduire l’erreur de restitution dans les basses couches. On peut remarquer un gain allant jusqu’à 0,5 %, à 700 hPa. Dans les plus hautes couches, l’erreur de restitution est déjà faible, elle peut donc difficilement être diminuée. On note cependant une légère augmentation de l’erreur aux alentours de 30 hPa. Cette structure reste marginale et est dans une partie de l’atmosphère où l’humidité est très faible<sup>1</sup>.

Combiner les mesures d’AMSU-A et MHS aux émissivités correspondantes et à la température de surface a donc permis de réduire les erreurs de restitution, à la fois du profil de température et du profil de vapeur d’eau.

## 7.2.2 Dans l’infrarouge

Nous avons mené une étude similaire dans l’infrarouge. Pour les restitutions du profil de température ou de vapeur d’eau, deux réseaux de neurones sont construits. L’un a comme entrées 20 composantes de l’ACP sur le spectre de températures de brillance de IASI, l’autre ajoute à ces 20 entrées, 9 composantes de l’ACP sur le spectre d’émissivité issues des moyennes mensuelles de restitution et la température de surface.

La partie gauche de la Figure 7.4 présente les statistiques de la RMS de l’erreur sur la restitution du profil de température en utilisant uniquement les données IASI (en bleu

---

1. On peut constater sur ces deux courbes des oscillations des erreurs de restitution au long du profil. Ces oscillations sont dues à l’indépendance entre les sorties du réseau de neurones. Du fait de la complexité de la restitution du profil de vapeur d’eau, le réseau de neurones converge difficilement (voir Section 7.1.3, page 188). Il s’agit de faire un compromis entre les erreurs sur les différentes sorties du réseau. Diminuer l’erreur sur une couche peut signifier de l’augmenter sur une autre. On retrouve de telles oscillations sur tous les profils d’erreur de restitution de la vapeur d’eau. Ces oscillations sont également présentes sur les erreurs de restitution de la température mais avec une amplitude plus faible. Ils sont donc moins visibles.

sombre) et en utilisant également l'émissivité et la température de surface (en bleu clair).

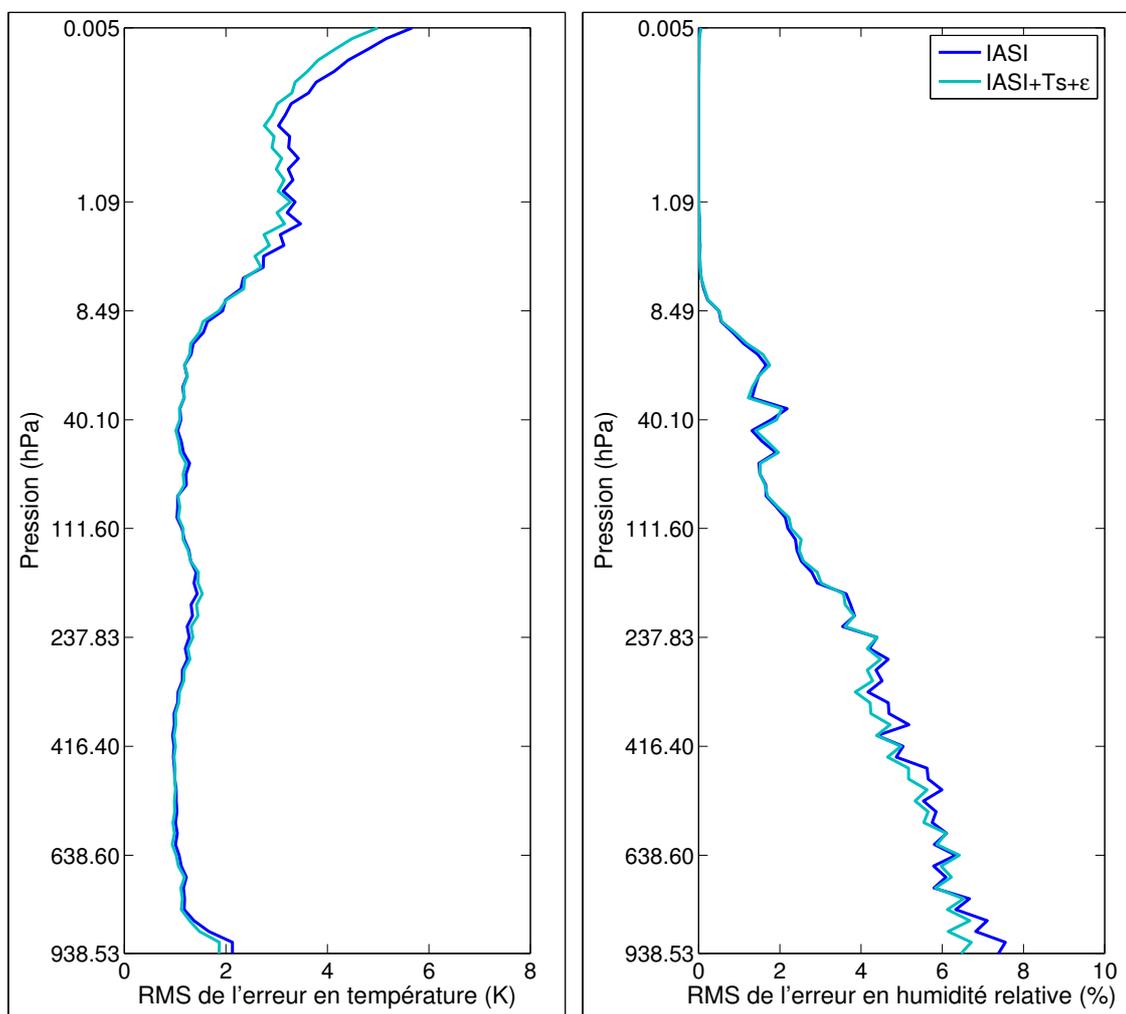


FIGURE 7.4 – Partie gauche : RMS de l'erreur de restitution du profil de température en utilisant IASI (en bleu sombre) et en utilisant également les informations de surface (en bleu clair). Partie droite : représentation similaire pour le profil d'humidité relative.

La structure générale du profil d'erreur est assez semblable à ce qu'on obtenait en utilisant les instruments micro-onde, avec des valeurs plus élevées dans la haute atmosphère et proche de la surface. L'erreur de restitution est assez proche (en valeur) des résultats obtenus avec les instruments micro-ondes, sauf dans la haute atmosphère où elle semble être plus faible.

L'ajout en entrée du réseau de neurones des neuf composantes de l'émissivité et de la température de surface (courbe bleue claire) diminue l'erreur de restitution. Comme avec le micro-onde, les différences notables sont dans le bas de l'atmosphère et dans les plus hautes couches.

Ici encore, l'amélioration de la restitution dans les plus hautes couches est probablement

due à un artifice statistique. Le gain dans les basses couches peut paraître faible : 0,2 K par rapport à une erreur de 2 K. Mais cela signifie une diminution de la RMS de l'erreur de 10 % dans les plus basses couches qui sont les plus compliquées à restituer et les plus importantes pour la météorologie.

La partie droite de la Figure 7.4 présente les statistiques de la RMS de l'erreur de restitution du profil d'humidité relative en utilisant en entrée du réseau de neurones uniquement IASI (en bleu sombre) ou IASI, la température de surface et les émissivités correspondantes (en bleu clair).

On retrouve, comme avec les instruments micro-ondes, un profil d'erreur avec une valeur élevée proche de la surface et une décroissance jusqu'à des valeurs très faibles dans le haut de l'atmosphère. Les erreurs sont ici nettement plus faibles qu'avec les instruments micro-ondes. L'erreur proche de la surface est de l'ordre de 7 %, contre 9 % précédemment.

On peut constater que la prise en compte des informations de surface a permis de diminuer l'erreur de restitution tout au long du profil. Cet apport est plus important plus proche de la surface. L'erreur sur la plus basse couche diminue de 1 %, soit plus de 10 % en erreur relative.

On a montré que, quelques soient les instruments considérés, il faut utiliser les émissivités et la température de la surface en entrée de l'algorithme de restitution. Leur utilisation diminue l'erreur de restitution dans les plus basses couches de l'atmosphère, qui sont les couches les plus compliquées à restituer au-dessus des continents, mais aussi les plus importantes. Des résultats similaires sont observés en assimilation dans les NWP. Pavelin and Candy (2013) ont montré que la prise en compte des émissivités infrarouges permet l'assimilation des canaux sensibles à la surface. Ceci entraîne une diminution de l'erreur sur les basses couches du profil de température et dans le milieu de la troposphère pour le profil d'humidité et vient confirmer les résultats que nous obtenons ici.

### 7.3 Synergie infrarouge et micro-onde

Nous allons désormais étudier la synergie entre les différents instruments de la plateforme MetOp : AMSU-A, MHS et IASI. La méthode que nous utilisons est semblable au travail que nous avons effectué au Chapitre 5 (page 121). Nous considérons ici, comme dans la section précédente, les deux instruments micro-ondes, AMSU-A et MHS de façon conjointe. Nous allons donc comparer des restitutions de profil de température ou de vapeur d'eau en utilisant AMSU-A et MHS, seulement IASI ou les trois instruments conjointement. Dans un premier temps ce travail sera effectué sans prendre en compte les informations de surface. Nous étudierons ensuite l'apport des informations de surface à la restitution synergique. Nous utilisons ici le même codage couleur que précédemment afin de faciliter la lecture.

## 7.3.1 Sans informations de surface

La partie gauche de la Figure 7.5 présente les statistiques de la RMS de l'erreur de restitution du profil de température en utilisant les mesures de IASI (courbe bleue), celles de AMSU-A et MHS (courbe verte) et celles des trois capteurs simultanément (courbe rouge).

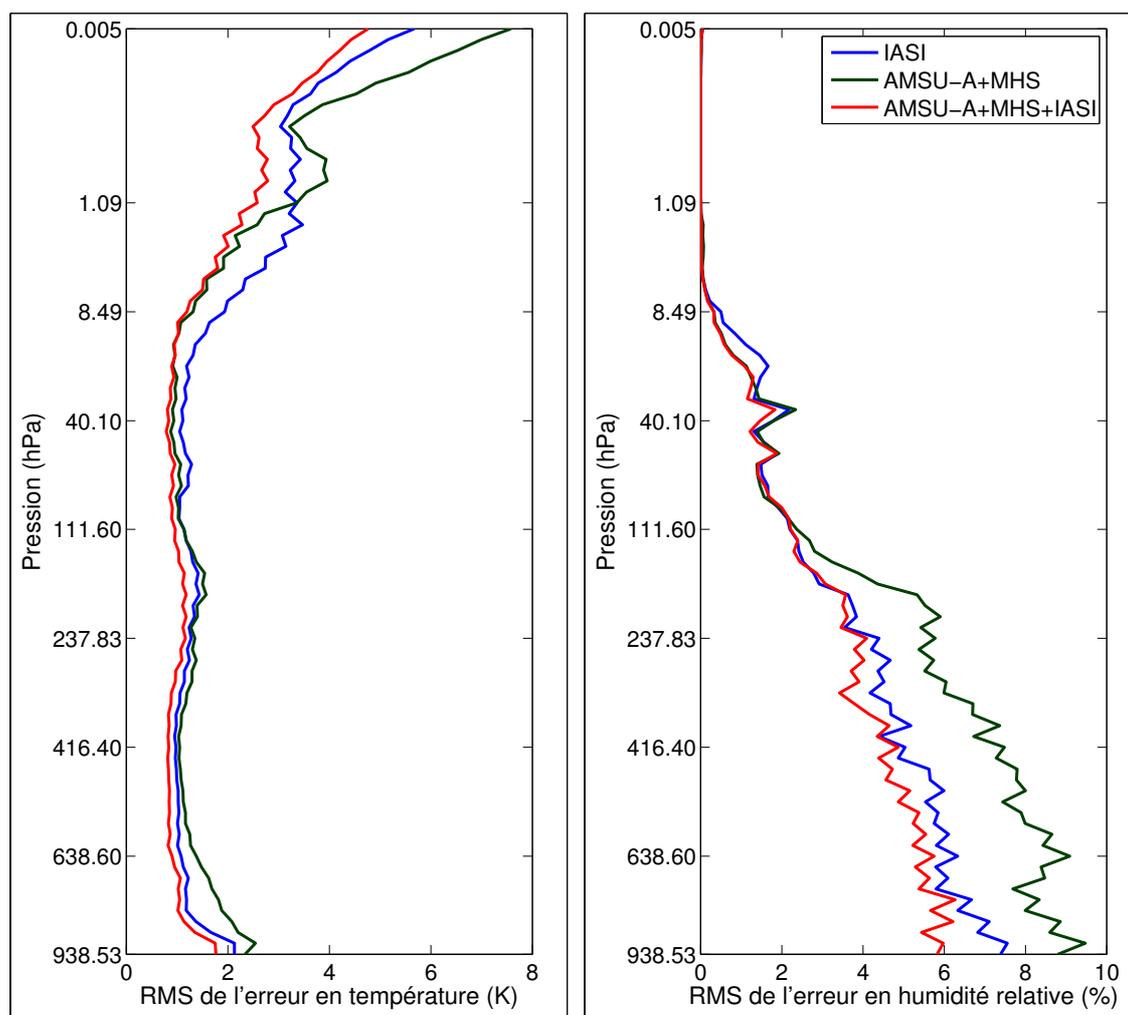


FIGURE 7.5 – Partie gauche : RMS de l'erreur de restitution du profil de température en utilisant IASI (courbe bleue), AMSU-A et MHS (courbe verte) ou les trois instruments conjointement (courbe rouge). Partie droite : représentation similaire pour le profil d'humidité relative.

Dans la partie la plus basse de l'atmosphère, l'erreur de restitution en utilisant IASI est plus faible que l'erreur de restitution utilisant les deux instruments micro-ondes. Entre 240 et 1 hPa, c'est la restitution utilisant AMSU-A et MHS qui présente de meilleures statistiques. Dans les hautes couches, c'est à nouveau pour les restitutions utilisant IASI que l'erreur est plus faible.

On note que tout au long du profil, la restitution utilisant les trois instruments simultanément présente une erreur inférieure aux deux autres restitutions. Il y a donc bien une synergie entre les rayonnements infrarouges et micro-ondes au-dessus des surfaces continentales, pour la restitution du profil de température. L'erreur de restitution de la température est inférieure à 2 K proche de la surface et inférieure à 5 K dans le haut de l'atmosphère.

La partie droite de la Figure 7.5 présente le profil de la RMS d'erreur de restitution du profil d'humidité relative en utilisant IASI (courbe bleue), AMSU-A et MHS (courbe verte) ou les trois capteurs simultanément (courbe rouge). L'erreur de restitution du profil d'humidité en utilisant IASI est plus faible que celle utilisant les instruments micro-ondes, sauf dans le haut de l'atmosphère où il y a très peu de vapeur d'eau. Malgré cela, combiner l'infrarouge et le micro-onde améliore la restitution de façon notable. L'erreur diminue de plus de 1 % dans les basses couches de l'atmosphère.

La synergie existe belle et bien car l'erreur de restitution en utilisant simultanément l'infrarouge et le micro-onde est plus faible que la meilleure des deux restitutions séparées, et ce pour le profil de température ou d'humidité relative. Cette synergie entre les instruments micro-ondes et infrarouges de la plateforme MetOp nous a permis de diminuer l'erreur de restitution sur les profils de température et de vapeur d'eau, au-dessus des continents<sup>2</sup>.

### 7.3.2 Synergie complète entre toutes les données disponibles

Après avoir montré l'importance de la prise en compte des informations de surface (émissivité et température), nous avons montré que la synergie entre IASI, AMSU-A et MHS permet de réduire l'erreur de restitution. Nous allons désormais utiliser la totalité des informations à notre disposition pour restituer les profils de température et de vapeur d'eau.

Nous avons rassemblé sur la partie gauche de la Figure 7.6 les profils de RMS d'erreur pour toutes les différentes configurations, en gardant le même codage couleur que précédemment. Nous avons ajouté sur cette figure la configuration où nous utilisons en entrées les mesures des trois instruments, l'émissivité de la surface dans l'infrarouge et dans le micro-onde et la température de surface (courbe jaune).

---

2. Nous avons montré au Chapitre 5 (page 121) qu'il était compliqué de mettre en valeur la synergie indirecte, du fait de la complexité du problème à restituer. Nous avons fait le choix ici de comparer des réseaux de neurones à l'architecture équivalente (*i.e.*, autant de neurones sur la couche cachée). Ce choix nous permet de comparer facilement l'apport des différentes entrées. Pour mesurer la synergie indirecte, il faut doubler le nombre de neurones sur la couche de sortie (les deux profils, de température et d'humidité relative). En gardant l'architecture que nous avons ici, nous ne parvenons pas à mettre en avant la synergie indirecte. On peut s'attendre à la retrouver si on complexifie d'avantage le réseau en augmentant le nombre de neurones. Cependant, le nombre de neurones sur la couche cachée augmente de façon drastique le temps d'apprentissage du réseau et la mémoire vive nécessaire pour calculer la descente du gradient (voir Annexe B.2.2, page 215). Nous n'avons donc pas mis en place de tels réseaux. La synergie indirecte n'est pas mise en avant dans cette étude. On s'attend cependant à la retrouver grâce à des réseaux de neurones adaptés.

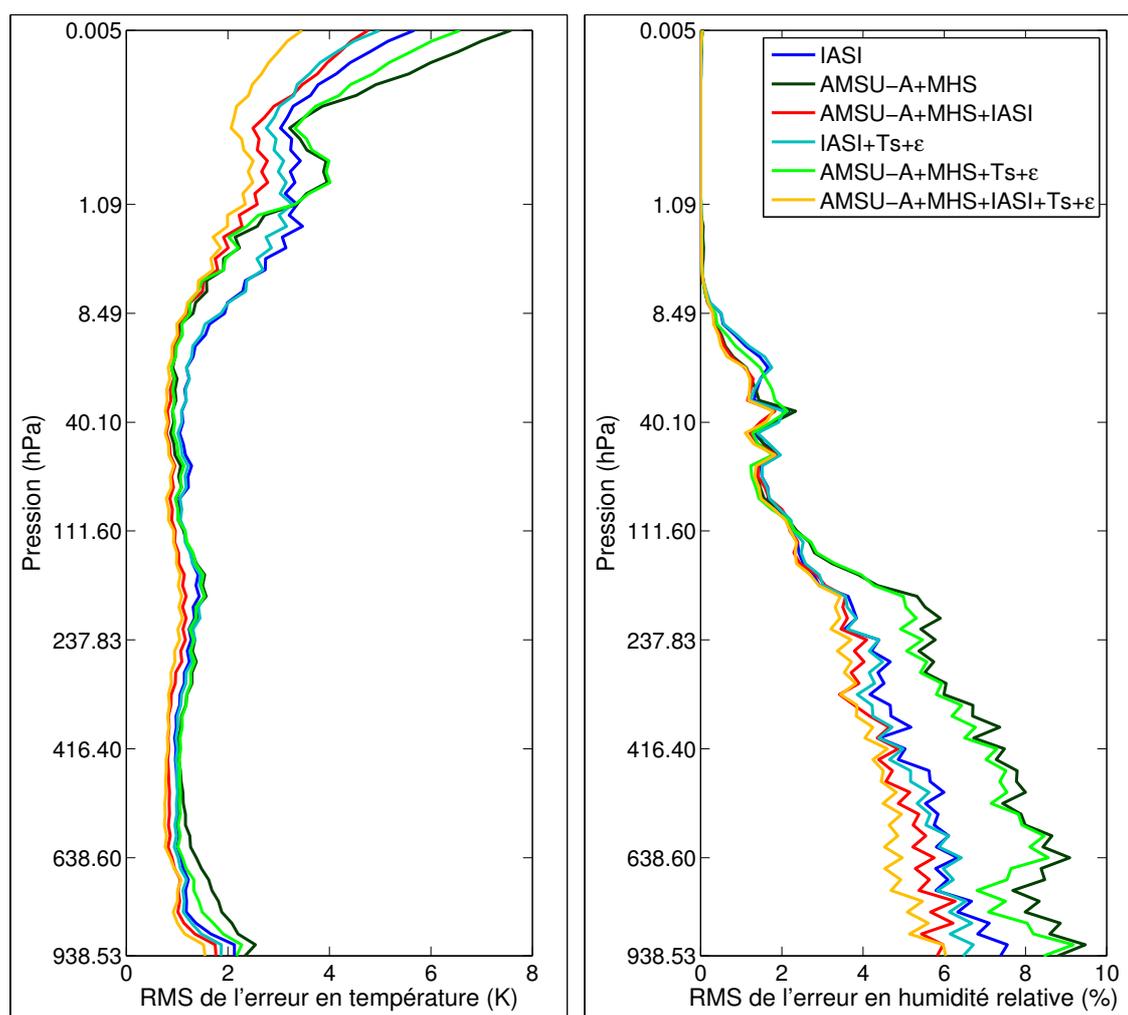


FIGURE 7.6 – Partie gauche : RMS de l’erreur de restitution du profil de température en utilisant IASI (bleu sombre); AMSU-A et MHS (vert sombre); les trois capteurs (rouge); IASI, la température de surface et les émissivités correspondantes (bleu clair); AMSU-A, MHS, la température de surface et les émissivités correspondantes (vert clair); IASI, AMSU-A, MHS, la température de surface et les émissivités correspondantes (jaune). Partie droite : représentation similaire pour le profil d’humidité relative.

L’émissivité utilisée dans les différentes configurations est celle correspondant aux capteurs considérés. Lorsqu’on utilise IASI, on utilise les 9 composantes d’émissivité infrarouge. Lorsqu’on utilise AMSU-A et MHS, on utilise les 20 émissivités correspondant aux 20 canaux. Lorsqu’on utilise les trois capteurs, toutes les émissivités sont utilisées en entrée du réseau de neurones.

On constate que la restitution utilisant la totalité des informations en entrées présente une erreur plus faible que toutes les autres restitutions et ce tout au long du spectre. Une nouvelle fois, l’amélioration est surtout visible dans les plus hautes et les plus basses couches de l’atmosphère.

La partie droite de la Figure 7.6 correspond à la même figure pour la restitution du profil d'humidité relative. Les courbes colorées correspondent aux mêmes configurations d'entrées que celles présentées ci-dessus. Seules les sorties sont modifiées.

La restitution la plus précise est celle utilisant le maximum d'informations en entrée. Une fois encore, la synergie entre les différents instruments ainsi que celle avec les émissivités et la température de surface permet d'améliorer la connaissance du profil d'humidité relative.

## 7.4 Mesure de la synergie infrarouge/micro-onde et de l'apport des informations de surface

Nous avons donc prouvé que la prise en compte des informations de surface permettait de réduire l'erreur de restitution des profils de vapeur d'eau et de température au-dessus des surfaces continentales. Cette diminution a lieu majoritairement dans les basses couches de l'atmosphère.

Nous avons également prouvé que la prise en compte simultanée des différents capteurs diminuait l'erreur de restitution sur la quasi-totalité de la colonne atmosphérique. Il est donc nécessaire de prendre en compte, de façon combinée, toutes les informations, afin de restituer les profils atmosphériques de la façon la plus précise possible.

La Figure 7.7 représente les facteurs de synergie dans le cas où on prend en compte les émissivités et la température de surface (jaune) et dans le cas sans ces informations (rouge). Pour rappel, le facteur de synergie correspond au rapport entre l'erreur de la meilleure restitution indépendante (IASI ou AMSU-A+MHS) et la restitution simultanée. Un facteur de synergie de 100 % correspond au cas où il n'y a pas de synergie. Un facteur supérieur signifie que la synergie a réduit l'erreur de restitution. Les deux facteurs sont calculés par rapport à la meilleure restitution indépendante, sans prise en compte de la surface. La courbe jaune correspond donc à un facteur de synergie et de contribution de surface.

Nous calculons de façon analogue un facteur d'apport des informations de surface correspondant au rapport entre l'erreur de restitution, sans prise en compte des informations de surface, et l'erreur de restitution en les utilisant. Là encore, un facteur de 100 % signifie qu'il n'y a aucun apport et un facteur supérieur signifie qu'il y a une diminution de l'erreur de restitution grâce aux informations de surface. Les courbes correspondantes sont tracées pour IASI (courbe bleu) et pour les instruments micro-ondes (courbe verte). Pour pouvoir plus facilement mesurer les différents apports, les droites noires correspondent donc à la meilleure restitution entre IASI seul et AMSU-A + MHS. Il n'y a alors pas de synergie, c'est pourquoi la droite est à 100%.

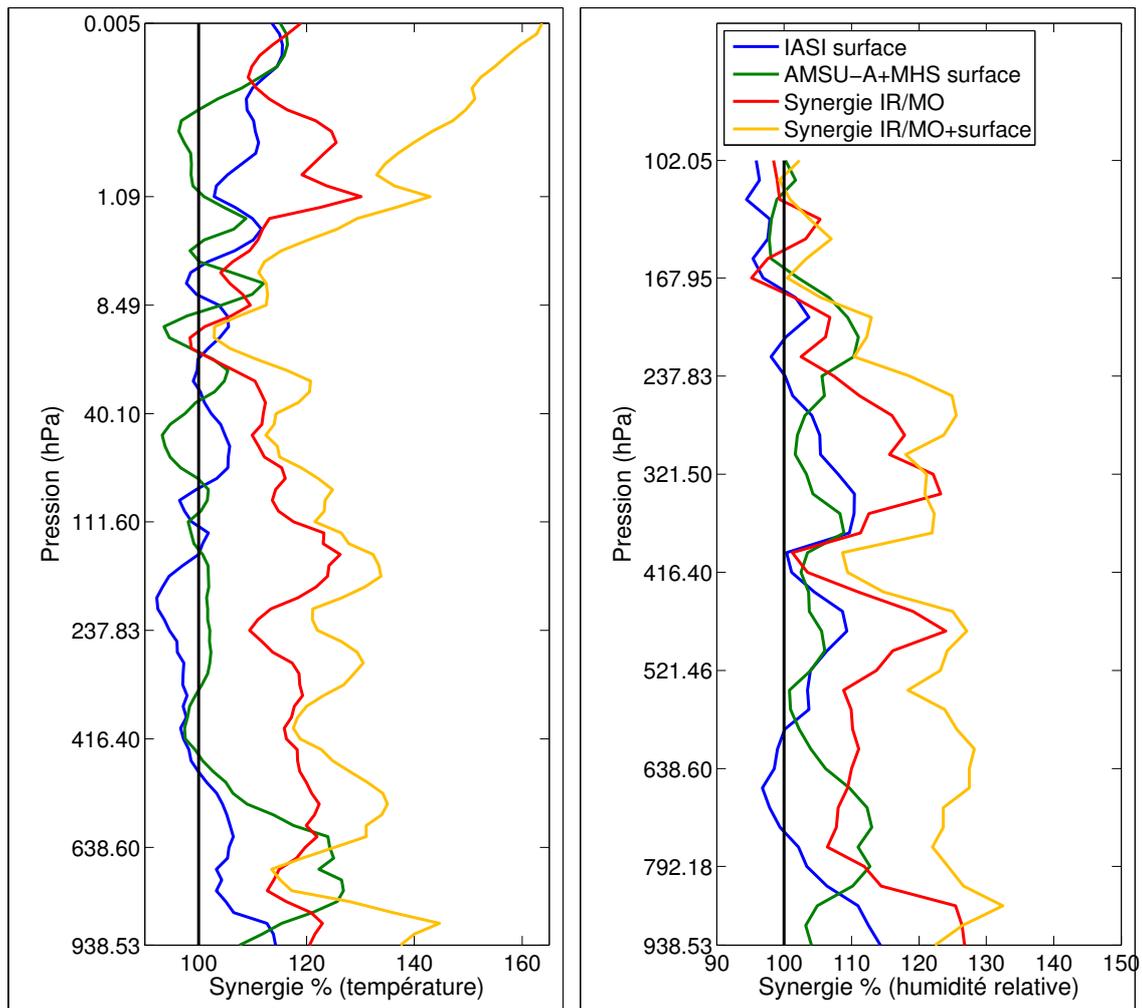


FIGURE 7.7 – Apport des émissivités infrarouges (bleu) et micro-ondes (vert) et facteur de synergie infrarouge/micro-onde (IR/MO) pour les restitutions sans prise en compte des informations de surface (rouge), en les prenant en compte (jaune).

Nous ne représentons que les couches inférieures de l'atmosphère pour le profil de synergie de restitution de l'humidité relative. Les faibles erreurs dans le haut de l'atmosphère entraînent des facteurs élevés qui nuisent à la compréhension de la figure.

On retrouve sur cette figure tous les résultats que l'on a énoncé. Les émissivités et la température de surface améliorent la restitution dans le bas de l'atmosphère, et la synergie améliore la restitution tout le long de l'atmosphère.

Le facteur de synergie est plus élevé, dans le cas où on prend en compte les informations de surface. En effet, plus la quantité d'information est élevée, plus la synergie aura de l'effet. Il est donc nécessaire de prendre en compte toutes les informations disponibles en amont de l'algorithme de restitution, afin de diminuer l'erreur sur l'estimation des profils atmosphériques.

## 7.5 Conclusion

L'utilisation conjointe de tous les instruments satellites est effectuée au sein des algorithmes d'assimilation dans les NWP. L'algorithme que nous avons mis en place ici permet d'obtenir des résultats indépendants de ces modèles. Nous avons montré par ces algorithmes tout l'intérêt de la prise en compte de la synergie, qui entraîne, par exemple, une diminution relative de près de 20 % de l'erreur de restitution de la température dans les bases couches de l'atmosphère. Il serait donc intéressant que la communauté utilise cette synergie pour optimiser les restitutions. La synergie est particulièrement facile à prendre en compte lorsque les capteurs sont à bord de la même plateforme satellite, ce qui est le cas d'AMSU-A, MHS et IASI à bord de MetOp.

Au-dessus des continents, la prise en compte des émissivités de surface permet de diminuer l'erreur de restitution des profils atmosphériques. Si ce résultat était déjà relativement bien connu pour le sondage dans le micro-onde, nous avons montré ici que c'est également le cas dans l'infrarouge. De plus, la prise en compte simultanée des émissivités dans les deux domaines spectraux permet d'affiner les restitutions.

Nous avons démontré ici l'intérêt d'utiliser toutes les informations simultanément pour restituer les profils atmosphériques. La synergie indirecte (entre les variables à restituer) n'a pas été mise évidence. Il serait néanmoins intéressant de développer des schémas d'inversion restituant simultanément les profils de température et de vapeur, ainsi que la température et les émissivités de surface. Un tel schéma profiterait des corrélations entre les variables à restituer.

Au cours de cette étude, nous n'avons utilisé que des méthodes statistiques pour restituer les profils atmosphériques. Les résultats que nous avons obtenus sont très proches des résultats obtenus avec des méthodes d'assimilation (au sein des NWP). Ceci vient justifier notre démarche et appuyer plus encore notre conclusion : il faut prendre en compte simultanément les différentes sources d'information pour améliorer la restitution.

# CONCLUSION ET PERSPECTIVES

Ce travail de thèse peut être dissocié en deux parties. Il y a, d'une part, des avancées techniques très innovantes, comme des algorithmes d'inversion des données satellites, des outils d'échantillonnage de bases hétérogènes de grande dimension... D'autre part, de nouvelles applications sont traitées, comme la restitution conjointe de la température et de l'émissivité infrarouge de surface, la restitution de profils atmosphériques au-dessus des continents ou la fusion d'observations infrarouges et micro-ondes.

Ces deux aspects de la télédétection satellite moderne doivent évoluer en parallèle. Nous avons montré au cours de cette thèse à quel point il est essentiel d'améliorer les techniques utilisées et d'étendre ainsi les domaines pour lesquels les observations satellites peuvent être utiles.

Cette thèse a pour objectif de contribuer à une meilleure exploitation des observations satellites pour l'étude de l'atmosphère, des surfaces et de la climatologie globale.

## Compression et échantillonnage des bases de données

### Conclusion

Au cours de cette étude, nous avons dû faire face à de nombreuses limitations techniques, tant du fait de la grande quantité de données à notre disposition que du nombre élevés de variables sur lesquelles nous avons travaillé. Cette tendance pour le traitement de données satellites toujours plus volumineuses devrait s'accroître dans le futur. Les instruments sont de plus en plus résolus, spectralement et spatialement (on peut, par exemple, citer SWOT (Surface Water and Ocean Topography), un satellite envisagé par la NASA et le CNES qui accumule plusieurs TerraBytes de données par jour). Il est nécessaire de faire évoluer la télédétection satellite actuelle face à ces nouvelles données. Nous avons donc développé des outils capables de diminuer la taille et le nombre de ces données. Parmi les nombreux outils de traitement de données que nous avons présentés, deux sont particulièrement utiles, dans des contextes différents.

Pour le premier, nous avons montré les différents avantages que présente l'utilisation de l'Analyse en Composantes Principales (ACP), particulièrement pour l'utilisation de capteurs hyperspectraux (ici nous avons utilisé IASI). L'ACP nous a permis de compresser et de traiter plus rapidement un grand volume de données. Cette technique de compression permet de conserver, avec un nombre limité de variables, un maximum de la variabilité totale. Grâce à cet algorithme, nous avons pu obtenir des émissivités hyperspectrales et utiliser la totalité

des canaux IASI pour restituer des profils atmosphériques.

Le deuxième outil développé pour réduire la taille des bases de données utilisées est un algorithme d'échantillonnage. Parmi ces algorithmes, nous en avons distingué deux types. Ceux dits "statistiques" qui conservent la variabilité naturelle des différentes variables au sien de la base, et ceux dit "uniformes" qui cherchent au contraire à représenter toutes les situations, en particulier les situations extrêmes, dans des proportions uniformes. Si la méthode des  $k$ -moyennes présentée pour effectuer un échantillonnage statistique est relativement connue, nous avons conçu une méthode d'échantillonnage "uniforme" innovante et performante. L'échantillonnage par entropie nous permet de nous affranchir des problèmes d'hétérogénéité entre les variables et de bien représenter les situations, même les plus extrêmes (Paul and Aires 2013a).

## Perspectives

Ces bases de données échantillonnées sont utilisées, dans cette étude, pour faire de l'apprentissage de méthodes statistiques d'inversion, mais elles peuvent être ré-utilisées dans différents contextes. Elles peuvent, par exemple, servir de première ébauche pour des méthodes d'inversion itératives, ou pour des études de contenues en information si l'on veut tester diverses configurations instrumentales sur des bases atmosphériques limitées mais le plus variées possible. Ces bases peuvent également servir à la même application, c'est-à-dire l'apprentissage d'algorithme d'inversion, c'est pourquoi elles seront fournies à la communauté.

L'échantillonnage par entropie est une méthode relativement novatrice qui permet de résoudre de nombreuses limitations de l'échantillonnage. Il est très souple et facilement adaptable à d'autres applications. Cet échantillonnage peut être réutilisé simplement. Il sera proposé à la communauté, directement sous forme de bases échantillonnées ou la méthodologie générale sera présentée afin de l'adapter à un autre usage. Il sera notamment implémenté dans le centre de calcul de l'EUMETSAT afin de leur fournir une base de données atmosphériques uniforme et multivariée.

De telles méthodes sont cruciales aujourd'hui dans tous les domaines, car le volume de données à traiter, lié à l'utilisation croissante des technologies informatiques, ne cesse d'augmenter. Ces problématiques sont bien identifiées et constituent à elles seules un nouveau domaine appelé le *big data*. Il faut être capable d'assimiler le nombre croissant d'informations disponibles. Dans ce cadre, les méthodes présentées ci-dessus sont particulièrement importantes.

## Restitutions de caractéristiques de surface

### Conclusion

Nous avons mis au point, au cours de cette étude, un algorithme de restitution des caractéristiques de surface dans l'infrarouge. La connaissance précise de la température et de l'émissivité de la surface est nécessaire pour caractériser le rayonnement du sol et donc pouvoir mieux interpréter les mesures effectuées par les capteurs à bord des satellites. En utilisant le savoir-faire développé par les nombreuses équipes travaillant sur les émissivités hyperspectrales infrarouges, nous avons mis au point un algorithme de restitution analytique basé sur un estimateur bayésien qui utilise une première ébauche issue de restitutions indépendantes, à partir de MODIS (Paul et al. 2012). Ces restitutions de surface ont été validées par de multiples méthodes. Nous avons donc décidé d'en faire un algorithme opérationnel fonctionnant en pseudo-temps réel. Une chaîne opérationnelle a été développée, et une base de données de restitutions de la température et de l'émissivité hyperspectrale infrarouge de la surface de 2007 à nos jours a été construite. Du fait de l'utilisation de l'estimateur bayésien, la base de données d'émissivités restituées contient une estimation des matrices de covariance d'erreur des inversion.

### Perspectives

Dans le cadre d'un partenariat avec EUMETSAT, cet algorithme sera installé directement sur le centre de calcul de Darmstadt. Cette chaîne opérationnelle permettra d'obtenir une plus longue série temporelle d'émissivité de surface, particulièrement intéressante à étudier, comme nous l'avons montré au cours de cette étude. L'émissivité est une variable relativement directe pour l'observation satellite. C'est donc un bon moyen de caractériser les surfaces continentales et de suivre leurs évolution. La connaissance globale de l'émissivité hyperspectrale infrarouge de la surface permet, par exemple, de connaître la composition du sol ou de suivre l'évolution du couvert végétal. Allier les émissivités infrarouges à d'autres informations comme les émissivités micro-ondes, ou des indicateurs de végétation issus du visible, permettrait d'obtenir une description relativement complète des surfaces continentales, qui sont au coeur des problématiques de certaines populations.

Après avoir validé nos restitutions de surface avec d'autres données indépendantes, il faut les comparer aux autres restitutions de surface à partir des mesures IASI. Pour cela, un programme d'intercomparaisons des différentes bases de données est actuellement mis en place. L'objectif est d'identifier et d'analyser les différences entre ces émissivités de surface, tant spatialement que spectralement ou temporellement.

D'un point de vue technique, l'inversion bayésienne correspond à ce qui est fait dans les centres opérationnels. L'algorithme mis en place ici permet donc de se rapprocher de ce qui est fait au sein de ces centres (l'estimation optimale), en restant indépendants de leurs modèles. De telles inversions sont très intéressantes pour la communauté car elles

permettront notamment l'assimilation des données satellites sensibles au bas de l'atmosphère (qui est la partie de l'atmosphère la plus importante pour la population) au sein des centres de prévision numérique du temps.

## Synergie infrarouge et micro-onde pour la restitution de profils atmosphériques

### Conclusion

La restitution de profils de température ou de vapeur d'eau au-dessus des océans en utilisant IASI, AMSU-A et MHS simultanément résulte en une erreur plus faible que si ces instruments sont considérés séparément. L'utilisation combinée des mesures infrarouges et micro-ondes nous a donc permis de diminuer l'erreur de restitution des profils atmosphériques (Aires et al. 2011a). Nous avons montré que, dans ce contexte, l'utilisation de réseaux de neurones était préférable, car elle permet d'obtenir une erreur plus faible, tant sur le profil de température que sur celui d'humidité relative. La même méthodologie au-dessus des continents a été mise en place. Une fois encore, l'utilisation simultanée des différents capteurs a conduit à une erreur de restitution plus faible que s'ils sont utilisés séparément. De plus, adjoindre à ces mesures les caractéristiques de surface précédemment restituées (émissivité et température) permet de réduire plus encore l'erreur de restitution des profils de température et d'humidité relative (Paul and Aires 2013b).

La synergie des instruments infrarouges et micro-ondes est déjà utilisée au sein des centres de prévision dans les algorithmes d'assimilation. Nous avons montré ici qu'elle pouvait être également exploitée dans le cadre d'inversion satellites pures. Nous avons démontré, de plus, que les émissivités infrarouges et micro-ondes permettent une bonne caractérisation des surfaces et donc une diminution de l'erreur sur l'estimation des paramètres atmosphériques, notamment dans les basses couches de l'atmosphère, particulièrement variables mais importantes pour la population.

On a donc montré qu'il était possible de mieux exploiter les observations satellites existantes, et ainsi mieux décrire l'atmosphère. Alors que la température et la vapeur d'eau sont deux des variables les plus importantes pour la caractérisation de l'atmosphère, et ont donc été largement étudiées par la communauté, nous avons mis en lumière des pistes d'amélioration de leur restitution.

### Perspectives

Nous avons considéré, au cours de cette étude, des restitutions en ciel clair uniquement. La complémentarité évidente entre l'infrarouge et le micro-onde rend l'étude de cette synergie par temps nuageux intéressante. Si la synergie entre les rayonnements est présente dans le cas le plus simple (*i.e.*, sans nuages), elle devrait être encore plus forte pour les restitutions

plus complexes. Par exemple, le sondage dans l'infrarouge donne accès à la température et à la pression au sommet du nuage alors que le sondage dans le micro-onde permet de pénétrer dans le nuage. Cette complémentarité entre les observations nous incite à penser que l'on retrouvera une synergie plus importante. La synergie peut également être utilisée pour restituer d'autres variables, comme l'ozone, les nuages, les propriétés de surface...

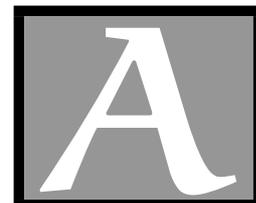
Nous avons été à la pointe de la recherche en initiant, il y a quelques années, ce travail sur la synergie infrarouge et micro-onde pour la restitution de profils atmosphériques au-dessus des océans, dans le cadre d'un projet avec l'ESA (European Space Agency). Ces développements sont maintenant mûrs pour être utilisés dans des contextes opérationnels. L'EUMETSAT a montré récemment un intérêt pour cette thématique de synergie infrarouge et micro-onde, une implémentation opérationnelle est même envisagée.

Ce travail a été effectué sur des capteurs existants (AMSU-A, MHS et IASI sur la plateforme MetOp). Au vu des résultats que nous avons obtenus et des améliorations potentielles, il nous semble que la définition de nouveaux instruments serait plus pertinente si elle était réalisée globalement pour tous les instruments d'une même plateforme, plutôt qu'indépendamment suivant chaque instrument. L'approche que nous avons développée et la mise en valeur de la synergie entre les différents rayonnements justifie ce choix.



---

# ACRONYMES



Liste des différents acronymes utilisés dans cette étude :

- **A-DCS** : Argos Advanced Data Collection System
- **ACP** : Analyse en Composantes Principales
- **AIRS** : Atmospheric InfraRed Sounder
- **AMSR-E** : Advanced Microwave Scanning Radiometer-EOS
- **AMSU-A** : Advanced Microwave Sounding Unit A
- **ARA** : Analyse du Rayonnement Atmosphérique
- **ASCAT** : Advanced SCATterometer
- **ASTER** : Advanced Spaceborne Thermal Emission and Reflection radiometer
- **ATOVS** : Advanced TIROS Operational Vertical Sounder
- **AVHRR** : Advanced Very High Resolution Radiometer
- **CrIS** : Cross-Track Infrared Sounder
- **CRTM** : Community Radiative Transfer Model
- **EADS** : European Aeronautic Defence and Space company
- **EOS** : Earth Observing System
- **ECMWF** : European Center fo Medium-range Weather Forecast
- **EOF** : Empirical Orthogonal Function
- **EPS** : EUMETSAT Polar System
- **ESA** : European Space Agency
- **EUMETSAT** : EUropean organization for the exploitation of METeorological SA-Tellites
- **FASTEM** : FAST microwave Emissivity Model
- **FRTM** : Fast Radiative Transfer Model
- **GCM** : Global Circulation Model
- **GOME** : Global Ozone Monitoring Experiment
- **GOS** : Global Observing System

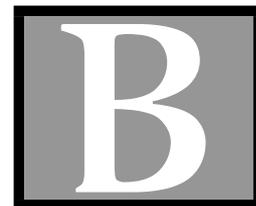
- **GPS** : Global Positioning System
- **GRAS** : Global navigation satellite system Receiver for Atmospheric Sounding
- **HAMSTRAD** : H<sub>2</sub>O Antarctica Microwave Stratospheric and Tropospheric RADiometers
- **HIRS** : High-resolution Infrared Radiation Sounder
- **HSB** : Humidity Sounder for Brazil
- **IASI** : Infrared Atmospheric Sounding Interferometer
- **IGBP-DIS** : International Geosphere Biosphere Program Data and Information System
- **IPSL** : Institut Pierre Simon Laplace
- **ISCCP** : International Satellite Cloud Climatology Project
- **LATMOS** : Laboratoire Atmosphères, Milieux, Observations Spatiales
- **LERMA** : Laboratoire d'Etudes du Rayonnement et de la Matière en Astrophysique
- **LIDAR** : Light Detection And Ranging
- **LMD** : Laboratoire de Météorologie Dynamique
- **LSA/SAF** : Land Surface Analysis/Satellite Applications Facility
- **LUT** : Look-Up Table
- **MADRAS** : Microwave Analysis and Detection of Rain and Atmospheric Structures
- **MBA** : Master of Business Administration
- **MHS** : Microwave Humidity Sounder
- **MODIS** : MODerate-resolution Imaging Spectroradiometer
- **MSG** : Meteosat Second Generation
- **NASA** : National Aeronautics and Space Administration
- **NOAA** : National Oceanic and Atmospheric Administration
- **NWP** : Numerical Weather Prediction center
- **RFI** : Radio Frequency Interference
- **RMS** : Root Mean Square, on utilise l'acronyme anglais pour la racine de l'erreur quadratique
- **RTTOV** : Radiative Transfer for TOV
- **S&R** : Search and Rescue
- **SAF** : Satellite Applications Facility
- **SEM-2** : Space Environment Monitor
- **SEVIRI** : Spinning Enhanced Visible and InfraRed Imager

- 
- **SMOS** : Soil Moisture and Ocean Salinity
  - **SSM/I** : Special Sensor Microwave/Imager
  - **SSMIS** : Special Sensor Microwave Imager/Sounder
  - **SWOT** : Surface Water Ocean Topography
  - **TELSEM** : a Tool to Estimate Land Surface Emissivities in the Microwave
  - **TIGR** : Thermodynamic Initial Guess Retrieval
  - **TIR** : Transformed InfraRed
  - **TMI** : TRMM Microwave Image
  - **TRMM** : Tropical Rainfall Measuring Mission
  - **TIROS** : Television and InfraRed Observational Satellite
  - **TOVS** : TIROS Operational Vertical Sounder
  - **UCSB** : University of California, Santa Barbara
  - **UIT** : Union Internationale des Télécommunications
  - **VIRS** : Visible and Infrared Scanner



---

# MATHÉMATIQUES



## Sommaire

<b>B.1 L'Analyse en Composantes Principales</b> . . . . .	<b>211</b>
<b>B.2 Les algorithmes d'inversion</b> . . . . .	<b>213</b>
B.2.1 La restitution bayésienne . . . . .	213
B.2.2 Le réseau de neurones . . . . .	215

Cet annexe revient sur différents outils mathématiques qui ont été utilisés au long de cette étude. Ces algorithmes peuvent être distingués en deux catégories :

- les algorithmes de pré-traitement des données
- les algorithmes d'inversion.

Les algorithmes de pré-traitement peuvent être très variés. Il peut s'agir d'algorithmes de calibration qui permettent de contraindre les données suivant un modèle précis ou d'algorithmes de compression qui permettent de réduire la dimension des données. Dans le cadre de cette étude, l'algorithme qui nous intéresse est un algorithme de compression des données. L'objectif est de réduire les dimensions d'une variable complexe. Cette dimension trop élevée alourdit le temps de calcul. Nous utiliserons ici l'analyse en composantes principales afin de conserver le maximum de variabilité possible.

Les algorithmes d'inversion numériques sont utilisés lorsque la fonction à inverser est trop complexe pour pouvoir être inversée de façon analytique. Nous utiliserons au cours de cette étude quatre différentes méthodes afin d'approximer une fonction inverse : la restitution bayésienne, les plus proches voisins, la régression linéaire et le réseau de neurones. Nous approfondissons ici la méthode bayésienne et les réseaux de neurones, les autres méthodes étant relativement communes et déjà présentée au cours de l'étude.

## B.1 L'Analyse en Composantes Principales

L'Analyse en Composantes Principales (ACP) consiste à changer l'espace sur lequel on travaille, on projette, pour cela, les données sur un nouvel espace plus pertinent. Cette projection, si elle est effectuée correctement, permet de diminuer la dimension de la base de données (*i.e.*, diminuer le nombre de variables, qu'il s'agisse de variables géophysiques ou

radiométriques). Elle permet également d'obtenir des variables non corrélées, ce qui conduit à une analyse des facteurs qui expliquent la variabilité des données considérées.

Prenons un échantillon de  $k$  réalisations de  $n$  variables, structuré dans une matrice  $X$  :

$$X = \begin{pmatrix} x_{1,1} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{k,1} & \dots & x_{k,n} \end{pmatrix}$$

Chaque variable aléatoire  $X_n = (x_{1,1}, \dots, x_{k,n})$  a une moyenne  $\overline{X_n}$  et un écart-type  $\sigma_n$ .

Supposons que la matrice  $V$  est centrée, si ce n'est pas le cas il suffit de considérer la matrice  $X - \overline{X_n} \cdot \mathbb{1}$ . On cherche alors une matrice  $M$  de taille  $p \times n$  de projection de  $\mathbb{R}^n$  dans  $\mathbb{R}^p$ , de telle sorte que  $Y_n = M \cdot X_n$  soit une compression de  $V$  et que  $M^T \cdot M = \mathbb{1}$ . Chaque ligne de  $M$  est donc un filtre linéaire sur lequel on projette  $X_n$ .

Soit  $\bar{M}$  une matrice  $n \times n$  telle  $\bar{M}^T \cdot \bar{M} = \mathbb{1}$ . Notons  $\bar{M}_j$  la  $j^{\text{ème}}$  ligne de  $\bar{M}$ . On peut alors prendre comme matrice de compression  $M$  les  $p$  premières lignes de  $\bar{M}$ .  $M$  est donc la sous-matrice  $p \times n$  de  $\bar{M}$  constituées de ses  $p$  premières lignes.

Le critère de qualité de la compression est un critère des moindres carrés sur l'erreur commise sur les  $n - p$  autres filtres soit :

$$\begin{aligned} \frac{1}{2 \cdot k} \sum_{i=1}^k \sum_{j=p+1}^n (\bar{M}_j x_j^i - y_j^i)^2 &= \frac{1}{2 \cdot k} \sum_{i=1}^k \sum_{j=p+1}^n \left( \bar{M}_j (x_j^i - \overline{x_j^i}) \right)^2 \\ &= \frac{1}{2} \sum_{j=p+1}^n (\bar{M}_j^T \cdot C \cdot \bar{M}_j)^2 \end{aligned}$$

Où  $C = \frac{1}{k} \sum_{i=1}^k (X^i - \overline{X^i}) \cdot (X^i - \overline{X^i})^T$  est donc la matrice de covariance de  $X$ .

En notant  $\mu_{ij}$  les multiplicateurs de Lagrange, le critère d'orthogonalité de  $M$  nous donne :

$$\frac{1}{2} \sum_{i=p+1}^n \sum_{j=p+1}^n \mu_{ij} (\bar{M}_i \bar{M}_j^T - \delta_{ij})$$

L'objectif est donc la minimisation du terme  $Trace(\bar{M} \cdot C \bar{M}^T) - Trace(\bar{M} \bar{M}^T - \mathbb{1})$ .  $\bar{M}$  étant symétrique, on peut réécrire cela comme :

$$\begin{aligned} C \cdot \bar{M}^T &= \bar{M} \\ \bar{M} \cdot C \cdot \bar{M}^T &= \mathbb{1} \end{aligned}$$

Chaque ligne de  $\bar{M}$  représente donc un vecteur propre de  $C$ , matrice de covariance de  $X$ . La variance autour de chacun des  $X_n$  autour de chacun de ces nouveaux axes  $M_j$  est la valeur propre associée à ce vecteur propre.

Pour calculer l'ACP d'une base de données, il suffit donc de diagonaliser la matrice de covariance  $C$  des données  $X$ . La base de projection est alors directement donnée par les vecteurs propres de la matrice de covariance. La variabilité expliquée par chaque composante de l'ACP est la valeur propre associée. Trier les vecteurs propres par valeur propre décroissante permet alors de trier les axes de projection de la nouvelle base suivant leur ordre d'importance.

## B.2 Les algorithmes d'inversion

Il existe diverses méthodes d'approximation de fonction en mathématiques. Soit  $M$  un système physique, défini par ses caractéristiques  $f$ , que l'on ne peut observer qu'à partir des mesures  $F$  :

$$M(f) = F \tag{B.1}$$

Il y a alors deux méthodes pour inverser ce problème et retrouver  $f$  à partir de  $F$ . On peut chercher à inverser cette équation de façon locale, et trouver la variable  $f$  qui correspond à une mesure  $F$  donnée. On peut également chercher à inverser le système de façon globale afin d'estimer l'opérateur  $M^{-1}$ . Nous développons dans cette partie deux méthodes distinctes d'inversion d'un tel problème : une localisée, l'estimateur bayésien du maximum *a priori* ; l'autre global, les réseaux de neurones. Nous nous plaçons ici dans le cas où les problèmes sont discrets et le nombre de paramètres est fini, ce qui est toujours le cas dans le domaine de la météorologie (Tarantola 1987; Rodgers 1990, 2000; Aires 1999).

### B.2.1 La restitution bayésienne

Nous cherchons donc à inverser l'équation (B.1). On suppose ici que cette relation est localement linéaire. Elle peut alors s'écrire :

$$M \cdot f = F$$

Nous supposons également que les erreurs sur les mesures  $F$  suivent une loi gaussienne de matrice de covariance  $S_F$  et qu'elles ne sont pas biaisées, on peut alors écrire :

$$P(F|f) = \mathcal{N}\left(M \cdot f, S_F^{-\frac{1}{2}}\right) = \frac{1}{(2\pi)^{\frac{n}{2}} \cdot S_F^{-\frac{1}{2}}} \cdot e^{(-\frac{1}{2} \cdot (F - M \cdot f)^T \cdot S_F^{-1} \cdot (F - M \cdot f))} \tag{B.2}$$

où  $n$  est la dimension de  $F$ .  $P(F|f)$  est la probabilité que la mesure vaille  $F$  lorsque les paramètres valent  $f$ .

Nous supposons enfin que nous avons une première ébauche de la solution notée  $f_g$ . La solution  $f$  suit alors une loi gaussienne autour de cette première ébauche, définie par une

matrice de covariance d'erreur  $S_{fg}$ . On a donc :

$$P(f) = \mathcal{N}\left(f_g, S_{fg}^{-\frac{1}{2}}\right) = \frac{1}{(2\pi)^{\frac{n}{2}} \cdot S_{fg}^{-\frac{1}{2}}} \cdot e^{(-\frac{1}{2} \cdot (f-f_g)^T \cdot S_{fg}^{-1} \cdot (f-f_g))} \quad (\text{B.3})$$

L'inversion bayésienne correspond à chercher  $\hat{f}$  qui maximise  $P(f|F)$ . Or le théorème de Bayes nous dit que :

$$P(f|F) = \frac{P(F|f) \cdot P(f)}{P(F)}$$

En injectant les équations (B.2) et (B.3) dans la formule de Bayes, on obtient :

$$P(f|F) = \frac{\left( \frac{1}{(2\pi)^{\frac{n}{2}} \cdot S_F^{-\frac{1}{2}}} \cdot e^{(-\frac{1}{2} \cdot (F-M \cdot f)^T \cdot S_F^{-1} \cdot (F-M \cdot f))} \right) \cdot \left( \frac{1}{(2\pi)^{\frac{n}{2}} \cdot S_{fg}^{-\frac{1}{2}}} \cdot e^{(-\frac{1}{2} \cdot (f-f_g)^T \cdot S_{fg}^{-1} \cdot (f-f_g))} \right)}{P(F)}$$

que l'on peut réécrire comme :

$$\ln(P(f|F)) = -\frac{1}{2} \left( (F - M \cdot f)^T \cdot S_F^{-1} \cdot (F - M \cdot f) + (f - f_g)^T \cdot S_{fg}^{-1} \cdot (f - f_g) \right) + K \quad (\text{B.4})$$

où  $K$  est une constante. Il s'agit d'une forme quadratique, qu'on peut donc écrire sous la forme :

$$-2 \cdot \ln(P(f|F)) = (f - \hat{f})^T \cdot S^{-1} \cdot (f - \hat{f}) + K'$$

où  $K'$  est une constante. C'est-à-dire que la solution recherchée suit également une loi gaussienne de moyenne  $\hat{f}$  et de matrice de covariance d'erreur  $S$ . En utilisant l'équation (B.4), on peut égaliser les deux termes :

$$(f - \hat{f})^T \cdot S^{-1} \cdot (f - \hat{f}) + K' = (F - M \cdot f)^T \cdot S_F^{-1} \cdot (F - M \cdot f) + (f - f_g)^T \cdot S_{fg}^{-1} \cdot (f - f_g) + K \quad (\text{B.5})$$

On peut alors égaliser dans un premier temps tous les termes quadratique en  $f$  de l'équation (B.5) :

$$f^T \cdot S^{-1} \cdot f = f^T \cdot M^T \cdot S_F^{-1} \cdot M \cdot f + f \cdot S_{fg}^{-1} \cdot f$$

Cette équation est vraie pour tout  $f$ , donc on obtient donc la matrice de covariance d'erreur de la solution  $S$  (qui est notée  $Q$  à la l'équation (2.9), page 65 dans la Section 2.3.1) :

$$Q = S = \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \quad (\text{B.6})$$

En revenant à l'équation (B.5), on peut désormais égaliser les termes linéaires en  $f$  :

$$-\hat{f}^T \cdot S^{-1} \cdot f = -F^T \cdot S_F^{-1} \cdot M \cdot f - f_g^T \cdot S_{fg}^{-1} \cdot f$$

Ici encore, cette équation est vraie pour tout  $f$ , on obtient alors, en prenant la transposée :

$$S^{-1} \cdot \hat{f} = M^T \cdot S_F^{-1} \cdot F + S_{fg}^{-1} \cdot f_g$$

car les matrices de covariance  $S$ ,  $S_F$  et  $S_{fg}$  sont symétriques. En remplaçant  $S$  par son expression à l'équation B.6, on a :

$$\left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right) \cdot \hat{f} = M^T \cdot S_F^{-1} \cdot F + S_{fg}^{-1} \cdot f_g$$

On a alors :

$$\begin{aligned} \hat{f} &= \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \cdot \left( M^T \cdot S_F^{-1} \cdot F + S_{fg}^{-1} \cdot f_g \right) \\ &= \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \cdot \left( M^T \cdot S_F^{-1} \cdot F + \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right) \cdot f_g \right) \\ &\quad - \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \cdot \left( M^T \cdot S_F^{-1} \cdot M \right) \cdot f_g \\ \hat{f} &= f_g + \left( M^T \cdot S_F^{-1} \cdot M + S_{fg}^{-1} \right)^{-1} \cdot M^T \cdot S_F^{-1} \cdot (F - M \cdot f_g) \end{aligned} \quad (\text{B.7})$$

Cette équation nous donne directement l'expression de la solution du problème inverse en fonction de la première ébauche  $f_g$ , du modèle direct  $M$ , de la matrice de covariance d'erreur des observations  $S_F$ , de la matrice de covariance d'erreur sur la première ébauche  $S_{fg}$  et de la mesure  $F$ . Il faut tout de même garder à l'esprit que la solution décrit une loi gaussienne centrée sur  $\hat{f}$ , défini à l'équation (B.7), et dont la matrice de covariance est définie à l'équation (B.6).

Il s'agit ici d'une inversion locale, les calculs matriciels devant être effectués à chaque restitution. Il faut disposer d'une première ébauche de la solution et des matrices de covariance d'erreur de la première ébauche et des mesures.

### B.2.2 Le réseau de neurones

Une présentation des réseaux de neurones est fournie à la Section 2.2.1 (page 50). Nous ne reviendrons pas ici sur la définition d'un réseau de neurones. L'objectif ici est d'expliquer la rétropropagation des erreurs au sein d'un réseau, qui a permis à ces derniers de devenir aujourd'hui une technique très courante et utilisée dans différents domaines (médecine, mathématiques, physique, ingénierie) (Rumelhart et al. 1986).

En reprenant les notations définies à la Section 2.2.1 (page 50), considérons un réseau de neurones à  $n$  sorties, appelées  $y$ . L'apprentissage du réseau va consister à modifier les poids synaptiques des différents neurones de façon à minimiser la somme des erreurs quadratiques (*i.e.*, le critère des moindres carrés) :

$$W_{final} = \min_w \left( \frac{1}{2 \cdot p} \sum_{e=1}^p \sum_{i=1}^n (\hat{y}_i^e(w) - y_i^e)^2 \right)$$

où  $w$  correspond aux poids synaptiques,  $p$  correspond au nombre de situations dans la base d'apprentissage,  $n$  correspond au nombre de neurones sur la dernière couche du réseau.  $\hat{y}_i^e(w)$  correspond à la  $i^{\text{ème}}$  composante de la sortie calculée avec les poids synaptiques  $w$  et  $y_i^e$  correspond à celle désirée.

On cherche alors à minimiser l'écart entre les sorties estimées par le réseau et celles désirées au sein de la base d'apprentissage. Dans ce but, on utilise une descente de gradient. Afin d'éviter de rencontrer un minimum local, chaque situation de la base d'apprentissage est présentée de façon aléatoire au réseau. À chaque étape (*i.e.*, à chaque fois qu'une nouvelle situation  $e$ , tirée aléatoirement dans  $\{1 \dots p\}$ , est présentée au réseau pendant la phase d'apprentissage), on minimise le le critère :

$$C^e = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i^e(w) - y_i^e)^2$$

À partir de là, on effectue une rétropropagation de l'erreur pour adapter les différents poids synaptiques du réseau. Soit  $w_{i,j}$  un poids synaptique du réseau que l'on veut modifier. On a, pour le neurone  $i$ , l'entrée totale  $a_i = \sum_{l=1}^o w_{i,l} \cdot x_l$  et la sortie  $x_i = \sigma(a_i)$ . On a alors :

$$\frac{\partial C^e}{\partial w_{i,j}} = \underbrace{\frac{\partial C^e}{\partial a_i}}_{=\delta_i} \cdot \underbrace{\frac{\partial a_i}{\partial w_{i,j}}}_{=x_j}$$

On peut déterminer  $\delta_i$  pour les couches de sortie :

$$\begin{aligned} \delta_i &= \frac{\partial C^e}{\partial a_i} \\ &= \sigma'(a_i) \frac{\partial C^e}{\partial y_i} \\ &= \sigma'(a_i) (\hat{y}_i^e - y_i^e) \end{aligned}$$

Pour les neurones situés sur la couche cachée, on a une relation de récurrence avec la couche suivante :

$$\begin{aligned} \delta_i &= \frac{\partial C^e}{\partial a_i} \\ &= \sum_j \left( \frac{\partial C^e}{\partial a_k} \cdot \frac{\partial a_j}{\partial a_i} \right) \\ &= \sum_j (\delta_k \cdot w_{i,j} \cdot \sigma'(a_i)) \\ &= \sigma'(a_i) \cdot \sum_j (\delta_j \cdot w_{i,j}) \end{aligned}$$

On peut alors déterminer tous les  $\delta_i$  à partir de ceux de la dernière couche.

La modification des poids synaptiques du réseau se fait ensuite par une descente de gradient, sous la forme :

$$\Delta w = -\gamma w \cdot \frac{\partial C}{\partial w}$$

Cet algorithme de rétropropagation est adapté à l'architecture du perceptron multicouche que l'on utilise ici. La structure en couches successives du réseau permet de ne pas avoir à recalculer les dérivés partielles des couches cachées du réseau. Cet algorithme permet de garder une relation linéaire entre le temps de calcul pour la mise à jour des poids synaptiques et le nombre de paramètres du réseau. L'apprentissage du réseau peut donc être décrit en sept étapes successives (Aires 1999) :

1. On initialise aléatoirement les poids  $w_{i,j}$  du réseau ;
2. On choisit une entrée  $x^e$  parmi la base d'apprentissage ;
3. On propage cette entrée parmi les  $m$  couches du réseau :  $x_i^o = \sigma(a_i^m) = \sigma\left(\sum_j w_{i,j} x_j^{m-1}\right)$  ;
4. On calcule les erreurs  $\delta_i$  pour la couche de sortie :  $\delta_i = \sigma'(a_i)(\hat{y}_i^e - y_i^e)$  ;
5. On rétropropage ces erreurs à travers les  $m$  couches :  $\delta_i^{m-1} = \sigma'(a_i) \cdot \sum_j \left(\delta_j^m \cdot w_{i,j}^m\right)$  ;
6. On mets à jour les poids synaptiques du réseau :  $\Delta w_{i,j}^m = \rho_{i,j} \delta_i^m x_j^{m-1}$ , où  $\rho_{i,j}$  est le pas d'apprentissage pour le poids synaptique  $w_{i,j}$  (méthode classique de descente de gradient stochastique) ;
7. On recommence à l'étape 2, jusqu'à atteindre un critère de convergence.

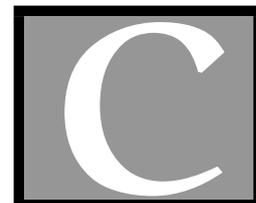
De nombreux critères de convergence peuvent être choisis : le nombre d'itérations, la valeur du critère à minimiser, la modification des paramètres ou un test sur la capacité de généralisation du réseau. Nous utilisons ici la dernière possibilité grâce au calcul de l'erreur du réseau de neurones considéré sur la base de validation.

La rétropropagation présentée ici utilise une descente de gradient stochastique simple. Il existe de nombreuses techniques d'optimisation plus sophistiquées. La méthode que nous utilisons au cours des apprentissages des différents réseaux est celle de Levenberg-Marquardt. La différence consiste en l'introduction d'un terme de régularisation (aussi appelé amortissement) qui va être modifié à chaque itération et permettre de mieux conditionner le choix des poids synaptiques (Hagan and Menhaj 1994).

Il s'agit cette fois d'une inversion globale. Une fois paramétré, le réseau n'évolue pas. On obtient l'inversion désirée directement en sortie.



## RE-ANALYSES DE L'ECMWF



Le tableau ci-dessous décrit toutes les variables des re-analyses de l'ECMWF ([Berrisford et al. 2011](#); [Dee et al. 2011](#); [Uppala et al. 2005](#); [Simmons and Gibson 2000](#)). Il s'agit des données ERA-40 que l'on a remises en forme pour notre utilisation.

Numéro	Description	Unité
Localisation		
1	Numéro de cellule	
2	Latitude	°
3	Longitude	°
Date		
4	Année	
5	Mois	
6	Jour	
7	Heure	
8	Minute	
Attributs		
9	Type de surface	
10	Altitude	m
Profils atmosphériques		
11-53	Température	K
54-96	Humidité spécifique	ppmv
97-139	Humidité relative	%
140-182	Ozone	ppmv
Profils de nuages		
183-242	Couverture nuageuse	
243-302	Eau liquide dans le nuage	$\text{g}\cdot\text{kg}^{-1}$
303-362	Glace dans le nuage	$\text{g}\cdot\text{kg}^{-1}$
Données de pluie		
363	Précipitations à grande échelle	m
364	Précipitations convectives	m

ANNEXE C. RE-ANALYSES DE L'ECMWF

Numéro	Description	Unité
Données à 2 m		
365	Température à 2 m	K
366	Pression à 2 m	hPa
367	Vent Ouest-Est	$\text{m}\cdot\text{s}^{-1}$
368	Vent Nord-Sud	$\text{m}\cdot\text{s}^{-1}$
369	Fetch du vent (étendue où il peut souffler sans rencontrer d'obstacle)	m
Donnée de surface		
370	Température de surface	K
Colonne totale		
371	Couverture nuageuse totale	%
372	Quantité de vapeur d'eau intégrée	$\text{kg}\cdot\text{m}^{-2}$
373	Quantité d'eau liquide	$\text{kg}\cdot\text{m}^{-2}$
374	Quantité de glace	$\text{kg}\cdot\text{m}^{-2}$
Données de nuages		
375	Haute couverture nuageuse	
376	Couverture nuageuse intermédiaire	
377	Basse couverture nuageuse	
Émissivité		
378-384	Émissivité de surface pour les canaux SSM/I	

# TABLE DES MATIÈRES

<b>Sommaire</b>	<b>vi</b>
<b>Introduction</b>	<b>1</b>
<b>1 Le sondage atmosphérique</b>	<b>5</b>
1.1 Le transfert radiatif . . . . .	7
1.1.1 Le rayonnement électromagnétique . . . . .	7
1.1.1.1 L'émission . . . . .	8
1.1.1.2 L'absorption . . . . .	9
1.1.1.3 La diffusion . . . . .	11
1.1.2 L'équation de transfert radiatif . . . . .	12
1.2 Le sondage satellite : principe . . . . .	18
1.2.1 Le spectre électromagnétique . . . . .	21
1.2.2 Le sondage dans le micro-onde . . . . .	24
1.2.3 Le sondage dans l'infrarouge . . . . .	25
1.3 La plate-forme MetOp . . . . .	25
1.3.1 IASI . . . . .	29
1.3.2 AMSU-A . . . . .	31
1.3.3 MHS . . . . .	32
1.4 Le modèle de transfert radiatif RTTOV . . . . .	33
1.5 Les centres de prévision météorologique . . . . .	34
1.5.1 Dénomination des données satellites . . . . .	34
1.5.2 Les centres de prévision numérique : NWP . . . . .	34
<b>2 Restitutions de surface</b>	<b>37</b>
2.1 L'émissivité . . . . .	40
2.1.1 Émissivité de surface dans le micro-onde . . . . .	41
2.1.2 Émissivité de surface dans l'infrarouge . . . . .	42
2.1.2.1 Bases de données d'émissivités mesurées en laboratoire . . . . .	43
2.1.2.2 L'émissivité restituée à partir de MODIS . . . . .	44
2.1.2.3 La base de données UWIRemis . . . . .	46
2.1.2.4 L'émissivité IASI de la NASA . . . . .	48

2.1.2.5	L'émissivité IASI du groupe ARA . . . . .	48
2.1.3	Mise en coïncidence des bases . . . . .	49
2.2	Construction d'une base de données d'émissivités infrarouges . . . . .	49
2.2.1	Les réseaux de neurones . . . . .	50
2.2.2	Base hyperspectrale d'émissivité . . . . .	51
2.2.2.1	Une base de première ébauche indépendante de IASI . . . . .	51
2.2.2.2	Une classification de surface . . . . .	51
2.2.2.3	Création de la base à partir des observations MODIS . . . . .	52
2.2.3	Analyse en composantes principales des spectres d'émissivité . . . . .	53
2.2.3.1	Exemple simple d'analyse en composantes principales . . . . .	53
2.2.3.2	Méthodologie . . . . .	53
2.2.3.3	Résultats obtenus sur les spectres d'émissivité . . . . .	54
2.2.3.4	Compression des émissivités infrarouges . . . . .	55
2.2.4	Interpolation spectrale de l'émissivité infrarouge . . . . .	56
2.2.5	La première ébauche . . . . .	59
2.3	Un nouvel algorithme de restitution de surface . . . . .	60
2.3.1	Inversion bayésienne du transfert radiatif . . . . .	61
2.3.2	Sélection des paramètres de la restitution . . . . .	65
2.3.2.1	La matrice de covariance d'erreur de $F$ . . . . .	65
2.3.2.2	La matrice de covariance d'erreur de la première ébauche . . . . .	66
2.3.2.3	Les canaux sélectionnés . . . . .	67
2.3.2.4	Le nombre de composantes de l'ACP . . . . .	68
2.3.2.5	Conclusion . . . . .	70
2.3.3	Influence de la première ébauche . . . . .	70
2.3.4	Résultats et évaluation . . . . .	71
2.3.4.1	Conditions expérimentales . . . . .	71
2.3.4.2	Analyse spectrale . . . . .	72
2.3.4.3	Évaluation de la température de surface . . . . .	76
2.4	Conclusion . . . . .	82
<b>3</b>	<b>Chaîne opérationnelle de restitution de surface</b> . . . . .	<b>83</b>
3.1	Mise en place d'une chaîne opérationnelle . . . . .	84
3.2	Restitutions en pseudo temps réel . . . . .	84
3.2.1	Étude spatiale . . . . .	85
3.2.2	Étude spectrale . . . . .	88
3.3	Moyennes mensuelles . . . . .	90
3.3.1	Variabilité spatiale . . . . .	91
3.3.2	Variabilité temporelle . . . . .	95
3.4	Comparaison avec des radiosondages au-dessus du dôme C . . . . .	98
3.4.1	Données disponibles et objectif . . . . .	98

3.4.2	Méthode . . . . .	98
3.4.3	Résultats . . . . .	99
3.5	Perspectives . . . . .	104
3.6	Conclusion . . . . .	105
<b>4</b>	<b>Synergie : Principes généraux</b>	<b>107</b>
4.1	Principe . . . . .	108
4.2	Exemple sans synergie . . . . .	111
4.3	Synergie additive . . . . .	111
4.4	Synergie indirecte . . . . .	114
4.5	Synergie de débruitage . . . . .	117
4.6	Généralisation . . . . .	118
4.7	Conclusion . . . . .	119
<b>5</b>	<b>Synergie au-dessus des océans, en ciel clair</b>	<b>121</b>
5.1	Base d'apprentissage . . . . .	123
5.1.1	Données atmosphériques utilisées . . . . .	123
5.1.2	Échantillonnage de la base . . . . .	124
5.1.3	Simulation de transfert radiatif et bruit instrumental . . . . .	127
5.1.4	Réduction de la dimension des données IASI . . . . .	129
5.1.4.1	Sélection de canaux . . . . .	129
5.1.4.2	Compression de canaux . . . . .	130
5.1.4.3	Statistiques de compression . . . . .	131
5.1.4.4	Statistiques de débruitage . . . . .	133
5.2	Restitutions atmosphériques . . . . .	134
5.2.1	Corrélations entre les variables et les mesures . . . . .	136
5.2.2	Restitutions par $k$ -plus proches voisins . . . . .	139
5.2.2.1	Méthode . . . . .	139
5.2.2.2	Résultats . . . . .	141
5.2.3	Restitutions par régression linéaire . . . . .	145
5.2.3.1	Méthode . . . . .	145
5.2.3.2	Résultats . . . . .	145
5.2.4	Restitutions grâce à un réseau de neurones . . . . .	147
5.2.4.1	Méthode . . . . .	147
5.2.4.2	Résultats . . . . .	148
5.2.5	Comparaison des trois méthodes . . . . .	150
5.2.6	Évaluation de la synergie . . . . .	152
5.3	Conclusion . . . . .	154

<b>6</b>	<b>Échantillonnage par entropie</b>	<b>155</b>
6.1	L'entropie au sens de Shannon . . . . .	157
6.1.1	Historique et définition . . . . .	158
6.1.2	Exemple de calcul d'une entropie . . . . .	159
6.2	Base de données atmosphériques initiale . . . . .	160
6.2.1	Les produits L2 IASI d'EUMETSAT . . . . .	160
6.2.2	Niveaux de pression . . . . .	162
6.3	Échantillonnage . . . . .	164
6.3.1	Méthode d'échantillonnage : l'entropie . . . . .	165
6.3.1.1	Discrétisation . . . . .	166
6.3.1.2	Variabilité de la vapeur d'eau . . . . .	166
6.3.1.3	Pondération de chaque variable . . . . .	167
6.3.1.4	Optimisation du calcul . . . . .	169
6.3.1.5	Convergence de l'algorithme d'échantillonnage . . . . .	170
6.3.2	Apport de l'entropie pour l'échantillonnage multivarié . . . . .	171
6.3.2.1	Exemple d'échantillonnage sur une seule variable . . . . .	171
6.3.2.2	Prise en compte de plusieurs variables grâce à l'entropie . . . . .	174
6.4	Résultats . . . . .	176
6.4.1	Représentation des différentes variables . . . . .	176
6.4.1.1	Base de données complète . . . . .	176
6.4.1.2	Base de données échantillonnée . . . . .	177
6.4.2	Variabilité spatiale . . . . .	179
6.5	Conclusion . . . . .	180
<b>7</b>	<b>Synergie au-dessus des continents, en ciel clair</b>	<b>183</b>
7.1	Configuration des restitutions . . . . .	185
7.1.1	Base d'apprentissage . . . . .	186
7.1.2	Méthode de restitution . . . . .	187
7.1.3	Convergence des différents réseaux de neurones . . . . .	188
7.2	Apports de l'émissivité et de la température de surface . . . . .	190
7.2.1	Dans le micro-onde . . . . .	190
7.2.2	Dans l'infrarouge . . . . .	192
7.3	Synergie infrarouge et micro-onde . . . . .	194
7.3.1	Sans informations de surface . . . . .	195
7.3.2	Synergie complète entre toutes les données disponibles . . . . .	196
7.4	Synergie et informations de surface . . . . .	198
7.5	Conclusion . . . . .	200
	<b>Conclusion et perspectives</b>	<b>201</b>

---

<b>A</b>	<b>Acronymes</b>	<b>207</b>
<b>B</b>	<b>Mathématiques</b>	<b>211</b>
B.1	L'Analyse en Composantes Principales . . . . .	211
B.2	Les algorithmes d'inversion . . . . .	213
B.2.1	La restitution bayésienne . . . . .	213
B.2.2	Le réseau de neurones . . . . .	215
<b>C</b>	<b>Re-analyses de l'ECMWF</b>	<b>219</b>
	<b>Table des matières</b>	<b>225</b>
	<b>Table des figures</b>	<b>229</b>
	<b>Bibliographie</b>	<b>243</b>



# TABLE DES FIGURES

1.1	Mesures météorologiques par Florin Périer et Gay-Lussac . . . . .	6
1.2	Spectre d'émission du Soleil et de la Terre . . . . .	9
1.3	Schéma des niveaux d'excitation d'une molécule . . . . .	10
1.4	Schéma des différents types de diffusion . . . . .	12
1.5	Variation de la luminance spectrale traversant une souche simple . . . . .	13
1.6	Luminance spectrale reçue par un satellite . . . . .	14
1.7	Rayonnement reçu par un satellite au sommet de l'atmosphère . . . . .	16
1.8	Exemple de fonction de poids . . . . .	17
1.9	Le spectre électromagnétique . . . . .	21
1.10	Représentation de l'orbite polaire . . . . .	26
1.11	Schéma des géométries de balayage pour les capteurs satellite . . . . .	28
1.12	Représentation de l'orbite de MetOp avec une fauchée . . . . .	29
1.13	Jacobiens de IASI . . . . .	31
1.14	Jacobiens de AMSU-A et MHS . . . . .	33
2.1	Sensibilité de IASI aux paramètres de surface . . . . .	40
2.2	Carte de variabilité annuelle des émissivités MODIS . . . . .	46
2.3	Exemple de spectres d'émissivité . . . . .	47
2.4	Schéma d'un réseau de neurones . . . . .	50
2.5	L'analyse en composante principale . . . . .	53
2.6	Les composantes principales . . . . .	55
2.7	Erreur spectrale de compression . . . . .	56
2.8	Erreur spectrale d'interpolation . . . . .	58
2.9	Erreur spatiale d'interpolation . . . . .	58
2.10	Exemple d'interpolation . . . . .	59
2.11	Sélection des canaux pour la restitution . . . . .	68
2.12	Sélection du nombre de composantes pour la restitution . . . . .	69
2.13	Statistiques d'écart en température de brillance . . . . .	73
2.14	Statistiques d'écart en température de brillance pour les données restituées . . . . .	75
2.15	Carte d'écart entre les différentes températures de surface . . . . .	78
3.1	Carte d'émissivités au-dessus de l'Afrique du nord . . . . .	85

3.2	Carte imagée d'émissivités restituées au-dessus de l'Afrique du Nord . . . . .	87
3.3	Histogrammes des émissivités à 833 et 1160 $\text{cm}^{-1}$ . . . . .	89
3.4	Émissivité moyenne en janvier à 1120 $\text{cm}^{-1}$ . . . . .	92
3.5	Erreur théorique et dispersion moyenne en janvier à 1120 $\text{cm}^{-1}$ . . . . .	94
3.6	Variations temporelles de l'émissivité et de la pluviométrie . . . . .	97
3.7	Colocalisation des sondages IASI au-dessus du Dôme C . . . . .	99
3.8	Comparaison des températures de surface au-dessus de l'Antarctique . . . . .	100
3.9	Comparaison des profils atmosphériques au dessus de l'Antarctique . . . . .	102
3.10	Émissivités restituées au-dessus du Dôme C . . . . .	103
4.1	Schéma de restitution synergique . . . . .	110
4.2	Variation de la synergie additive en fonction du jacobien de restitution . . . . .	112
4.3	Variation de la synergie additive en fonction du bruit instrumental . . . . .	114
4.4	Variation de la synergie indirecte en fonction de la première ébauche . . . . .	116
4.5	Variation de la synergie de débruitage en fonction des observations . . . . .	118
5.1	Variabilité de certains représentants de la base échantillonnée par $k$ -moyennes	126
5.2	Répartition géographique de la base échantillonnée par $k$ -moyennes . . . . .	126
5.3	Bruit instrumental de l'instrument IASI . . . . .	128
5.4	Statistiques d'erreur de débruitage de l'ACP . . . . .	132
5.5	Évolution de l'erreur de débruitage de l'ACP . . . . .	133
5.6	Corrélation entre les mesures micro-ondes et les variables atmosphériques . . . . .	136
5.7	Corrélation entre les mesures infrarouges et les variables atmosphériques . . . . .	137
5.8	Corrélation entre les mesures infrarouges et micro-ondes . . . . .	138
5.9	Corrélation entre la température et la vapeur d'eau atmosphérique . . . . .	139
5.10	Variation de l'erreur de restitution en fonction du nombre de voisins . . . . .	142
5.11	Schéma des configurations de restitution . . . . .	142
5.12	Restitution de la température par $k$ -plus proches voisins . . . . .	143
5.13	Restitution de l'humidité relative par $k$ -plus proches voisins . . . . .	144
5.14	Restitution de la température par régression linéaire . . . . .	146
5.15	Restitution de l'humidité relative par régression linéaire . . . . .	147
5.16	Restitution de la température par réseau de neurones . . . . .	149
5.17	Restitution de l'humidité relative par réseau de neurones . . . . .	150
5.18	Comparaison des restitutions de la température . . . . .	151
5.19	Comparaison des restitutions de la l'humidité relative . . . . .	152
5.20	Facteur de synergie pour la restitution de la température . . . . .	153
5.21	Facteur de synergie pour la restitution de l'humidité relative . . . . .	154
6.1	Répartition spatiale des différentes bases de données . . . . .	164
6.2	Densité de probabilité de la vapeur d'eau par niveau de pression . . . . .	167

---

6.3	Convergence de l'échantillonnage par entropie . . . . .	171
6.4	Variabilité des profils des bases de l'ECMWF . . . . .	173
6.5	Variabilité des profils des bases échantillonnées par entropie . . . . .	175
6.6	Densité de probabilités des variables de la base de données complète . . . . .	177
6.7	Densité de probabilités des variables de la base de données échantillonnée . . . . .	178
6.8	Répartition spatiale des bases de données échantillonnées . . . . .	180
7.1	Schéma des configurations de restitutions au-dessus des continents . . . . .	185
7.2	Convergence des réseaux de neurones . . . . .	189
7.3	Apport de la surface dans le micro-onde pour la restitution des profils . . . . .	191
7.4	Apport de la surface dans l'infrarouge pour la restitution des profils . . . . .	193
7.5	Restitution synergique des profils . . . . .	195
7.6	Restitution complète des profils . . . . .	197
7.7	Apport synergique et de surface . . . . .	199



# BIBLIOGRAPHIE

- F. Aires. *Problèmes inverses et réseaux de neurones : application à l'interféromètre haute résolution IASI et à l'analyse de séries temporelles*. PhD thesis, Université Paris IX - Dauphine, mars 1999. [30](#), [64](#), [213](#), [217](#)
- F. Aires. Neural network uncertainty assessment using Bayesian statistics with application to remote sensing : 1. Network weights. *Journal of Geophysical Research (Atmospheres)*, 109(D18) :10303–+, May 2004. doi : 10.1029/2003JD004173. [135](#)
- F. Aires. Measure and exploitation of multi-sensor and multi-wavelength synergy for remote sensing : Part I - Theoretical considerations. *Journal of Geophysical Research*, 116 (D02301), 2011. doi : 10.1029/2010JD014701. [108](#)
- F. Aires and C. Prigent. Sampling techniques in high-dimensional spaces for satellite remote sensing databases generation. *Journal of Geophysical Research*, 112 :D20301, 2007. doi : 10.1029/2007JD008391. [124](#), [136](#), [156](#), [157](#)
- F. Aires, C. Prigent, W. Rossow, and M. Rothstein. A new neural network approach including first-guess for retrieval of atmospheric water vapour, cloud liquid water path, surface temperature and emissivities over land from satellite microwave observations. *Journal of Geophysical Research*, 106(D14) :14887–14907, 2001. [19](#), [38](#), [39](#), [156](#)
- F. Aires, A. Chedin, N. Scott, and W. Rossow. A regularized neural net approach for retrieval of atmospheric and surface temperatures with the IASI instrument. *Journal of Applied Meteorology*, 41(2) :144–159, 2002a. [129](#), [157](#)
- F. Aires, W. Rossow, and A. Chedin. Rotation of EOFs by the independent component analysis : Toward a solution of the mixing problem in the decomposition of geophysical time series. *Journal of the Atmospheric Sciences*, 59(1) :111–123, 2002b. [55](#), [130](#)
- F. Aires, W. Rossow, N. Scott, and A. Chedin. Remote sensing from the infrared atmospheric sounding interferometer instrument 1. Compression, denoising, and first-guess retrieval algorithms. *Journal of Geophysical Research (Atmospheres)*, 107 :6–1, November 2002c. doi : 10.1029/2001JD000955. [127](#), [130](#)
- F. Aires, W. Rossow, N. Scott, and A. Chedin. Remote sensing from the infrared atmospheric sounding interferometer instrument 2. Simultaneous retrieval of temperature, water vapor,

- and ozone atmospheric profiles. *Journal of Geophysical Research (Atmospheres)*, 107 :7–1, November 2002d. doi : 10.1029/2001JD001591. [130](#)
- F. Aires, C. Prigent, and W. Rossow. Neural network uncertainty assessment using Bayesian statistics with application to remote sensing : 2. Output errors. *Journal of Geophysical Research (Atmospheres)*, 109(D10), 2004a. [135](#)
- F. Aires, C. Prigent, and W. Rossow. Neural network uncertainty assessment using Bayesian statistics with application to remote sensing : 3. Network Jacobians. *Journal of Geophysical Research (Atmospheres)*, 109(D10), 2004b. [135](#)
- F. Aires, F. Bernardo, H. Brogniez, and C. Prigent. An innovative calibration method for the inversion of satellite observations. *Journal of Applied Meteorology and Climatology*, 49(12) :2458–2473, 2010. [157](#)
- F. Aires, M. Paul, C. Prigent, B. Rommen, and M. Bouvet. Measure and exploitation of multi-sensor and multi-wavelength synergy for remote sensing : Part II - An application for the retrieval of atmospheric temperature and water vapour from MetOp. *Journal of Geophysical Research*, 116(D02302), 2011a. doi : 10.1026/2010JD014702. [19](#), [122](#), [131](#), [134](#), [204](#)
- F. Aires, C. Prigent, F. Bernardo, C. Jiménez, R. Saunders, and P. Brunel. A Tool to Estimate Land-Surface Emissivities at Microwave frequencies (TELSEM) for use in numerical weather prediction. *Quarterly Journal of the Royal Meteorological Society*, 137 :690–699, 2011b. doi : 10.1002/qj.803. [34](#), [38](#), [40](#), [42](#), [186](#)
- F. Aires, O. Aznay, C. Prigent, M. Paul, and F. Bernardo. Synergetic multi-wavelength remote sensing versus a posteriori combination of retrieved products : Application for the retrieval of atmospheric profiles using MetOp. *Journal of Geophysical Research*, 117 (D18304), 2012. doi : 10.1029/2011JD017188. [20](#), [122](#)
- R. Amer, T. Kusky, and A. Ghulam. Lithological mapping in the Central Eastern Desert of Egypt using ASTER data. *Journal of African Earth Sciences*, 56(2) :75–82, 2010. [41](#)
- T. August, D. Klaes, P. Schlüssel, T. Hultberg, M. Crapeau, A. Arriaga, A. O’Caraoll, D. Coppens, R. Munro, and X. Calbet. IASI on Metop-A : Operational Level 2 retrievals after five years in orbit. *Journal of Quantitative Spectroscopy and Radiative Transfer*, 113 (11) :1340–1371, 2012. doi : 10.1016/j.jqsrt.2012.02.028. [104](#)
- A.M. Baldridge, S.J. Hook, C.I. Grove, and G. Rivera. The ASTER spectral library version 2.0. *Remote Sensing of Environment*, 113(4) :711–715, 2009. [44](#)
- M. Barral. On the rain-water collected at the observatory at Paris. *Philosophical Magazine Series 4*, 4(26) :396–398, 1852. doi : 10.1080/14786445208647147. [6](#)

- 
- P. Berrisford, D. Dee, P. Poli, R. Brugge, K. Fielding, M. Fuentes, P. Kallberg, S. Kobayashi, S. Uppala, and A. Simmons. The era-interim archive version 2.0. Technical report, ECMWF, November 2011. [123](#), [219](#)
- W. Blackwell. A neural-network technique for the retrieval of atmospheric temperature and moisture profiles from high spectral resolution sounding data. *IEEE Transactions on Geoscience and Remote Sensing*, 43(11) :2535–2546, 2005. [131](#)
- E. Borbas and B. Ruston. *The RTTOV UWiremis module IR land surface emissivity module*, volume AS09-04. EUMETSAT, 2010. [47](#)
- E. Borbas, R. O. Knuteson, S. W. Seemann, E. Weize, L. Moy, and H. L. Huang. A high spectral resolution global land surface infrared emissivity database. *Joint 2007 EUMETSAT Meteorological Satellite Conference and the 15th Satellite Meteorology and Oceanography Conference of the American Meteorological Society*, 2007. [47](#)
- V. Capelle, A. Chedin, E. Pequignot, P. Schluessel, and S. M. Newman. Infrared continental surface emissivity and skin temperature retrieved from IASI observations over the tropics. *Journal of Applied Meteorology and Climatology*, 51(6) :1164–1179, 2012. [48](#), [104](#)
- V. Caselles, E. Valor, C. Coll, and E. Rubio. Thermal band selection for the prism instrument 1. analysis of emissivity-temperature separation algorithms. *Journal of Geophysical Research*, 102(D10) :11145–11164, 1997. [76](#)
- G. Chalon, F. Cayla, and D. Diebel. IASI - An advanced sounder for operational meteorology. *Proceedings of the 52<sup>nd</sup> Congress of IAF, Toulouse France, 1-5 Oct. 2001*, 2001. [29](#)
- A. Chedin, N. Scott, C. Wahiche, and P. Moulinier. The improved initialization inversion method : A high resolution physical method for temperature retrievals from satellites of the TIROS-N series. *Journal of Applied Meteorology and Climatology*, 24(2) :128–143, 1985. [157](#)
- A. Chedin, E. Pequignot, S. Serrar, and N. Scott. Simultaneous determination of continental surface emissivity and temperature from NOAA 10/HIRS observations : Analysis of their seasonal variations. *Journal of Geophysical Research*, 109(D20110) : 10.1029/2004JD004886, 2004. [43](#)
- F. Chevallier, F. Cheruy, N. Scott, and A. Chedin. A neural network approach for a fast and accurate computation of a longwave radiative budget. *Journal of Applied Meteorology*, 37(11) :1385–1397, 1998. [157](#)
- F. Chevallier, J-J. Morcrette, A. Chedin, and F. Cheruy. TIGR-like atmospheric-profile databases for accurate radiative-flux computation. *Quarterly Journal of the Royal Meteorological Society*, 126(563) :777–785, 2000. [157](#)

- F. Chevallier, S. Di Michele, and A. McNally. *Diverse profile datasets from the ECMWF 91-level short-range forecasts*. European Centre for Medium-Range Weather Forecasts, 2006. [157](#), [171](#)
- C. Cho and D. H. Staelin. Could clearing of Atmospheric Infrared Sounder hyperspectral infrared radiances using stochastic methods. *Journal of Geophysical Research*, 111(D9) : D09S18, 2006. [122](#)
- A. Cholesky. Sur la résolution numérique des systèmes d'équations linéaires. Archives de l'École Polytechnique, December 1910. [54](#)
- E. Clapeyron. Mémoire sur la puissance motrice de la chaleur. *Journal de l'École Polytechnique*, XXIII<sup>ème</sup> cahier(Tome XIV), 1834. [20](#)
- R. Clausius. "Ueber die bewegende Kraft der Wärme und die Gesetze, welche sich daraus für die Wärmelehre selbst ableiten lassen. *Annalen der Physik*, 173 :441–476, 1856. doi : 10.1002/andp.18561730306. [20](#)
- P. Courtier, E. Andersson, W. Heckley, D. Vasiljevic, M. Hamrud, A. Hollingsworth, F. Rabier, M. Fisher, and J. Pailleux. The ECMWF implementation of three-dimensional variational assimilation (3D-Var). I : Formulation. *Quarterly Journal of the Royal Meteorological Society*, 124(550) :1783–1807, 1998. [35](#)
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4) :303–314, 1989. [148](#)
- D. Dee, S. Uppala, A. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae, M. Balsameda, G. Balsamo, P. Bauer, A. Beljaars, L. van de Berg, J. Bidlot, N. Bormann, C. Delsol, R. Dragani, M. Fuentes, A. Geer, L. Haimberger, S. Healy, H. Hersbach, E. Holm, L. Isaksen, P. Kallberg, M. Kohler, M. Matricardi, A. McNally, B. Monge-Sanz, J.-J. Morcrette, B.-K. Park, C. Peubey, P. de Rosnay, C. Tavolato, J.-N. Thépaut, and F. Vitart. The ERA-Interim reanalysis : configuration and performance of the data assimilation system. *Quarterly Journal of the Royal Meteorological Society*, 137 :553–597, April 2011. doi : 10.1002/qj.828. [123](#), [219](#)
- M. Divakarla, C. Barnet, M. Goldberg, L. McMillin, E. Maddy, L. Zhou, and X. Liu. Validation of Atmospheric Infrared Sounder temperature and water vapor retrievals with matched radiosonde measurements and forecasts. *Journal of Geophysical Research (Atmospheres)*, 111(D9), 2006. doi : {10.1029/2005JD006116}. [151](#)
- S. English. Estimation of temperature and humidity profile information from microwave radiances over different surface types. *Journal of Applied Meteorology*, 38(10) :1526–1541, 1999. [184](#)

- S. English and T. Hewison. A fast generic millimeter-wave emissivity model. *Proceedings SPIE 3503*, Microwave Remote Sensing of the Atmosphere and Environment :288, August 1998. doi : doi:10.1117/12.319490. [41](#)
- P. Eriksson, C. Jiménez, S. Bühler, and D. Murtagh. A Hotelling transformation approach for rapid inversion of atmospheric spectra. *Journal of Quantitative Spectroscopy and Radiative Transfer*, 73(6) :529–543, 2002. [130](#)
- J. Escobar. *Base de données pour la restitution de paramètres atmosphériques à l'échelle globales ; étude sur l'inversion par réseaux de neurones des données des sondeurs verticaux atmosphériques satellitaires présents et à venir*. PhD thesis, Université de Paris Diderot, Paris VII, 1993. [157](#)
- EUMETSAT Document. IASI Level 2 Product Guide. Technical report, EUMETSAT, 2012. [161](#)
- K. Evans, J. Wang, P. Racette, G. Heymsfield, and L. Li. Ice cloud retrievals and analysis with the compact scanning submillimeter imaging radiometer and the cloud radar system during CRYSTAL FACE. *Journal of Applied Meteorology*, 44(6) :839–859, 2005. [140](#)
- J.R. Eyre. A fast radiative transfer model for satellite sounding systems. Technical Report 176, ECMWF Research Department Technical Memo, 1991. [33](#)
- S. Freitas, I. Trigo, J. Bioucas-Dias, and F. Götsche. Quantifying the uncertainty of land surface temperature retrievals from seviri/meteosat. *IEEE Transactions on Geoscience and Remote Sensing*, 48(1) :523–534, January 2010. doi : 10.1109/TGRS.2009.2027697. [78](#)
- L. Garand, M. Buehner, and N. Wagner. Background error correlation between surface skin and air temperatures : estimation and impact on the assimilation of infrared window radiances. *Journal of Applied Meteorology*, 43(12) :1853–1863, 2004. [184](#)
- S. Geman, E. Bienenstock, and R. Doursat. Neural networks and the bias/variance dilemma. *Neural computation*, 4(1) :1–58, 1992. [140](#)
- G. Grandjean, P. Paillou, N. Baghdadi, E. Heggy, T. August, and Y. Lasne. Surface and subsurface structural mapping using low frequency radar : A synthesis of the Mauritanian and Egyptian experiments. *Journal of African Earth Sciences*, 44(2) :220–228, 2006. [41](#)
- P. Griffiths. Fourier transform infrared spectrometry. *Science*, 222(4621) :297–302, 1983. doi : 10.1126/science.6623077. [42](#)
- M. T. Hagan and M. Menhaj. Training feedforward networks with the Marquardt algorithm. *IEEE Transactions on Neural Networks*, 5(6) :989–993, 1994. [217](#)

- T. Hewison and R. Saunders. Measurements of the AMSU-B antenna pattern. *IEEE Transactions on Geoscience and Remote Sensing*, 34(2) :405–412, March 1996. doi : 10.1109/36.485118. [32](#)
- K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5) :359–366, 1989. [50](#), [148](#)
- H.-L. Huang and P. Antonelli. Application of principal component analysis to high-resolution infrared measurement compression and retrieval. *Journal of Applied Meteorology*, 40(3) :365–388, 2001. [130](#)
- H. L. Huang and R. Purser. Objective measures of the information density of satellite data. *Meteorology and Atmospheric Physics*, 60(1-3) :105–117, 1996. [129](#)
- G.C Hulley and S.J. Hook. The North American ASTER Land Surface Emissivity Database (NAALSED) Version 2.0. *Remote Sensing of Environment*, 113 :1967–1975, 2009. [39](#), [71](#)
- P. A. E. M. Janssen, S. Abdalla, H. Hersbach, and J. R. Bidlot. Error estimation of buoy, satellite and model wave height data. *Journal of Atmospheric and Oceanic Technology*, 24 :1665–1677, September 2007. doi : 10.1175/JTECH2069.1. [79](#)
- C. Jiménez, S. Bühler, B. Rydberg, P. Eriksson, and K. Evans. Performance simulations for a submillimetre-wave satellite instrument to measure cloud ice. *Quarterly Journal of the Royal Meteorological Society*, 133(S2) :129–149, 2007. [140](#)
- C. Jiménez, J. Catherinot, C. Prigent, and J. Roger. Relations between geological characteristics and satellite-derived infrared and microwave emissivities over deserts in northern Africa and the Arabian Peninsula. *Journal of Geophysical Research*, 115, October 2010. doi : 10.1029/2010JD013959. [42](#), [86](#), [183](#)
- I. T. Jolliffe. *Principal Component Analysis : second edition*. Springer, 2002. [53](#), [131](#)
- D. Klaes, M. Cohen, Y. Buhler, P. Schlüssel, R. Munro, A. Engeln, E. Clérigh, H. Bonekamp, J. Ackermann, J. Schmetz, and J-P. Luntama. An Introduction to the EUMETSAT Polar system. *Bulletin of the American Meteorological Society*, 88(7) :1085–1096, September 2007. doi : 10.1175/BAMS-88-7-1085. [26](#)
- J. Kornfeld and J. Susskind. On the effect of surface emissivity on temperature retrievals. *Monthly Weather Review*, 105(12) :1605–1608, 1977. [184](#)
- C. Kummerow, W. Barnes, T. Kozu ans J. Shiue, and J. Simpson. The tropical rainfall measuring mission (TRMM) sensor package. *Journal of Atmospheric and Oceanic Technology*, 15(3) :809–817, 1998. [95](#)

- C. Kummerow, Y. Hong, W. Olson, S. Yang, R. Adler, J. McCollum, R. Ferraro, G. Petty, D. Shin, and T. Wilheit. The evolution of the Goddard Profiling Algorithm (GPROF) for rainfall estimation from passive microwave sensors. *Journal of Applied Meteorology*, 40(11) :1801–1820, 2001. [156](#)
- V. Lakshmi, E. Wood, and B. Choudhury. A soil-canopy-atmosphere model for use in satellite microwave remote sensing. *Journal of Geophysical Research*, 102(D6) :6911, 1997. [41](#)
- J. Y. Lettvin, H. R. Maturana, W. S. McCulloch, and W. H. Pitts. What the frog’s eye tells the frog’s brain. *Proceedings of the Institute of Radio Engineers*, 47(11) :1940–1951, November 1959. doi : 10.1109/JRPROC.1959.287207. [50](#)
- J. Li, W. Menzel, W. Zhang, F. Sun, T. Schmit, J. Gurka, and E. Weisz. Synergistic use of MODIS and AIRS in a variational retrieval of cloud parameters. *Journal of Applied Meteorology*, 43(11) :1619–1634, 2004. [122](#)
- J. Li, J. Li, E. Weize, and D. K. Zhou. Physical retrieval of surface emissivity spectrum from hyperspectral infrared radiances. *Geophysical Research Letters*, 34 : 10.1029/2007GL030543, 2007. [19](#)
- Z. Li, J. Li, X. Jin, T. J. Schmit, and E. Borbas. An objective methodology for infrared land surface emissivity evaluation. *Journal of Geophysical Research*, 115(D22308), November 2010. doi : 10.1029/2010JD014249. [71](#)
- Z. Li, J. Li, Y. Li, Y. Zhang, T. Schmit, L. Zhou, M. Goldberg, and W. Menzel. Determining diurnal variations of land surface emissivity from geostationary satellites. *Journal of Geophysical Research (Atmospheres)*, 117(D23), 2012. [184](#)
- Q. Liu, F. Weng, and S. English. An improved fast microwave water emissivity model. *IEEE Transactions on Geoscience and Remote Sensing*, 49(4) :1238–1250, 2011. doi : 10.1109/TGRS.2010.2064779. [41](#)
- S. Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2) :129–137, 1982. [124](#)
- P. Mahalanobis. On the generalised distance in statistics. *Proceedings of the National Institute of Science of India*, 2(1) :49–55, 1936. [140](#)
- R. F. Marsden. A proposal for a neutral regression. *Journal of Atmospheric and Oceanic Technology*, 16 :876–883, 1999. [80](#)
- G. Masiello and C. Serio. Simultaneous physical retrieval of surface emissivity spectrum and atmospheric parameters from infrared atmospheric sounder interferometer spectral radiances. *Applied Optics*, 52(11) :2428–2446, 2013. doi : 10.1364/AO.52.002428. [104](#)

- 
- M. Matricardi. Technical Note : An assessment of the accuracy of the RTTOV fast radiative transfer model using IASI data. *Atmospheric Chemistry and Physics*, 9(2), 2009. [65](#)
- M. Matricardi, F. Chevallier, G. Kelly, and J.-N. Thépaut. An improved general fast radiative transfer model for the assimilation of radiance observations. *Quarterly Journal of the Royal Meteorological Society*, 130(596) :153–173, January 2004. doi : 10.1256/qj.02.181. [33](#)
- W. Menke. *Geophysical data analysis : Discrete inverse theory*. Academic Press, New York, 1984. [129](#)
- G. Mie. Beiträge zur Optik trüber Medien, speziell kolloidaler Metallösungen. *Annalen der Physik*, 330(3) :377–445, 1908. doi : 10.1002/andp.19083300302. [12](#)
- M. Mira, E. Valor, R. Boluda, V. Caselles, and C. Coll. Influence of soil water content on the thermal infrared emissivity of bare soils : Implication for land surface temperature determination. *Journal of Geophysical Research*, 112(F04003), 2007. doi : 10.1029/2007JF000749. [76](#), [93](#)
- T. Mo. Prelaunch calibration of the Advanced Microwave Sounding Unit-A for NOAA-K. *IEEE Transactions on Microwave Theory and Techniques*, 44(8) :1460–1469, August 1996. doi : 10.1109/22.536029. [31](#)
- M. Paul and F. Aires. Using Shannon’s entropy to sample heterogeneous and high-dimensional atmospheric dataset. *Quarterly Journal of the Royal Meteorological Society*, Submitted, 2013a. [157](#), [202](#)
- M. Paul and F. Aires. Infrared and microwave synergy for the retrieval of atmospheric profiles over land - Impact of surface emissivities and temperature. *Journal of Geophysical Research*, In preparation, 2013b. [185](#), [204](#)
- M. Paul, F. Aires, C. Prigent, I. Trigo, and F. Bernardo. An innovative physical scheme to retrieve simultaneously surface temperature and emissivities using high spectral infrared observations from IASI. *Journal of Geophysical Research*, 117(D11), 2012. [19](#), [40](#), [104](#), [203](#)
- E. Pavelin and B. Candy. Assimilation of surface-sensitive infrared radiances over land : Estimation of land surface temperature and emissivity. *Quarterly Journal of the Royal Meteorological Society*, 2013. doi : 10.1002/qj.2218. [194](#)
- J. Pearl. *Causality : models, reasoning and inference*, volume 29. MIT Press, Cambridge, 2000. [147](#)

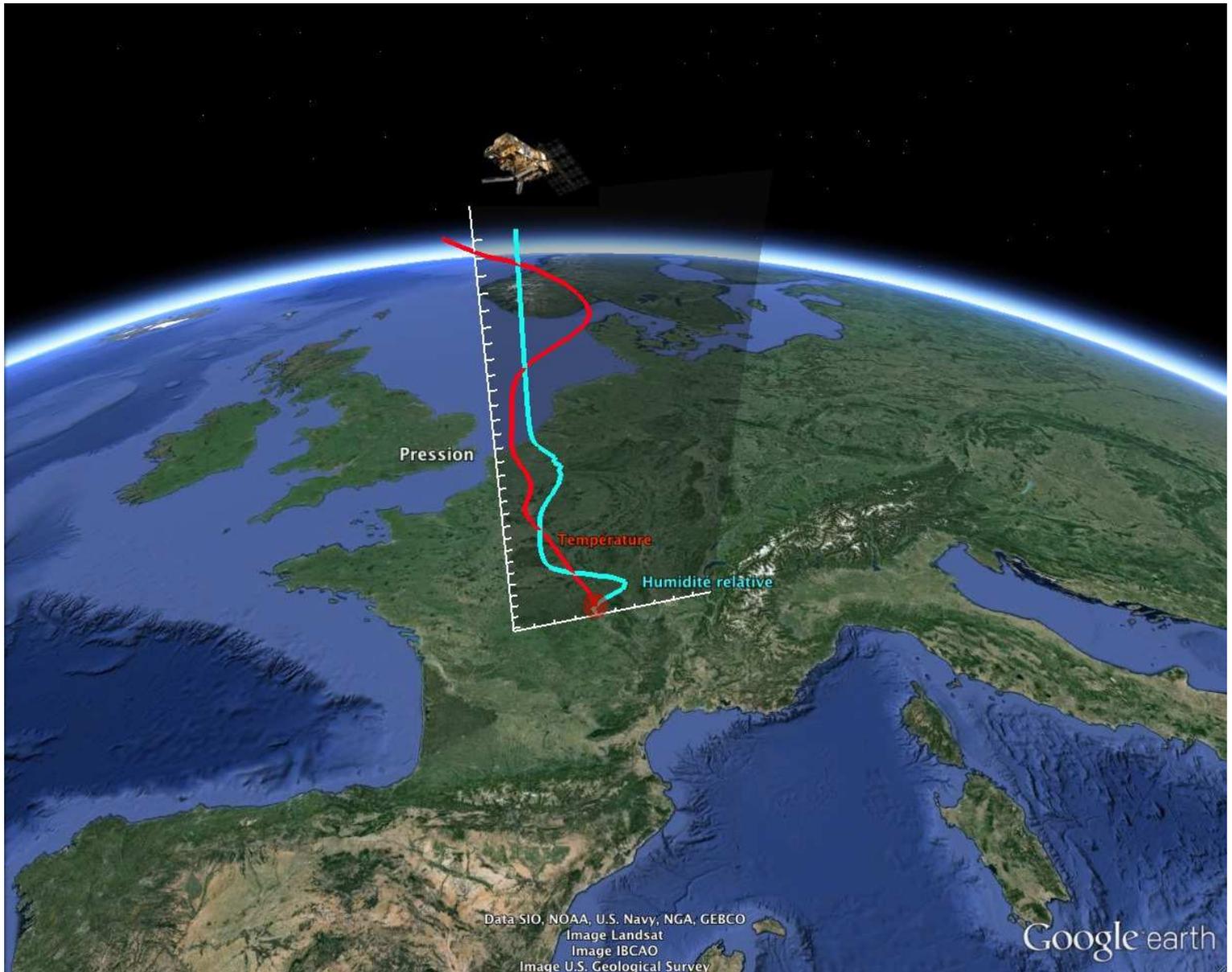
- E. Pequignot. *Détermination de l'émissivité et de la température des surfaces continentales. Application aux sondeurs spatiaux HIRS et AIRS/IASI*. PhD thesis, Ecole Polytechnique, 2006. [39](#), [48](#)
- E. Pequignot, A. Chedin, and N. Scott. Infrared continental surface emissivity spectra retrieved from AIRS hyperspectral sensor. *Journal of Applied Meteorology and Climatology*, 47 :1619–1633, 2008. [48](#)
- C. Prigent, W. Rossow, E. Matthews, and B. Marticonera. Microwave radiometric signatures of different surface types in deserts. *Journal of Geophysical Research*, 104(12) :147–12, 1999. [42](#)
- C. Prigent, F. Aires, and W. Rossow. Land surface microwave emissivities over the globe for a decade. *Bulletin of the American Meteorological Society*, 87(11) :1572–1584, 2006. doi : 10.1175/BAMS-87-11-1573. [39](#), [41](#), [183](#)
- C. Prigent, E. Jaumouille, F. Chevallier, and F. Aires. A parameterization of the microwave land surface emissivity between 19 and 100 GHz, anchored to satellite-derived estimates. *IEEE Transactions on Geoscience and Remote Sensing*, 46 :344–352, 2008. [40](#), [42](#)
- P. Prunet, J.-N. Thépaut, and V. Cassé. The information content of clear sky IASI radiances and their potential for numerical weather prediction. *Quarterly Journal of the Royal Meteorological Society*, 124(545) :211–241, 1998. [129](#)
- F. Rabier, N. Fourrié, D. Chafäi, and P. Prunet. Channel selection methods for infrared atmospheric sounding interferometer radiances. *Quarterly Journal of the Royal Meteorological Society*, 128(581) :1011–1027, 2002. [129](#)
- W. Rankine. On an equation between the temperature and the maximum elasticity of steam and other vapours. *Edinburgh New Philosophical Journal*, 28, 1849. [20](#)
- P. Ricaud, B. Gabard, S. Derrien, J-L. Attié, T. Rose, and H. Czekala. Validation of tropospheric water vapor as measured by the 183 GHz HAMSTRAD Radiometer over the Pyrenees Mountains, France. *IEEE Transactions on Geoscience and Remote Sensing*, 48(5) :2189–2203, 2010a. [98](#)
- P. Ricaud, B. Gabard, S. Derrien, J-P. Chaboureau, T. Rose, A. Mombauer, and H. Czekala. A 183 GHz Radiometer Dedicated to Sound Tropospheric Water Vapor over Concordia Station, Antarctica. *IEEE Transactions on Geoscience and Remote Sensing*, 48(3) :1365–1380, 2010b. [98](#)
- P. Ricaud, C. Genthon, J-L. Attié, F. Caminati, G. Canut, P. Durand, J-F. Vanacker, L. Moggio, Y. Courcoux, A. Pelligrini, and T. Rose. Summer to Winter Diurnal Variabilities of Temperature and Water Vapour in the lowermost troposphere as observed by

- HAMSTRAD over Dome C, Antarctica. *Boundary-Layer Meteorology*, 143(1) :227–259, 2012. doi : 10.1007/s10546-011-9673-6. [98](#)
- K. De Ridder. Surface soil moisture monitoring over Europe using Spectral Sensor Microwave/Imager (SSM/I) imagery. *Journal of Geophysical Research*, 108(D14) :4422, 2003. doi : 10.1029/2002JD002796. [41](#)
- C. D. Rodgers. Characterization and error analysis of profiles retrieved from remote sounding measurements. *Journal of Geophysical Research*, 95 :5587–5595, April 1990. doi : 10.1029/JD095iD05p05587. [213](#)
- C. D. Rodgers. *Inverse methods for atmospheric sounding : Theory and practise*. World Scientific Publishing, 1<sup>st</sup> edition, 2000. [51](#), [129](#), [213](#)
- W. Rossow and R. Schiffer. Advances in understanding clouds from ISCCP. *Bulletin of the American Meteorological Society*, 80 :2261–2288, 1999. [135](#)
- D. E. Rumelhart, G. E. Hinton, and R. J. Williams. *Learning Internal Representations by Error Propagation*. MIT Press, Cambridge, 1986. [51](#), [147](#), [187](#), [215](#)
- B. Rydberg, P. Eriksson, S. Bühler, and D. Murtagh. Non-Gaussian Bayesian retrieval of tropical upper tropospheric cloud ice and water vapour from Odin-SMR measurements. *Atmospheric Measurements Techniques*, 2(2) :621–637, 2009. [140](#)
- J. Salisbury and D. D’Aria. Emissivity of terrestrial materials in the 8-14  $\mu\text{m}$  atmospheric window. *Remote Sensing of Environnement*, 42(2) :83–106, 1992. [42](#)
- J. Salisbury and D. D’Aria. Emissivity of terrestrial materials in the 3-5  $\mu\text{m}$  atmospheric window. *Remote Sensing of Environnement*, 47(3) :345–361, 1994. [42](#)
- J. Salisbury, A. Wald, and D. D’Aria. Thermal-infrared remote sensing and Kirchhoff’s law 1. Laboratory measurements. *Journal of Geophysical Research*, 99 :11897–11911, 1994. [44](#)
- R. Saunders, M. Matricardi, and P. Brunel. An improved fast radiative transfer model for assimilation of satellite radiance observations. *Quarterly Journal of the Royal Meteorological Society*, 125 :1407–1425, April 1999. doi : 10.1256/smsqj.55614. [33](#)
- R. Saunders, J. Hocking, P. Rayer, M. Matricardi, A. Geer, N. Bormann, P. Brunel, F. Karbou, and F. Aires. RTTOV 10 science and validation report. Technical report, NWP SAF, 2012. [34](#)
- P. Schlüssel, T. Hultberg, P. Phillips, T. August, and X. Calbet. The operational IASI level 2 processor. *Advances in space research*, 36(5) :982–988, 2005. [104](#)

- S. W. Seemann, E. Borbas, R. O. Knuteson, G. R. Stephenson, and H.-L. Huang. Development of a global infrared surface emissivity database for application to clear sky retrievals from multispectral satellite radiance measurements. *Journal of Applied Meteorology and Climatology*, 47 :108–123, 2008. [34](#), [46](#), [47](#), [49](#), [184](#)
- C. Shannon. Communication in the presence of noise. *Proceedings of the IRE*, 37(1) :10–21, 1949. [129](#)
- C. Shannon. Prediction and entropy of printed English. *Bell System Technical Journal*, 30 (1) :50–64, 1951. [158](#)
- C. Shannon, W. Weaver, R. Blahut, and B. Hajek. *The mathematical theory of communication*, volume 117. University of Illinois press Urbana, 1949. [158](#)
- A. Simmons and J. K. Gibson. *The ERA-40 project plan*. European Center for Medium-Range Weather Forecasts, 2000. [123](#), [219](#)
- J. Smith and J. Taylor. Initial cloud detection using the EOF components of high-spectral-resolution infrared sounder data. *Journal of Applied Meteorology*, 43(1) :196–210, 2004. [130](#)
- W. Snyder, Z. Wan, Y. Zhang, and Y-Z. Feng. Classification-based emissivity for land surface temperature measurement from space. *International Journal of Remote Sensing*, 19(14) :2753–2774, 1998. [103](#)
- A. Stoffelen. Toward the true near-surface wind speed : Error modeling and calibration using triple collocation. *Journal of Geophysical Research*, 103(C4) :7755–7766, April 1998. [79](#)
- J. Susskind, C. Barnet, and J. Blaisdell. Retrieval of atmospheric and surface parameters from AIRS/AMSU/HSB data in the presence of clouds. *IEEE Transactions on Geoscience and Remote Sensing*, 41(2) :390–409, 2003. [122](#)
- A. Tarantola. *Inverse problem theory. Models for data fitting and model parameter estimation*. Elsevier, Amsterdam, 1987. [213](#)
- J-C Thelen, S. Havemann, S. Newman, and J. Taylor. Hyperspectral retrieval of land surface emissivities using ARIES. *Quarterly Journal of the Royal Meteorological Society*, 135(645) :2110–2124, 2009. [104](#)
- P. Tremblin, V. Minier, N. Schneider, G. Al. Durand, M. Ashley, J. Lawrence, D. Luong-Van, J. Storey, G. An Durand, Y. Reinert, C. Veyssiere, C. Walter, P. Ade, P. Calisse, Z. Challita, E. Fossat, L. Sabbatini, A. Pellegrini, P. Ricaud, and J. Urban. Site testing for submillimetre astronomy at Dome C in Antarctica. *Astronomy and Astrophysics*, 535 A112, 2011. doi : 10.1051/0004-6361/201117345. [98](#)

- I. Trigo, C. Dacamara, P. Viterbo, J. L. Roujean, F. Olesen, C. Barroso, F. Camacho-De-Coca, D. Carrer, S. Freitas, J. Garcia-Haro, B. Geiger, F. Gellens-Meulenberghs, N. Ghilain, J. Melia, L. Pessanha, N. Siljamoy, and A. Arboleda. The satellite application facility for land surface analysis. *International Journal of Remote Sensing*, 32(10) :2725–2744, May 2011. doi : 10.1080/01431161003743199. 76
- B-J. Tsuang, M-D. Chou, Y. Zhang, A. Roesch, and K. Yang. Evaluations of land-ocean skin temperatures of the ISCCP satellite retrievals and the NCEP and ERA reanalyses. *Journal of Climate*, 21(2) :308–330, 2008. 66
- S. Uppala, P. Kallberg, A. Simmons, U. Andrae, V. Da Costa Bechtold, M. Fiorino, J. Gibson, J. Haseler, A. Hernandez, G. Kelly, X. Li, K. Onogi, S. Saarinen, N. Sokka, R. Allan, E. Andersson, K. Arpe, M. Balmaseda, A. Beljaars, L. Van De Berg, J. Bidlot, N. Bormann, S. Caires, F. Chevallier, A. Dethof, M. Dragosavac, M. Fisher, M. Fuentes, S. Hagemann, E. Holm, B. Hoskins, L. Isaksen, P. Janssen, R. Jenne, A. McNally, J-F. Mahfouf, J-J. Morcrette, N. Rayner, R. Saunders, P. Simon, A. Sterl, K. Trenberth, A. Untch, D. Vasiljevic, P. Viterbo, and J. Woollen. The era-40 re-analysis. *Quarterly Journal of the Royal Meteorological Society*, 131 :2961–3012, 2005. 123, 219
- H. C. van de Hulst. *Light scattering by small particles*. Dover Publications, New York, 1981. 12
- K. Vinnikov, A. Robock, S. Qiu, J. Entin, M. Owe, B. Choudhury, S. Hollinger, and E. Njoku. Satellite remote sensing of soil moisture in Illinois, United States. *Journal of Geophysical Research*, 104(D4) :4145–4168, 1999. 41
- R. L. Vogel, Q. Liu, Y. Han, and F. Weng. Evaluating a satellite-derived global infrared land surface emissivity data set for use in radiative transfer modeling. *Journal of Geophysical Research*, 116(D08105), 2011. doi : 10.1029/2010JD014679. 39, 72, 74
- Z. Wan. New refinements and validation of the MODIS Land-Surface Temperature/Emissivity products. *Remote Sensing of Environnement*, 112(1) : 10.1016/j.rse.2006.06.026, 2008. 44
- Z. Wan and J. Dozier. A generalised split-window algorithm for retrieving land-surface temperature from space. *IEEE Transactions on Geoscience and Remote Sensing*, 34 : 892–905, 1996. 76
- Z. Wan and Z.-L. Li. A physics-based algorithm for retrieving land-surface emissivity and temperature from EOS/MODIS data. *IEEE Transactions on Geoscience and Remote Sensing*, 35 :980–996, July 1997. doi : 10.1109/36.602541. 39, 44
- J. Welsh. An Account of Meteorological Observations in Four Balloon Ascents, Made under the Direction of the Kew Observatory Committee of the British Association for the Ad-

- vancement of Science. *Philosophical Transactions of the Royal Society of London*, 143 : 311–346, 1853. [6](#)
- F. Weng, B. Yan, and N. C. Grody. A microwave land emissivity model. *Journal of Geophysical Research*, 106(D17) :20115–20113, 2001. doi : 10.1029/2001JD900019. [38](#)
- Yunyue Yu, Jeffrey L. Privette, and Ana C. Pinheiro. Evaluation of split-window land surface temperature algorithms for generating climate data records. *IEEE Transactions on Geoscience and Remote Sensing*, 46(1) :179–192, january 2008. [38](#)
- D. Zhou, R. Dickinson, K. Ogawa, Y. Tian, M. Jin and T. Schmugge, and E. Tsvetsinskaya. Relations between albedos and emissivities from MODIS and ASTER data over North African Desert. *Geophysical Research Letters*, 30(20), 2003. [43](#)
- D. Zhou, A. Larar, X. Liu, W. Smith, L. Strow, P. Yang, P. Schlüssel, and X. Calbet. Global land surface emissivity retrieved from satellite ultraspectral IR measurements. *IEEE Transactions on Geoscience and Remote Sensing*, 49 :1277–1290, 2011. doi : 10.1109/TGRS.2010.2051036. [38](#), [48](#), [71](#), [104](#)
- J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3) :337–343, 1977. [130](#)



---

Écrit par Maxime PAUL, du 01/09/2010 au 31/08/2013. Édité et imprimé par le service de reprographie de l'Observatoire de Paris à Meudon le 10/07/2013. Ce document a été préparé grâce au logiciel de composition typographique L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>, du logiciel de calcul et de rendu graphique MATLAB, du logiciel de cartographie tridimensionnelle GOOGLE EARTH et du logiciel de création picturale GIMP.

# Synergie infrarouge et micro-onde pour la restitution atmosphérique

---

## Résumé

---

L'étude du climat et la météorologie nécessitent des modèles, mais également des bases de données indépendantes, issues d'observations *in situ* ou satellites. L'importance de ces dernières grandit. Nous proposons, dans cette thèse, d'optimiser leur utilisation pour restituer, à l'échelle globale, des profils atmosphériques de température et de vapeur d'eau. La contribution des surfaces terrestres sur le rayonnement conditionne la qualité des restitutions au-dessus des continents. Un schéma d'inversion bayésienne de l'équation de transfert radiatif a donc été mis au point. Il permet de restituer simultanément la température et l'émissivité hyper-spectrale infrarouge de la surface, à partir des mesures de l'instrument IASI. Une chaîne opérationnelle a été construite, entraînant la création d'une base de données d'émissivités infrarouges et de températures de surface depuis 2007. Ces informations de surface sont ensuite intégrées dans un algorithme de restitution de profils atmosphériques. Elles permettent une diminution notable de l'erreur, notamment dans les basses couches de l'atmosphère, cruciales en météorologie. Les réseaux de neurones utilisés pour les restitutions nécessitent des bases d'apprentissage. Nous avons donc mis au point une méthode d'échantillonnage de variables hétérogènes et de grande dimension. Enfin, nous avons montré que l'utilisation conjointe des observations infrarouges et micro-ondes est une source prometteuse d'amélioration de la télédétection satellite. La synergie entre des instruments tels IASI, AMSU-A et MHS sur la plateforme MetOp permet de mieux restituer les profils atmosphériques.

**Mots clés :** télédétection, synergie, émissivité de surface, analyse statistique de données, réseau de neurones, échantillonnage

---

## Infrared and microwave synergy applied to atmospheric retrievals

---

### Abstract

---

Climatology and meteorology are mainly based on numerical models, but they also need independent data from *in situ* measurements or satellite observations. In this thesis, we attempt to optimize the use of the satellite observations in order to globally retrieve atmospheric profiles of temperature and water vapor. Knowledge about the impact of land surfaces on the radiation measured by a satellite is crucial to be able to determine the quality of the retrieved profiles. A Bayesian estimator has been used to invert the radiative equation, leading to a simultaneous retrieval of surface temperature and emissivity in the infrared, based on IASI measurements. An operational algorithm has been built. It has allowed the creation of a surface emissivity and temperature database from 2007 to today. Those surface retrievals have been used in an atmospheric inversion scheme, which led to a global decrease in the error on the retrieval of temperature and water vapor profiles, especially in the troposphere, which is the most important in meteorology. The neural network-based algorithm used for the retrievals needs a representative learning database. To build such datasets, we created a multi-variate sampling method able to compute numerous non-homogeneous variables. Finally, we have shown that the simultaneous use of infrared and microwave observations is a promising way to increase the quality of the satellite retrievals. The synergy between instruments like IASI, AMSU-A and MHS on board MetOp decreases the error of the retrieved atmospheric profiles of temperature and water vapor.

**Keywords :** remote sensing, synergy, surface emissivity, statistical data analysis, neural network, sampling