



HAL
open science

Approche ensembliste et par logique floue pour le diagnostic causal de procédés de raffinage. Application à un pilote de FCC

Bruno Heim

► **To cite this version:**

Bruno Heim. Approche ensembliste et par logique floue pour le diagnostic causal de procédés de raffinage. Application à un pilote de FCC. Automatique / Robotique. Institut National Polytechnique de Grenoble - INPG, 2003. Français. NNT: . tel-00197531

HAL Id: tel-00197531

<https://theses.hal.science/tel-00197531>

Submitted on 14 Dec 2007

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.

Remerciements

Je tiens à remercier toutes les personnes que j'ai rencontrées au cours de ces trois années.

Plus précisément, je tiens à remercier les membre du jury.

Je remercie les personnes qui m'ont encadré pendant ma thèse, qui m'ont soutenu, aidé et transmis leurs connaissances. Je remercie Madame Sylviane GENTIL, professeur au LAG-INPG et directrice de thèse ainsi que Madame Louise TRAVÉ-MASSUYÈS, directeur de recherche au LAAS-CNRS et co-encadrant de la thèse. Je remercie aussi Madame Sylvie CAUVIN-DELAIDE, ingénieur de recherche à l'IFP et Monsieur Benoît CELSE, ingénieur de recherche à l'IFP. Je remercie tous les membres de mon jury qui m'ont apporté un regard neuf, intéressant et instructif sur mes travaux. Je remercie Madame Suzanne LESECQ, maître de conférence au LAG-UJF, Monsieur Belkacem OULD-BOUAMAMA, maître de conférence au LAIL-USTL, Monsieur Luis PUIGJANER, professeur au TQG-UPC, rapporteur de la thèse, et Monsieur José RAGOT, professeur à l'INPL, rapporteur de la thèse et président du jury.

Je remercie toutes les personnes en charge de l'unité pilote de FCC avec qui j'ai eu de nombreux contacts particulièrement instructifs.

Je remercie Messieurs Hervé CAUFFRIEZ, Jean-Luc DUPLAN, Stéphane GIRARDON, Marc REYMOND et Jan VERSTRAETE. Je remercie aussi tous les opérateurs qui ont apporté un regard pratique et critique sur mes travaux.

Je remercie finalement toutes les personnes de l'IFP qui m'ont apporté un support technique lors de mes travaux.

Chapitre 1 Introduction Générale	1
1.1 La supervision et le diagnostic de procédés	2
1.2 Objectifs de la thèse	5
1.3 Réalisation et cadre de l'étude	5
1.3.1 Introduction	5
1.3.2 Le modèle causal, support du module de diagnostic	7
1.3.3 Le module de génération d'alarmes	8
1.3.4 Le module de localisation	10
1.3.5 Le module d'identification de défauts	10
1.3.6 Conclusion	11
1.4 Organisation du mémoire	12
1.5 Brève revue bibliographique	13
1.5.1 Introduction	13
1.5.2 Approche basée sur le traitement du signal	13
1.5.3 Approche basée sur des modèles quantitatifs	15
1.5.4 Approches basées sur des modèles qualitatifs ou semi-qualitatifs	18
1.5.5 Conclusion	20

Chapitre 2	Modélisation causale	21
2.1	Introduction	22
2.2	Modèle structurel des relations	27
2.2.1	Introduction	27
2.2.2	Détermination du système physique Σ	28
2.2.3	Détermination des composants qui constituent Σ	29
2.2.4	Identification de la configuration de Σ	29
2.2.5	Identification des variables décrivant Σ	29
2.2.6	Expression formelle des relations	31
2.1.7	Composants, supports des relations	32
2.1.8	Conditions d'activation et configuration de Σ	32
2.1.9	Exemple de description d'un système	32
2.1.10	Exemples de relations de redondance analytique	34
2.1.11	Conclusion	35
2.3	Modèle causal structurel	36
2.3.1	Introduction	36
2.3.2	Obtention de la causalité selon Iwasaki et Simon	37
2.1.3	Obtention de la causalité via la théorie des graphes	44
2.1.4	Bilan sur la méthode d'obtention de la causalité	55
2.1.5	Information contenue dans le modèle causal structurel	56
2.1.6	Systèmes à plusieurs modes de fonctionnement	57
2.1.7	Conclusion	58
2.4	Modèle causal approché	59
2.4.1	Introduction	59
2.4.2	Réduction	60
2.4.3	Approximation	61
2.4.4	Conclusion	61
2.5	Quantification du modèle	62
2.5.1	Introduction	62
2.5.2	Représentation du modèle	63
2.1.3	Identification des paramètres du modèle	63
2.1.4	Conclusion	64
2.6	Conclusion	64

Chapitre 3	Méthode de diagnostic	67
3.1	Introduction	68
3.2	Génération des références	70
3.2.1	Introduction	70
3.2.2	Références globales	70
3.2.3	Génération des références locales	72
3.2.4	Conclusion	74
3.3	Génération des alarmes	75
3.3.1	Introduction	75
3.3.2	Les alarmes globales	75
3.3.3	Génération des alarmes locales	76
3.3.4	Conclusion	78
3.4	Techniques de détection	78
3.4.1	Introduction	78
3.4.2	Logique binaire	79
3.4.3	Logique floue	81
3.4.4	Méthode ensembliste	86
3.4.5	Conclusion	95
3.5	Comparaison des techniques de détection	96
3.5.1	Introduction	96
3.5.2	Comparaison qualitative des techniques	96
3.5.3	Comparaison des techniques sur des données simulées	100
3.5.4	Complémentarité des techniques	113
3.5.5	Conclusion	117
3.6	Localisation des défauts	117
3.6.1	Introduction	117
3.6.2	Stratégies de diagnostic	118
3.6.3	Conclusion	123
3.7	Identification des défauts	124
3.7.1	Introduction	124
3.7.2	Modèles qualitatifs de fonctionnement anormal	124
3.7.3	Conclusion	129
3.8	Conclusion	129

Chapitre 4	Mise en œuvre sur un pilote de FCC	127
4.1	Le procédé de craquage catalytique	132
4.1.1	Introduction	132
4.1.2	Unités industrielles de FCC	133
4.1.3	Le pilote de FCC du CEDI	134
4.1.4	Conclusion	137
4.2	Les dysfonctionnements du procédé de FCC	137
4.2.1	Introduction	137
4.2.2	Problèmes rencontrés sur les unités industrielles de FCC	138
4.2.3	Problèmes rencontrés sur le pilote de FCC du CEDI	139
4.2.4	Conclusion	110
4.3	Application au pilote de FCC du CEDI	140
4.3.1	Introduction	140
4.3.2	Le modèle causal du pilote de FCC	140
4.3.3	Résultats obtenus sur le pilote de FCC	144
4.3.4	Problèmes rencontrés et apports des travaux au domaine	152
4.3.5	Contexte de l'étude et du développement du module informatique	155
4.3.6	Résultats escomptés sur un FCC industriel	160
4.3.7	Conclusion	162
Chapitre 5	Conclusion et perspectives	161
Bibliographie		171
Annexe A	Algorithmes de réduction et d'approximation	181
Annexe B	Identification des paramètres	191
Annexe C	Ordre de calcul des noeuds	209

Nomenclature

Ensembles et Abréviations

A	Arcs dans G	$V_{\text{exo.causal}}$	Variables exogènes du MCS
A'	Arcs dans G'	$V_{\text{pseudo.exo}}$	Variables pseudo-exogènes du MCS
$COMPS$	Composants industriels de Σ	V^S	Capteur de la variable V
E	Relations du MCS	z	Opérateur de décalage
E_{Dyn}	Equations différentielles de E	z_k	Un instant discret
G	Graphe biparti	$\lambda_{\text{variable}}$	Résidu local de la <i>variable</i>
G'	Graphe orienté	ρ_{variable}	Résidu global de la <i>variable</i>
I	Influences dans le MCS	Σ	Système physique
n_E	Nombre de relations dans E	ASCO	Aide à la Supervision et à la Conduite pour les Opérateurs.
N_V	$\{1 \dots n_v\}$	MCA	Modèle Causal Approché
n_V	Nombre de variables dans V	MCR	Modèle Causal Réduit
OBS	Ensemble des variables connues	MCS	Modèle Causal Structurel
R	Relations de MCS	FCC	Fluid Catalytic Cracking
V	Variables décrivant Σ	MCS	Minimal Complete Subsystem
V_{Dyn}	Variables dérivées dans E	MSR	Modèle Structurel des Relations
V_{endo}	Variables endogènes de Σ	SCADA	System of Control And Data Acquisition
$V_{\text{endo.causal}}$	Variables endogènes du MCS	SDG	Signed Directed Graph
V_{exo}	Variables exogènes de Σ	SCM	Système Complet Minimal

Applications

<i>Application</i>	Variable générique valant <i>Exist</i> ou <i>Support</i> ou <i>Relation</i>
$A(V,e)$	Un arc non orienté entre la variable V et la relation e dans G
$A'(e,V)$	Un arc orienté entre la relation e et la variable V dans G
$A'(V,e)$	Un arc orienté entre la variable V et la relation e dans G'
$P(e,V)$	Un arc du couplage parfait entre la relation e et la variable V dans G
<i>Relation(arc)</i>	La relation associée à <i>arc</i>
<i>Exist(arc)</i>	Prend la valeur <i>vrai</i> si l' <i>arc</i> existe et <i>faux</i> sinon
<i>Support(relation)</i>	Composants industriels (capteurs exclus) associés a <i>relation</i>
<i>Connue(variable)</i>	Prend la valeur <i>vrai</i> si la <i>variable</i> est connue et <i>faux</i> sinon
$I_{\text{Model}}(V,Y)$	Un arc orienté entre la variable V et la variable Y dans <i>Model</i>
<i>Model</i>	Variable générique valant MCS ou MCR ou MCA
$Rel_{V \rightarrow E, \text{couplage}}(\text{variable})$	Prend pour valeur l'identité de la relation couplée avec <i>variable</i>
$Var_{E \rightarrow V, \text{intervenant}}(\text{relation})$	Ensemble des variables intervenant dans <i>relation</i>
$Var_{E \rightarrow V, \text{int ervenant, max. MCS}}(\text{relation})$	Variables intervenant et appartenant au MCS d'ordre le plus élevé
$Var_{E \rightarrow V, \text{couplage}}(\text{relation})$	Prend pour valeur l'identité de la variable couplée avec <i>relation</i>

$Var_{V \rightarrow V, causes}(variable)$	Ensemble des variables causes de <i>variable</i>
$Var_{V \rightarrow V, causes.backward}(variable, n)$	Ensemble des variables causes de <i>variable</i> de profondeur <i>n</i>
$Var_{V \rightarrow V, causes.directes}(variable)$	Ensemble des variables causes directes de <i>variable</i>
$Var_{V \rightarrow V, causes. \rightarrow directes}(variable)$	Ensemble des variables causes indirectes de <i>variable</i>
$Var_{V \rightarrow V, consequences.directes}(variable)$	Ensemble des variables conséquences de <i>variable</i>
$Capteur(variable)$	Prend pour valeur l'identité du capteur si la <i>variable</i> est mesurée sinon \emptyset
$Support_{influence}(arc)$	Ensemble de composants industriels et de capteurs associés à l' <i>arc</i>
$Support_{Relation}(relation)$	Ensemble de composants industriels et de capteurs associés à la <i>relation</i>

Variables et composants du procédé :

$P_{composant}$	Pression en ciel du <i>composant</i>
$F_{composant}$	Débit dans le <i>composant</i>
$\Delta P_{composant}$	Perte de charge aux bornes du <i>composant</i>
$L_{composant}$	Niveau de catalyseur dans le <i>composant</i>
$SP_{variable_régulée}$	Consigne de la <i>variable régulée</i>
$RV_{variable_régulée}$	Régulateur de la <i>variable régulée</i>

Chapitre 1

Introduction Générale

Ce chapitre pose le problème de la supervision et du diagnostic de procédés industriels. Il présente ensuite successivement :

- les objectifs de la thèse,
- le cadre de l'étude et les travaux réalisés,
- l'organisation du mémoire,
- une brève revue bibliographique des techniques de diagnostic.

1.1 La supervision et le diagnostic de procédés

La taille des installations industrielles, leur complexité et le nombre croissant de données à traiter rendent de plus en plus difficile la compréhension de leur fonctionnement. Dans ce contexte et en fonctionnement dégradé, les opérateurs doivent rapidement prendre une décision et mettre en œuvre des actions appropriées.

L'objectif de ce manuscrit est de proposer une méthode permettant de diagnostiquer les procédés de raffinage et de l'appliquer à un pilote de craquage catalytique. Le module informatique, que nous avons développé et testé sur des données réelles, aide l'opérateur à prendre sa décision en situation de crise.

L'objectif de la *supervision* est de surveiller et de contrôler le fonctionnement d'une installation pour qu'elle reste dans la plage de fonctionnement visée quelles que soient les perturbations extérieures. Elle est essentiellement effectuée par les opérateurs dans la salle de contrôle, véritable centre nerveux du complexe industriel [Bullemer et Nimmo 1996]. La supervision comprend la *surveillance*, la prise de décision et la mise en œuvre des actions appropriées pour maintenir l'opération (en mode normal ou en mode dégradé) et éviter des dégâts matériels ou humains. Certains auteurs développent des méthodes qui permettent de reconfigurer des procédés [Diego *et al.* 2001], [Sanmarti, Puigjaner et Friedler 1998].

Les informations nécessaires à la supervision sont traitées par des calculateurs numériques sophistiqués ou le SCADA¹⁻¹ qui inclut l'interface opérateur et des systèmes automatiques de régulation et d'arrêts d'urgence. Leur fonction n'est pas uniquement l'acquisition de données mais aussi la mise en œuvre de processus automatiques de *contrôle commande* et de *surveillance*.

La *commande* consiste à faire varier les consignes des régulateurs pour atteindre les objectifs de rendement, de qualité ou de productivité. La plus classique s'appuie sur des consignes locales pouvant être déterminées grâce à l'expérience de l'ingénieur

¹⁻¹ System of Control And Data Acquisition

de production ou par des outils d'optimisation. Les systèmes de commande garantissent la stabilité du procédé et sont très efficaces dans une grande partie des situations. La commande avancée intègre les techniques multivariables et propose des régulateurs robustes aux variations faibles des paramètres du procédé.

La *surveillance* de la majeure partie des procédés industriels se limite à des systèmes de traitement d'*alarmes*. Elles sont déclenchées lorsque les mesures de certaines variables clés sortent de la plage de fonctionnement souhaitée. Un *seuil haut* SH et un *seuil bas* SB sont fixés pour chaque variable et à chaque acquisition le système de surveillance vérifie si la mesure se situe dans l'intervalle [SB,SH]. Le *défaut* correspond alors à un écart inacceptable de la valeur de la variable par rapport à cette valeur souhaitée. Les valeurs des seuils SB et SH sont fixées par des experts du procédé selon des critères de sécurité des hommes, de l'installation et de son environnement.

Le fonctionnement dégradé peut être la conséquence d'une *défaillance* [Villemeur 1988] qui est une modification suffisante et permanente des caractéristiques d'un composant industriel pour qu'une fonction requise ne puisse plus être assurée dans les conditions souhaitées. Il peut être observé à travers un indicateur de défaut, grandeur qui s'écarte d'une valeur de référence sous l'effet d'un (ou de plusieurs) défaut(s). Il peut aussi être la conséquence d'une *panne* qui est une inaptitude d'un composant industriel à accomplir sa fonction requise. Certaines défaillances peuvent nécessiter soit d'arrêter immédiatement l'installation, soit de basculer la commande vers un mode de repli, ou encore vers un mode dégradé qui consiste à changer les consignes des boucles locales ou même la configuration de la commande ou du procédé.

Les opérateurs, confrontés à une défaillance, doivent rapidement l'identifier afin de maintenir le fonctionnement ou d'éviter toute prise de risque [Iserman et Ballé 1997]. Ils sont alors confrontés à des situations de crise et doivent déterminer l'origine des alarmes. La propagation des effets des défauts entraîne, via les phénomènes de causalité, des «cascades d'alarmes» qui rendent la supervision délicate et extrêmement critique [Cauvin 1995].

De nombreuses techniques ont été développées pour faciliter la supervision. Elles doivent permettre de diagnostiquer, c'est à dire de *détecter*, de *localiser*, *d'identifier* des défauts tout en générant des explications compréhensibles par les opérateurs. La détection consiste à déterminer si le procédé présente ou non des défauts. La localisation précise quel composant industriel est affecté par le défaut. L'identification de défaut détermine la taille, la nature et l'évolution du défaut.

L'objectif du diagnostic [Isermann et Ballé 1997] est de détecter, de localiser et d'identifier les défauts (des capteurs, des actionneurs et du procédé) dans les cas particuliers suivants :

- les défauts qui s'installent de manière progressive, abrupte ou intermittente,
- les défauts dans les boucles physiques fermées (boucles de rétroaction),
- les défauts s'installant lors de transitions entre états stationnaires.

Des applications industrielles de diagnostic ont été développés.

- Nous citons par exemple le « Système d'Aide à la Conduite des Hauts-fourneaux » élaboré par la société USINOR [Thirion et Lesaffre 1999].
- Au début des années 80, l'entreprise américaine Honeywell a lancé le groupe de travail de Gestion d'Alarme dont l'une des premières tâches a été d'associer des alarmes aux différents problèmes rencontrés sur une unité industrielle automatisée et d'assurer une gestion optimale de celle-ci. Ces travaux ont débouché sur un consortium de Gestion de Situations Anormales entre Honeywell et sept très grandes entreprises des Etats-Unis, tels qu'Amocco, Chevron, Exxon, Novaco, Shell, Texaco et deux fournisseurs de logiciels : Gensym et Elec. Dans ce cadre, nous citons les travaux [Mylaraswamy et Venkatasubramanian 1997] qui utilisent des classificateurs statistiques pour effectuer le diagnostic d'un procédé de FCC.

1.2 Objectifs de la thèse

L'objectif de la thèse est :

- de développer une méthodologie et un module informatique de supervision et de diagnostic pour des procédés de raffinage et de pétrochimie,
- d'appliquer cette méthode à un procédé pilote de raffinage, le craquage catalytique.

Ceci implique :

- de spécifier les méthodes choisies pour la modélisation puis pour la détection, la localisation et l'identification de défauts,
- d'étudier l'apport de ces méthodes.

Ce travail est conduit dans le cadre du projet Européen CHEM¹⁻². L'objectif du projet est de développer et de faire collaborer des outils de diagnostic et d'aide à la décision pour la conduite d'unités industrielles de chimie, de raffinage et de pétrochimie. Différents partenaires industriels et académiques participent à ce projet qui est plus largement décrit dans la section 4.3.5.2.

1.3 Réalisation et cadre de l'étude

1.3.1 Introduction

L'objectif de la thèse est de développer une méthode générale de supervision et de diagnostic puis de l'appliquer au pilote de FCC de l'IFP à Solaize (cf. Chapitre 4). Le pilote de FCC a été choisi comme cas d'application car sa supervision et son diagnostic sont particulièrement complexes :

- La propagation rapide des défauts dans le procédé de FCC rend difficile la localisation du défaut source. Cette localisation est rendue d'autant plus

¹⁻² CHEM: "Advanced Decision Support System for Chemical and Petrochemical Processes". Ce projet est financé par la communauté Européenne sous le "Competitive and Sustainable Growth programme of the Fifth RTD Framework Programme (1998-2002) G1RD-CT-2001-00466". (www.cordis.lu ou www.chem-dss.org).

difficile que les liens de causalité induits par les régulations en boucle fermée sont complexes.

- Le FCC contient une boucle physique qui complexifie la compréhension des phénomènes de causes à effets.
- En situation de crise, l'opérateur est très souvent confronté à une cascade d'alarmes dont l'interprétation est particulièrement complexe et requiert un temps souvent trop long de réflexion pour pouvoir agir suffisamment rapidement.

Le système **ASCO**, (**A**ide à la **S**upervision et à la **C**onduite pour les **O**opérateurs), que nous avons développé, fonctionne en boucle ouverte : la décision puis l'action sont prises par l'opérateur. Le système **ASCO** a été testé en ligne sur le pilote de FCC.

Le système ASCO présente l'originalité de faire appel successivement à trois techniques complémentaires. L'enchaînement de ces techniques permet d'obtenir, à partir des observations, un message qui identifie la panne sur un composant :

- La première étape dont les principes sont décrits dans la section 1.3.3, est effectuée par le **module de génération d'alarmes**.
- Les alarmes générées sont ensuite traitées par le **module de localisation** qui élabore une liste de composants physiques, tous suspectés d'être défectueux. Les grands principes de cette étape sont décrits dans la section 1.3.4.
- Finalement, la connaissance des experts sur ces composants industriels est analysée par le **module d'identification de défauts** (cf. section 1.3.5). Ce module identifie une panne sur un composant. Cette identification de défaut se traduit par un message pour l'opérateur.

La plupart des auteurs ne proposent pas de méthode combinant la génération d'alarme, puis la localisation et l'identification de défaut : ils ne décrivent qu'une ou deux de ces étapes. Il nous a été nécessaire, en pratique sur le pilote de FCC, d'utiliser ces trois modules qui apportent toutes les informations nécessaires pour la compréhension des phénomènes et le suivi de l'unité.

Le support des deux premiers modules qui est décrit dans la section 1.3.2, est un modèle de bon fonctionnement qui est causal.

1.3.2 Le modèle causal, support du module de diagnostic

Le modèle support du système ASCO est dynamique et causal : il décrit les influences entre les variables et permet d'expliquer le comportement des variables les unes par rapport aux autres. Les modèles causaux ont déjà été utilisés dans des travaux antérieurs pour la supervision [Leyval, Gentil et Feray-Beaumont 1994].

Les modèles causaux diffèrent des modèles analytiques par la notion de séquences de causes et d'effets [Rasmussen 1993] et [Travé-Massuyès *et al.* 2001]. Ils permettent d'organiser les équations du modèle.

Un modèle causal est facilement représenté par un graphe causal. Le graphe causal est constitué de nœuds et d'arcs orientés. Les nœuds représentent les variables et les arcs sont les influences entre ces variables. La Figure 1-1 présente le graphe causal du pilote de FCC que nous avons construit. Des fonctions de transfert classiques sont associées à chaque arc du graphe causal, le rendant ainsi quantitatif.

Le graphe causal est aussi utilisé en tant qu'interface : il permet de visualiser les alarmes et la propagation des défauts sous la forme de code de couleurs. Le chapitre 2 présente la méthode de modélisation causale pour le diagnostic que nous avons développée.

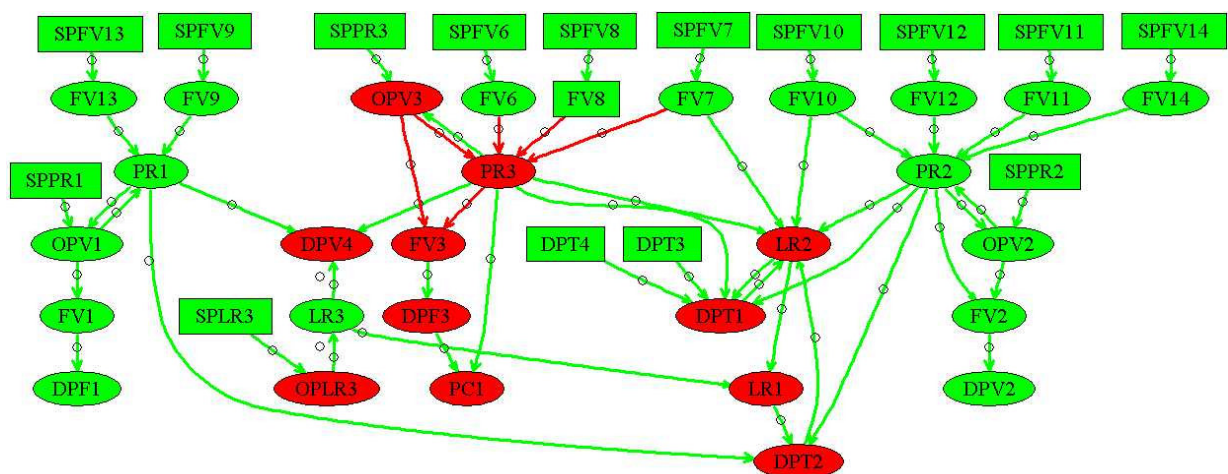


Figure 1-1: Le graphe causal du pilote de FCC

Des **résidus** spécifiques sont extraits de modèles causaux et sont utilisés comme indicateurs de défauts par le système ASCO pour la génération d'alarmes.

1.3.3 Le module de génération d'alarmes

Le **module de génération d'alarmes** compare les sorties du modèle causal avec les mesures et génère des alarmes.

Le modèle et les mesures sont souvent imprécis, par conséquent nous avons testé deux méthodes permettant de prendre en compte ces incertitudes. La première méthode est basée sur une approche par logique floue [Evsukoff 1998] et la seconde sur une approche ensembliste [Armengol 1999].

La méthode par logique floue est décrite dans la section 3.4.3 : une interprétation des résidus permet de générer une valeur d'appartenance à la classe AL (Alarme), comprise entre 0 et 1, ce qui permet de savoir si la mesure est anormale ou non. L'évolution graduelle de cette valeur de 0 à 1 représente l'évolution de la variable vers une valeur anormale. Cette évolution graduelle est exploitée par le système ASCO sous la forme d'un dégradé de couleurs des contours des nœuds du graphe causal. Le contour du nœud est rouge si un défaut est détecté ($AL = 1$). Il est vert si aucun défaut n'est détecté. Pour les valeurs intermédiaires de AL, il évolue du vert au jaune, à l'orange ... puis au rouge. L'approche par logique floue présente l'avantage de focaliser l'attention de l'opérateur sur des variables qui commencent à s'écarter de leur valeur sans toutefois être en alarme. Il peut aussi commencer à réfléchir à un problème bien avant que celui-ci ne soit clairement installé, et ainsi mieux préparer ses actions pour le contre-carrer. L'inconvénient principal de l'approche par logique floue est qu'elle s'appuie sur des seuils fixes de détection.

La Figure 1-2 montre un autre type d'information, suggéré par [Evsukoff 1998], que nous avons réalisé pour les opérateurs et qui permet de visualiser l'évolution de l'historique des couleurs du nœud PR3 sur une fenêtre temporelle glissante.

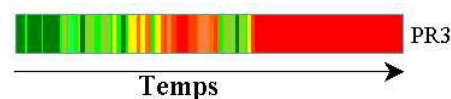


Figure 1-2 : Evolution temporelle de la couleur des nœuds du graphe causal

Nous avons aussi étudié l'apport d'une approche ensembliste pour la prise en compte des incertitudes : la sortie du modèle représente alors un intervalle contenant assurément la mesure en l'absence de défaut (cf. section 3.4.4). Une décision binaire est effectuée en étudiant l'intersection de la mesure par rapport à ces enveloppes. Les intervalles sur les paramètres du modèle et les incertitudes sur les mesures permettent de quantifier cet intervalle. Le système ASCO n'intègre pas cette méthode ensembliste mais nous avons effectué des tests des algorithmes dans le cadre du projet CHEM (cf. section 4.3.4.2). L'approche ensembliste permet de fournir des résultats garantis. Lorsque la mesure est en dehors de l'intervalle fourni par le modèle, il y a assurément un défaut.

Le principal avantage de l'approche ensembliste est de générer des enveloppes dont la taille dépend de l'excitation des entrées. L'approche par intervalles est binaire et elle garantit le résultat : s'il y a une alarme alors il y a assurément un défaut. Par contre, la décision est retardée et il y a des manques à la détection.

Les deux approches par logique floue et ensembliste, sont antagonistes et complémentaires. L'approche par logique floue ne garantit pas le résultat (des fausses détections sont probables), mais par contre elle permet d'anticiper les problèmes. La méthode ensembliste garantit le résultat mais cette contrainte engendre des manques à la détection. Nous avons donc jugé pertinent de comparer ces deux méthodes, dans un premier temps, puis d'élaborer une technique combinant ces deux méthodes dans un second temps.

Nous avons élaboré une méthode permettant de combiner le raisonnement par logique floue et l'approche ensembliste. Cette méthode, présentée dans la section 3.5.4, propose une interprétation par le raisonnement par logique floue des positions respectives de deux enveloppes générées via l'approche ensembliste. Nous n'avons pas implémenté cette méthode dans le système ASCO pour des raisons de difficulté de mise en œuvre pratique et de temps.

1.3.4 Le module de localisation

Pour effectuer la localisation des défauts nous avons suivi la démarche suivante : chaque arc du graphe causal est associé à un ensemble de composants industriels qui constituent son **support**. A partir de ces associations, et des alarmes précédemment générées, le **module de localisation** établit des diagnostics. Un diagnostic est un ensemble de composants dont toutes les entités sont suspectées d'avoir un comportement anormal.

Nous avons utilisé les méthodes proposées par [Cordier *et al.* 2000b] et [Reiter 1987] pour effectuer la localisation des défauts (cf. section 3.6). Les diagnostics sont générés à partir des **conflits** sur lesquels est appliqué un algorithme de «hitting-set». Un conflit est un ensemble de composants tels que l'on est sûr qu'au moins un d'entre eux se comporte anormalement pour expliquer les observations. Les conflits sont obtenus à partir des alarmes [Heim *et al.* 2002b].

Le système ASCO contient un module de localisation s'appuyant sur cette méthode. L'interface du système ASCO associe chaque composant à un rectangle. Un code de couleurs est utilisé pour symboliser le diagnostic. Le rectangle est initialement vert et, si le composant appartient à un diagnostic, le rectangle devient rouge. Les arcs dont le support contient ces composants apparaissent également en rouge sur le graphe causal (cf. Figure 1-1).



Figure 1-3 : Composants suspectés de se comporter anormalement

1.3.5 Le module d'identification de défauts

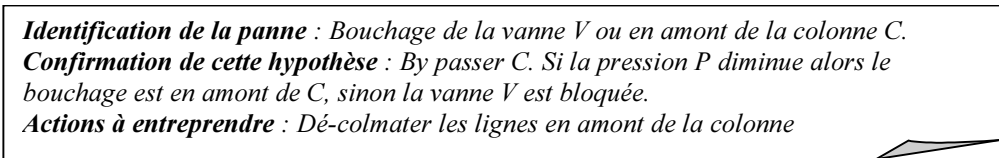
Le module ASCO s'appuie finalement sur un modèle de fonctionnement anormal (obtenu à partir de la connaissance experte) pour effectuer l'identification des défauts. Le **module d'identification de défauts** identifie la panne sur un composant et cela se traduit par un message pour l'opérateur.

Chaque composant est associé à une base de règles (ou système expert), [Heim, Cauvin et Gentil 2001]. Lorsque le composant appartient à un diagnostic, sa base de règles est activée. Chaque base de règle contient la connaissance experte sur le mauvais fonctionnement du composant. Elle associe des symptômes à des défaillances par des liens logiques.

Un traitement du signal de mesures du procédé permet de générer des symptômes (exemple : la pression P augmente, le débit D diminue). Puis si les symptômes sont observés selon des tests logiques, qui sont issus de la connaissance de l'expert (exemple : si P augmente et D diminue), alors un message qui identifie la défaillance (exemple : la vanne V est bouchée) sur un composant apparaît sur l'interface des opérateurs.

Ce message contient une phrase en langage naturel décrivant la défaillance supposée de ce composant, des actions à entreprendre pour valider l'information et pour rétablir le bon fonctionnement.

Un exemple de message du pilote de FCC est illustré sur la Figure 1-4. L'algorithme d'identification des défauts est présenté dans la section 3.7.



Identification de la panne : Bouchage de la vanne V ou en amont de la colonne C .
Confirmation de cette hypothèse : By passer C . Si la pression P diminue alors le bouchage est en amont de C , sinon la vanne V est bloquée.
Actions à entreprendre : Dé-colmater les lignes en amont de la colonne

Figure 1-4 : Identification de la défaillance

1.3.6 Conclusion

Le système **ASCO** a été développé sur le support informatique G2 de la société Gensym. Ce logiciel a été choisi car :

- il permet de développer rapidement des outils pour la supervision de procédés,
- il propose un langage de type objet déclaratif et graphique qui est simple d'utilisation,
- il a la capacité de fonctionner en temps réel.

Des tests en ligne ont permis de détecter les défauts plus rapidement que les opérateurs. Cette précocité permet de mettre en œuvre rapidement les actions pour

rétablir le bon fonctionnement. La combinaison des trois techniques complémentaires de détection, d'isolation et d'identification a donné des résultats très concluants.

L'utilisation prolongée du système **ASCO** en ligne permettra d'éprouver ses avantages et ses limites.

1.4 Organisation du mémoire

Le **second chapitre** de ce mémoire propose une brève revue bibliographique de la modélisation causale. Il présente ensuite la méthode de modélisation causale que nous avons élaborée. Comme nous l'avons appliquée au pilote de FCC, nous soulevons des problèmes concrets et proposons des solutions.

Le **troisième chapitre** présente le module de diagnostic s'appuyant sur ce modèle causal. Il présente successivement les méthodes de détection, de localisation et d'identification sur lesquelles nous nous sommes appuyés. Ces méthodes sont issues de divers travaux antérieurs mais nous avons développé la démarche qui permet de les intégrer et de les utiliser successivement. Nous présentons ensuite comment nous avons choisi de fixer les valeurs des paramètres des algorithmes de détection. Nous comparons finalement les apports de la méthode par logique floue et de la méthode ensembliste.

La **quatrième chapitre** présente les résultats que nous avons obtenus en appliquant cette méthode au pilote de FCC. Ce chapitre présente plus précisément :

- le procédé de craquage catalytique,
- une liste exhaustive des problèmes rencontrés sur l'unité pilote et sur les unités industrielles de FCC,
- le modèle causal du pilote de FCC,
- les résultats obtenus par le module de diagnostic sur des scénarios¹⁻³ de données réelles de mauvais fonctionnement,
- les problèmes rencontrés lors de ces travaux.

La section suivante présente tout d'abord une brève revue bibliographique des méthodes de diagnostic de procédés.

¹⁻³ Les scénarios ont été subis et n'ont pas été provoqués pour tester le module, bien évidemment.

1.5 Brève revue bibliographique

1.5.1 Introduction

Plusieurs méthodes de diagnostic de procédés sont distinguées et la plupart s'appuient sur un modèle. Ces modèles peuvent être de nature très différente (modèles heuristiques, statiques, dynamiques, hiérarchiques, causaux, de comportement normal ou anormal ou encore modèle du signal, etc.).

Les deux grandes approches à base de modèles sont l'approche analytique et l'approche heuristique. L'approche analytique est basée sur les mesures et les résultats des modèles analytiques qui engendrent les symptômes analytiques. L'approche heuristique utilise les symptômes représentant des informations quantitatives ou qualitatives (mesures ou expertise humaine) et une connaissance qualitative obtenue le plus souvent par des observations (sur le terrain) [Isermann et Ballé 1997].

Toutes les méthodes nécessitent une phase d'apprentissage qui permet de connaître une référence de fonctionnement normal ou dégradé. Cette référence permet de générer en temps réel les symptômes lors de la phase de diagnostic.

Les paragraphes suivants présentent les grandes familles de méthodes de diagnostic.

- Le paragraphe 1.5.2 présente des méthodes basées sur le traitement de données.
- Le paragraphe 1.5.3 présente des méthodes à base de modèles quantitatifs.
- Le paragraphe 1.5.4 présente des méthodes à base de modèles qualitatifs ou semi-qualitatifs.

1.5.2 Approche basée sur le traitement du signal

Principe :

Dans certains cas il est difficile, voire impossible, d'obtenir un modèle mathématique du (ou d'une partie du) procédé qui relie les signaux d'entrée aux

signaux de sortie. Les méthodes de traitement du signal permettent d'analyser les signaux de sortie observés et proposent pour chacun de ces signaux un modèle. Ce modèle est une référence qui peut caractériser soit un fonctionnement normal soit une défaillance particulière.

Exemples :

- Décision statistique

Une loi de densité de probabilité est obtenue en analysant les signaux pendant le fonctionnement normal. La détection consiste à déterminer si le signal obéit à cette loi ou s'en écarte significativement [Jia, Martin et Morris 1998]. Une difficulté de cette méthode est de disposer en temps réel d'un historique suffisamment grand permettant l'estimation de grandeurs statistiques telle qu'une moyenne ou un écart type. Une autre difficulté est d'être représentatif de tous les modes de fonctionnement. Il est aussi difficile d'interpréter physiquement le résultat de la détection.

- Approche fréquentielle

L'analyse spectrale des signaux de vibrations permet d'isoler des composantes spectrales : la répartition de l'énergie du signal est analysée en fonction de la fréquence. La variation de ce spectre permet de détecter des défauts voire d'identifier la défaillance. [Leseq et Barraud 2000] proposent d'utiliser des filtres d'ondelettes pour détecter des défauts respectivement dans un entraînement électrique. [Bakhtazad, Palazoglu et Romagnoli 1998] présentent une stratégie de représentation et de classement de données, en vue de la détection, en utilisant les coefficients de transformation en ondelettes. Une étude de cas porte sur deux réacteurs, à fonctionnement continu, en série et qui utilisent un mélangeur intermédiaire pour la seconde alimentation.

Les filtres d'ondelettes réalisent une projection du signal sur des espaces d'analyse possédant des propriétés particulières. Cela conduit à une présentation de l'information qui permet de mettre en évidence des comportements du signal qui n'étaient pas forcément décelables « à l'œil nu ». La difficulté de ces méthodes consiste à automatiser l'analyse du spectre ou des coefficients d'ondelettes.

- Classification de données

La classification de données permet de reconnaître des formes caractéristiques de diverses défaillances. La difficulté de cette méthode est qu'il faut disposer d'une base

de connaissance très complète du comportement défaillant [Dubuisson 2000] ou d'intégrer un système d'apprentissage en ligne. La classification peut porter sur des attributs numériques ou bien sur des informations qualitatives. Des méthodes permettent d'extraire des informations qualitatives, symptômes, à partir de données numériques. [Colomer, Mendelez et Gamero 2002] proposent de convertir la trajectoire numérique d'une variable en épisodes. Un épisode est une représentation qualitative d'un état de la variable sur un intervalle de temps donné (augmente, baisse ...).

1.5.3 Approche basée sur des modèles quantitatifs

1.5.3.1 Principe

L'utilisation de modèles mathématiques pour le diagnostic est très largement répandue [Frank 1990], [Frank 1996], [Isermann et Ballé 1997]. La comparaison de la sortie du modèle avec les mesures permet de générer un résidu constituant un indicateur de défauts. Les sorties du modèle se représentent classiquement sous la forme de valeurs numériques. Chaque relation du modèle peut être associée à un ou plusieurs composants industriels qui la sous tendent. Ils constituent le *support* de la relation [Cordier *et al* 2000a]. La localisation qui suit la détection peut être effectuée à l'aide de la *table de signature* ou *matrice(ou table) d'incidence*. Les colonnes de cette table (cf. Figure 1-5) sont représentatives des différents défauts et les lignes des différents résidus. La table est remplie avec les symptômes théoriques déduits des formes d'évaluation des résidus. Un 1 indique que le résidu est sensible au défaut, un 0 le contraire. La *signature* d'un défaut (colonne de la matrice) définit l'état des symptômes lorsque ce défaut affecte le système (l'état est une grandeur logique ou symbolique) [Busson *et al.* 1998]. [Gertler et Singer 1990] distinguent trois cas de matrices d'incidence :

- Non localisante (une colonne est nulle ou au moins deux colonnes sont identiques)
- Faiblement localisante (les colonnes sont non nulles et distinctes deux à deux),
- Fortement localisante (en plus d'être faiblement localisante, aucune colonne ne peut être obtenue à partir d'une autre en remplaçant un '1' par un '0').

	f_1	f_2	f_3
r_1	1	1	1
r_2	1	1	1
r_3	1	0	0

Non localisante

	f_1	f_2	f_3
r_1	1	1	1
r_2	1	0	1
r_3	1	1	0

Faiblement localisante

	f_1	f_2	f_3
r_1	1	1	0
r_2	1	0	1
r_3	0	1	1

Fortement localisante

Figure 1-5 : Table d'incidence

Une table non localisante ne permet pas de distinguer certains défauts entre eux. Une table faiblement localisante permet de localiser les défauts uniques sous hypothèse d'exonération. Une table fortement localisante garantit que les différentes sensibilités des résidus par rapport aux défauts ne conduisent pas à un diagnostic erroné. L'utilisation des tables de signatures établies pour des défauts simples pose des problèmes pour les situations de défauts multiples. Il faudrait en effet analyser toutes les combinaisons possibles des colonnes de la table.

1.5.3.2 Exemples

Les méthodes à base de modèles parallèles

Un modèle du procédé est alimenté en parallèle et reconstruit les sorties qui sont comparées aux mesures. Il est possible d'utiliser plusieurs modèles. Pour décrire des procédés non linéaires, certains auteurs proposent d'utiliser des réseaux de neurones comme modèles « boîtes noires » [Benquilou *et al.* 1999].

Les observateurs d'état

Les observateurs d'état sont des algorithmes, fondés sur un modèle du procédé, chargés de poursuivre l'état de celui-ci. On essaye, en général, de rendre l'observateur indépendant des perturbations non mesurées (entrées inconnues) et dépendant de certains défauts [Patton et Chen 1997].

L'identification

Les paramètres nominaux du modèle sont déterminés par des méthodes d'identification de paramètres pendant la période d'apprentissage hors ligne. Les paramètres estimés en ligne sont ensuite comparés aux paramètres nominaux pour la

détection de défaut. Cette méthode est surtout utilisée lorsque la structure du modèle est bien connue et que ses paramètres ont un sens physique. S'ils n'ont pas de sens physique, la modification des paramètres liée à une défaillance est difficile à interpréter. [Isermann 1993] propose l'application de cette méthode à un moteur à courant continu entraînant une machine outil. L'inconvénient de cette méthode est qu'il faut disposer de signaux suffisamment excités pour que l'algorithme d'identification converge.

Les équations de redondance analytique

La redondance d'information et de données est une source de modèles de bon ou de mauvais fonctionnement. Des **relations de redondance analytique** RRA sont extraites du modèle du système physique. Une RRA est une relation C (dédiuite par la combinaison de relations du modèle) qui lie entre elles des variables nécessairement mesurées, ici représentées par un ensemble Y et telle que $C(Y) = 0$. Une RRA peut être vérifiée pour diagnostiquer les composants physiques qui constituent son support. En effectuant des combinaisons judicieuses des relations du modèle, des RRA de supports différents (sensibles à des défauts particuliers) peuvent être générées.

Ces techniques sont très utilisées.

- [Ould Bouamama 2002] génèrent les relations de redondances analytiques (RRA) via les bond graphs (ou graphes de liaisons). Une application à un générateur de vapeur est proposée.
- [Krysanter et Nyberg 2002] ont développé une méthode de génération de RRA et l'ont appliquée sur un procédé complexe et non linéaire de fabrication de papier.
- [Evsukoff 1998] extrait des RRA spécifiques d'un modèle causal et l'applique à un atelier de colonnes pulsées. Le Chapitre 3 détaillera cette technique que nous avons suivie pour générer les alarmes.

1.5.4 Approches basées sur des modèles qualitatifs ou semi-qualitatifs

1.5.4.1 Principe

La connaissance à disposition sur les systèmes physiques est souvent partielle et imprécise et il n'est pas toujours possible d'en élaborer un modèle par simulation «classique». Dans ce cadre, deux approches sont distinguées, les approches purement qualitatives et les approches semi-qualitatives.

1.5.4.2 Exemple

Le raisonnement qualitatif

Le but du raisonnement qualitatif [Travé-Massuyès, Dague et Guerrin 1998] est de formaliser la connaissance intuitive du monde, de développer des méthodes de raisonnement qui utilisent cette connaissance dans des cas pratiques et de développer des modèles du mode de pensée humain à propos des systèmes physiques [Forbus 1984]. Cette connaissance peut se représenter par une tendance (généralement caractérisée par le signe des dérivées des mesures ou des résidus) ou un ordre de grandeur (faible, moyen, grand). Les modèles qualitatifs peuvent être utilisés lorsque aucun modèle mathématique suffisamment précis n'est disponible. Un exemple de connaissance qualitative est « quand le débit augmente, la température doit diminuer ». Cette connaissance est exploitée par des techniques informatiques relevant de l'intelligence artificielle [De Kleer 1987], [Travé-Massuyès *et al.* 2001]. Les modèles qualitatifs sont utilisés car ils présentent plusieurs intérêts :

- ils restent fidèles au mode de diagnostic des experts et leurs résultats sont donc facilement exploitables,
- dans beaucoup de cas, l'imprécision des connaissances rend l'obtention des modèles numériques difficile,
- dans un certain nombre de cas, associés à plus ou moins d'informations quantitatives, ils sont suffisants pour effectuer un diagnostic,
- les observations (mesures et observations humaines) sont souvent imprécises.

Parmi les méthodes de diagnostic qualitatif, nous citons GDE [De Kleer et Williams 1987], qui est utilisé pour les systèmes statiques et qui couple un propagateur de contraintes avec un ATMS (Assumption Thruth Maintenance System) pour déterminer des conflits. D'autres outils comme QSIM [Kuipers 1994] ont été développés pour effectuer la modélisation et la simulation qualitative de systèmes continus. Le simulateur calcule les évolutions possibles du système, mais certains comportements factices subsistent. QSIM a été utilisé à la base de plusieurs systèmes de monitoring tels que MIMIC [Dvorak et Kuipers 1989].

Cependant, certains inconvénients demeurent car les modèles qualitatifs fournissent souvent plusieurs résultats possibles. Pour traiter ces ambiguïtés qui proviennent de l'abstraction, il faut :

- choisir le bon modèle qualitatif,
- ajouter des connaissances quantitatives dans les modèles qualitatifs.

Les méthodes purement qualitatives sont peu appliquées en pratique. Les méthodes qui combinent des modèles qualitatifs et des modèles quantitatifs sont préférées. Nous avons choisi de suivre une approche semi-qualitative/semi-quantitative et nous verrons qu'elle apporte des informations complémentaires pour le diagnostic.

Les bases de règles

Ces systèmes ont donné naissance aux systèmes experts. Généralement, un modèle qualitatif de mauvais fonctionnement associe, via une base de règles, les symptômes aux défaillances selon, par exemple, le principe suivant :

Si [(Symptôme 1 et Symptôme 2) ou (Symptôme 3)] alors Défaillance (1-1)

[Pessi et Luparia 1992] propose de superviser une unité de distillation atmosphérique par un système à base de règles. Des techniques de traitement du signal permettent de générer les symptômes. L'inconvénient de ces méthodes est qu'il faut disposer d'une base riche de connaissance des défaillances du procédé. De plus, il est nécessaire de caractériser chaque défaillance de manière unique par des symptômes pertinents.

Les graphes causaux ou graphes d'influences

Les méthodes proposées dans [Montmain et Gentil 1993] et [Evsukoff, Montmain et Gentil 1997] qui utilisent un graphe causal (ou graphe d'influences) contenant des fonctions de transferts numériques sont une approche semi-qualitative/semi-quantitative. Par exemple les travaux [Travé-Massuyès *et al.* 2001] s'appuient aussi sur un modèle causal quantitatif intitulé Ca-En¹⁻⁴. L'approche développée dans ce manuscrit s'appuie en partie sur les travaux de [Evsukoff, Montmain et Gentil 1997] qui seront, par conséquent, largement discutés dans le Chapitre 3.

1.5.5 Conclusion

Les méthodes de diagnostic fondées sur des modèles, numériques ou qualitatifs, sont encore très peu utilisées dans le monde industriel. Les méthodes à base de réseaux de neurones sont en croissance ainsi que les combinaisons de différentes méthodes. En ce moment des projets Européens suivent cette tendance. Nous citons par exemple, le projet CHEM (cf. section 4.3.5.2) et le projet MAGIC¹⁻⁵ dont l'objectif est de développer un système multi-agents collaborant pour le diagnostic de procédés [Ploix, Gentil 2003], [voir la session MAGIC du workshop Safeprocess 2003 à Washington]. La tendance actuelle est à la mise en œuvre de méthodes combinant d'avantage d'expertise et de connaissance qualitative avec de la connaissance quantitative. C'est l'approche que nous proposons d'utiliser.

¹⁻⁴ (Causal Engine)

¹⁻⁵ Multi-Agents-Based Diagnostic Data Acquisition and Management in Complex Systems, Contrat N° : EU-IST-2000-30009, voir aussi : <http://magic.uni-duisburg.de>.

Chapitre 2

Modélisation causale

Ce chapitre propose une méthode de modélisation causale. L'objectif d'un modèle causal est de représenter les relations de causes à effets entre les variables. Chaque influence du modèle causal est associée à un ensemble de composants du procédé qui sous tendent cette influence.

La méthode proposée dans ce document est constituée de plusieurs étapes.

- Un modèle causal dit structurel est initialement élaboré à partir d'une analyse structurelle des relations du système.
- L'opération de réduction qui consiste à extraire les influences entre les variables **connues**²⁻¹ est ensuite appliquée sur le modèle causal structurel. Le modèle causal réduit ainsi obtenu contient des influences ne faisant intervenir que des variables connues.
- Finalement l'opération d'approximation est effectuée sur le modèle causal réduit. Elle permet de négliger certains phénomènes tout en garantissant une modélisation correcte. Le modèle causal approché est ainsi obtenu.

Le modèle causal approché contient la connaissance quantitative et qualitative qui est interprétée par le module de diagnostic présenté dans le Chapitre 3.

²⁻¹ Une variable connue peut être une mesure, une action ou une consigne transmise par le SCADA.

2.1 Introduction

Lorsqu'une défaillance se produit sur l'installation dont il a la charge, l'opérateur doit agir le plus souvent dans l'urgence et prendre une décision. Il doit localiser de manière plus ou moins précise la source des défauts qu'il observe. Ce raisonnement est souvent effectué en interprétant les alarmes classiques obtenues par dépassement de seuils et par l'observation des synoptiques. L'opérateur vérifie systématiquement si les mesures des variables régulées suivent bien leurs consignes et si elles sont cohérentes avec les actions appliquées. Il identifie quelles sont les variables connues qui ont un comportement normal ou anormal vis-à-vis soit de leur environnement local soit des consignes et des perturbations extérieures mesurées du procédé. La connaissance qu'il a acquise au cours de son expérience, conjuguée à des raisonnements physiques de causes à effets, lui permet de localiser la source de la défaillance qui explique le comportement anormal des autres variables.


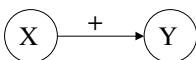
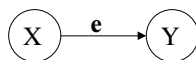
Les premiers systèmes de diagnostic étaient basés sur des systèmes experts. L'expérience a montré que ces techniques présentent des limites. Il est en particulier difficile d'introduire la notion de temps et d'assurer la cohérence d'une base de connaissance suffisamment importante [Travé-Massuyès et Gentil 1999]. Les méthodes à base de modèles sont de plus en plus retenues pour le diagnostic [Frank et Ding 2000]. Pour un procédé complexe, l'obtention d'un tel modèle est particulièrement difficile. De plus, ce modèle étant développé pour le diagnostic, des méthodes spécifiques de modélisation doivent être proposées. Ce chapitre propose des solutions à ce problème.

Le diagnostic causal consiste à déterminer quelle est la variable source qui explique les déviations sur toutes les autres variables [Montmain et Gentil 2000]. L'objectif de ce chapitre est de proposer une méthode pour l'obtention d'une représentation qualitative des relations normales de causes à effets. Cette représentation, le modèle causal, est utilisée de la manière suivante pour le diagnostic : si une variable X influence une variable Y et si la variable Y se comporte anormalement alors la raison est soit une perturbation sur Y , soit un changement de la relation entre X et Y , soit un comportement anormal de X . Ce raisonnement peut

être effectué de manière récursive afin de remonter des conséquences à la variable commune, source de tous les défauts observés.

Un modèle causal est qualitatif. Il contient l'identité des variables et les influences orientées entre les variables de la cause vers l'effet. Cette connaissance peut se représenter sous différentes formes (équations, matrices, graphe, etc.). Dans ce document, le modèle causal est représenté par un graphe causal. Par exemple, l'influence d'une variable X sur une variable Y est qualitativement représentée par le graphe causal de la Figure 2-1. Chaque modèle causal est donc représenté par un unique graphe causal.

Un modèle causal peut contenir des informations supplémentaires. Par exemple, chaque influence peut être labellisée par la relation «e» qui la supporte [Cordier *et al* 2000a] (cf. Figure 2-3). Les relations peuvent être signées (“+” signifie qu’une augmentation de X provoque une augmentation de Y et “-” qu’une augmentation de X provoque une diminution de Y). Cet exemple est illustré par la Figure 2-2. Une méthode de diagnostic basée sur des SDG (Signed Directed Graphs) est présentée dans [Iri *et al.* 1979]. Pour chaque arc, le test de la consistance entre le nœud cause et le nœud conséquence constitue la procédure basique de détection. Le sous graphe qui explique l'ensemble des défauts observés est identifié. L'algorithme proposé dans [Kramer et Palowitch 1987] s'appuie sur des graphes signés mais une information numérique supplémentaire est portée par les nœuds : la déviation normalisée de la mesure par rapport à la valeur du régime stationnaire. Les procédés présentent généralement un comportement dynamique. Par conséquent, la table des signatures évolue au cours du temps, ce qui nécessite une interprétation temporelle. [Montmain et Gentil 1993] proposent d'utiliser des fonctions de transfert continues et [Evsukoff 1998] propose d'utiliser des équations aux différences «e» pour prendre en compte la notion temporelle. [Mostreman, Biswas et Narasimham 1997] proposent de prédire le comportement futur de la variable pour chaque déviation anormale de leur dérivée qualitative. La plupart des procédés change de configuration, le modèle causal doit donc pouvoir s'adapter à ces changements. Certains sous graphes peuvent s'activer ou se désactiver en fonction de la configuration du procédé [Dziopa 1996], [Travé-Massuyès et Pons 1997]. Un des avantages des modèles causaux est d'exhiber l'ordre de résolution des équations. Le calcul des valeurs prises par les causes précède celui des effets.

		
<i>Figure 2-1. Influence simple</i>	<i>Figure 2-2. Influence signée</i>	<i>Figure 2-3. Influence & relation</i>

De nombreux travaux ont été effectués pour établir la causalité sous jacente à un phénomène physique. Nous distinguons les approches suivantes que nous présentons dans la suite de cette section :

- l'approche physique [Leyval, Gentil et Feray-Beaumont 1994],
- l'approche experte qui est heuristique [Heim, Cauvin et Gentil 2000a],
- la causalité mythique [De Kleer et Brown 1990],
- l'ordonnancement causal selon Iwasaki et Simon [Iwasaki et Simon 1986],
- les graphes de liaison ou bond graphs [Lucas 1994].

La première approche est physique : elle s'appuie sur la connaissance des ingénieurs responsables du procédé [Leyval, Gentil et Feray-Beaumont 1994]. Il s'agit alors de développer la causalité au sens physique du terme : quelle variable a une variation qui précède celle d'une autre, et donc l'influence causalement. Cette approche exige la disponibilité d'excellents ingénieurs, souvent ceux qui ont contribué à la conception du procédé.

Nous avons proposé une méthode experte permettant de représenter la causalité au sein des boucles complexes de régulation [Heim, Cauvin et Gentil 2001]. Nous nous sommes ensuite orientés vers la méthode décrite dans ce manuscrit car la méthode experte présentait les inconvénients suivants :

- la connaissance experte est difficile à formaliser et à organiser car les experts oublient souvent des variables intermédiaires lors de l'interprétation des phénomènes de causes à effets,
- la connaissance experte est subjective et dépend de la vision que l'expert a du procédé,
- la connaissance experte ne garantit pas une modélisation complète (oubli de phénomènes).

La causalité mythique, qui ne sera pas présentée dans ce document, consiste à postuler que le comportement du système résulte d'interactions ordonnées

temporellement entre composants voisins. Elle détermine le chemin suivi par une perturbation depuis une entrée du système [De Kleer et Brown 1990].

De leur côté, Iwasaki et Simon proposent une approche différente : elle se base sur une analyse purement structurelle d'un système d'équations. Elle ne nécessite cependant pas la résolution de ces équations. Iwasaki et Simon posent le problème de l'ordonnancement causal mais ils ne proposent pas d'algorithme pour l'obtenir. Nous verrons qu'il est nécessaire de faire appel à une méthode opérationnelle pour extraire cette causalité. Par exemple, elle peut être obtenue par une approche issue de la théorie des graphes [Porté *et al* 1988].

Les bonds graphs représentent les systèmes physiques sous la forme d'un intermédiaire entre la représentation physique et le modèle mathématique [Dauphin-Tanguy 2000]. Ils s'appuient sur la représentation des transferts de puissance en termes de variables d'effort et de flux (la puissance est le produit d'un effort par un flux). Les bonds graphs constituent un langage unifié à vocation pluridisciplinaire de représentation des systèmes. Les bonds graphs se représentent par un modèle graphique exhibant les relations de causes à effets. Les modèles mathématiques sous tendus s'extraient automatiquement. La causalité globale des systèmes statiques n'est toutefois pas facile à extraire [Ahriz 1998].

Un problème récurrent de l'ordonnancement causal doit être mentionné : la plupart des systèmes physiques inclue des boucles dites «de rétroaction». Les chercheurs ne s'accordent pas sur la signification de la causalité dans ce cas particulier. Iwasaki et Simon pensent que la notion de causalité n'a aucun sens dans les boucles de rétroaction. Ils considèrent que les variables d'une telle boucle sont mutuellement dépendantes. La méthode de diagnostic proposée dans le Chapitre 3 nécessite d'orienter les boucles de rétroaction selon une interprétation causale possible. Aussi nous adoptons l'approche [Travé-Massuyès et Pons 1997] qui extrait tous les ordonnancements causaux possibles dans les boucles de rétroaction [Porté *et al* 1988]. Le choix d'un ordonnancement causal particulier est à la charge de l'expert du procédé. Les experts d'un procédé ont souvent une interprétation causale préférée. Pour le diagnostic, contrairement à la simulation, on pourrait envisager d'utiliser toutes les relations issues de toutes les interprétations causales possibles, car elles ont généralement des supports différents [Flaus et Gentil 2003].

Contrairement aux travaux précédemment cités, ce chapitre ne s'attache pas à décrire uniquement la méthode d'ordonnancement causal, mais il propose une méthode plus générale qui décrit les étapes successives depuis la description du système physique jusqu'à l'obtention du modèle causal. Cette méthode ayant été appliquée au pilote de FCC, des problèmes concrets ont été soulevés et des solutions sont proposées.

Les principales étapes sont les suivantes :

- La section 2.2 présente la première étape : la génération du modèle structurel des relations (MSR). Le MSR est un système d'équations particulier obtenu à partir des principes premiers. Les relations sont élémentaires, répondent au principe de localité²⁻², leurs supports sont connus et des conditions d'activation sont appliquées en fonction de la configuration du système physique. Le MSR contient des variables mesurées et des variables non mesurées.
- Une méthode d'ordonnancement causal est ensuite appliquée au MSR. Le modèle causal structurel MCS est obtenu. La mise en œuvre suit une méthode issue de la théorie des graphes [Ford et Fulkerson 1956] et [Porté *et al* 1988]. La section 2.3 présente la technique d'ordonnancement causal.

Les deux étapes suivantes sont les opérations de réduction puis d'approximation.

- L'opération de réduction, appliquée au MCS, consiste à extraire les influences entre les variables connues puis à éliminer les variables non connues. Le modèle causal réduit (MCR) est ainsi obtenu.
- L'approximation (appliquée au MCR) génère le modèle causal approché MCA. Cette opération consiste à négliger les influences qui peuvent l'être : par exemple des influences qui sont négligeables vis à vis de l'objectif de diagnostic ou les influences de variables constantes dans le temps. Elle permet de conserver la connaissance qui garantira qu'un diagnostic, effectué sur la base du MCA, est correct. Le MCA contient donc des relations parfois approchées entre les variables connues.

²⁻² La description d'un composant ne doit jamais faire référence à un autre composant du système physique.

La section 2.4 présente les opérations puis les algorithmes de réduction et d'approximation. La méthodologie proposée est représentée sur la Figure 2-4.

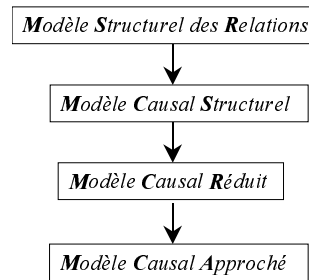


Figure 2-4: Méthodologie pour l'obtention d'un modèle causal approché

- Le MCA est un modèle dynamique de bon fonctionnement. Pour son utilisation en tant que simulateur, il est nécessaire d'effectuer préalablement sa quantification. Dans la section 2.5, une représentation linéarisée autour du point de fonctionnement est choisie et des méthodes d'identification des paramètres du modèle sont proposées.

L'ensemble de la méthodologie est illustrée dans ce chapitre sur un exemple académique simple. Le Chapitre 4 présentera le résultat de l'application de cette méthode au pilote de FCC.

2.2 Modèle structurel des relations

2.2.1 Introduction

Ce paragraphe présente la méthode d'élaboration du modèle structurel des relations (MSR). Le MSR est un système d'équations particulier obtenu en 7 étapes qui déterminent successivement :

- le système physique Σ étudié,
- les composants qui constituent Σ ,
- les configurations de Σ ,
- les variables pertinentes du modèle à vocation de diagnostic,
- la structure des relations,
- le support des relations,
- les conditions d'activation des relations en fonction des configurations de Σ .

Une grande majorité des auteurs ne détaille pas ces étapes ou seulement partiellement. Nous avons jugé nécessaire de les détailler afin de pouvoir appréhender la construction d'un modèle causal d'un procédé quelconque de manière systématique.

Remarque 2-1 : Il est important de remarquer qu'il n'est pas nécessaire de connaître les relations quantitativement mais que la structure des relations est suffisante pour obtenir le MSR²⁻³.

2.2.2 Détermination du système physique Σ

Il est indispensable de définir le système physique Σ sur lequel la modélisation causale va être appliquée. Σ peut représenter l'installation complète mais aussi un sous système de celle-ci. Ce choix est effectué en fonction des objectifs du modèle (par exemple le diagnostic). La Figure 2-5 illustre Σ au sein de son environnement. Il est important de remarquer à cette étape que la méthode de diagnostic présentée dans ce document s'applique à un sous système suffisamment connu de l'installation. D'autres méthodes de diagnostic, en particulier celles ne nécessitant pas de modèle, peuvent être appliquées sur d'autres sous-systèmes de l'installation. La méthode développée dans ce chapitre s'adresse prioritairement à un procédé continu (par opposition à un procédé batch).

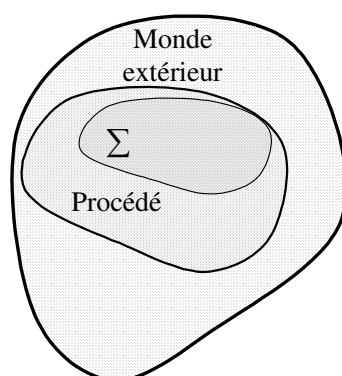


Figure 2-5. Position du système physique par rapport au monde extérieur

²⁻³ La structure d'une relation informe uniquement sur l'identité des variables en jeu dans cette relation.

2.2.3 Détermination des composants qui constituent Σ

Σ est divisé en composants. La granularité choisie pour représenter Σ (c'est-à-dire le nombre de composants) dépend des objectifs du modèle. Si un composant industriel contient des variables connues pertinentes pour le diagnostic de Σ alors il est judicieux de le diviser en d'autres composants. Par exemple, une vanne peut se représenter comme un tout ou il peut être intéressant de détailler les mécanismes de positionnement de celle-ci. D'un autre côté, si une vanne est remplacée quelle que soit sa défaillance, il n'est pas nécessaire de détailler son fonctionnement interne. Finalement, si le diagnostic est orienté pour la maintenance alors le bon niveau de granularité est celui du remplacement des composants. Par la suite, l'ensemble des composants de Σ est noté *COMPS* et l'ensemble des variables connues est *OBS*.

2.2.4 Identification de la configuration de Σ

La troisième étape consiste à identifier les configurations possibles de Σ . Une configuration est définie par les modes de fonctionnement des différents composants qui constituent Σ . En effet, certains composants peuvent avoir plusieurs modes de fonctionnement. Par exemple, une vanne tout ou rien présente deux modes de fonctionnement. Un autre exemple est une boucle de régulation qui est ouverte ou fermée. Il apparaît donc que l'ensemble des variables utilisées pour décrire Σ dépend du mode de fonctionnement de chaque composant. Les étapes suivantes (jusqu'à 2.2.7) se réfèrent à un unique mode de fonctionnement. La section 2.3.6 propose une méthode permettant de traiter les systèmes ayant plusieurs configurations.

2.2.5 Identification des variables décrivant Σ

Soit V l'ensemble de cardinal n_V des variables identifiées pour décrire Σ . Ces variables suffisent à décrire les phénomènes représentatifs de Σ . Le modèle doit être complet mais pas nécessairement détaillé : si un phénomène peut être décrit de manière simple, il n'est pas nécessaire de complexifier sa description.

Le monde extérieur a des influences sur Σ , les parties de l'installation industrielle qui ne sont pas incluses dans Σ échangent de la matière ou de l'énergie avec Σ . Ces influences peuvent aussi être dues à des perturbations extérieures ou encore à des actionneurs. Si ces influences sont pertinentes pour le module de diagnostic, il est

nécessaire de prendre en compte les variables à la source de ces influences. L'ensemble de ces variables constitue l'ensemble V_{exo} des variables exogènes à Σ .

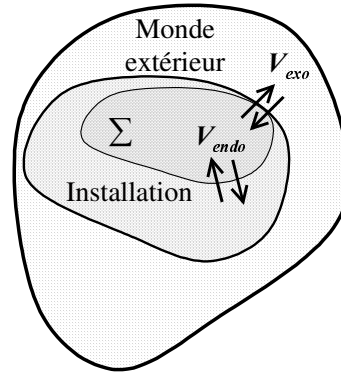


Figure 2-6. Variables décrivant Σ

Une variable de Σ est exogène si sa valeur est indépendante de celle des autres variables de Σ . Une variable qui n'est pas exogène appartient à l'ensemble V_{endo} des variables endogènes.

$$V = V_{endo} \cup V_{exo} \quad (2-1)$$

$$V_{exo} = \{V \in V / V \text{ est exogène pour } \Sigma\} \quad (2-2)$$

$$V_{endo} = \{V \in V / V \text{ est endogène pour } \Sigma\} \quad (2-3)$$

Dans la suite de ce document, une distinction sera faite entre les variables mesurées et les variables non mesurées. Nous définissons donc l'ensemble **CAPTEURS** qui est l'ensemble des capteurs du système physique. Le capteur d'une variable V est noté V^s . Nous définissons aussi l'application *Capteur* dont l'ensemble de départ est l'ensemble V et l'ensemble d'arrivée est l'ensemble **CAPTEURS**. Cette application permet de savoir si une variables est mesurée :

$$\begin{aligned} \forall V \in V, V \text{ mesurée} &\Rightarrow \text{Capteur}(V) = V^s \\ \forall V \in V, V \text{ non mesurée} &\Rightarrow \text{Capteur}(V) = \emptyset \end{aligned} \quad (2-4)$$

Remarque 2-2 : Une attention particulière doit être portée à la position²⁻⁴ des capteurs et des variables en général. En effet, cette connaissance qui est souvent implicite est primordiale pour le diagnostic car elle conditionne le support des relations. Elle doit être explicitement précisée.

²⁻⁴ La position désigne le lieu physique où la mesure est effectuée.

2.2.6 Expression formelle des relations

La cinquième étape consiste à définir les relations qui lient entre elles les variables de l'ensemble V . L'expression formelle des relations est suffisante à ce stade de la méthodologie. Les relations sont :

- des principes premiers (bilans de matière et d'énergie),
- des relations empiriques (obtenues via la connaissance de l'expert ou la littérature),
- des relations de connections,
- des relations réalisées par les calculateurs (régulations).

Ces relations permettent de décrire Σ dans une configuration donnée. Par conséquent, des conditions d'activation dépendant de la configuration de Σ sont associées aux relations. En d'autres termes, l'ensemble des relations décrivant Σ dépend de sa configuration, tout comme dans le système Ca-En [Travé-Massuyès *et al.* 2001].

La méthode d'ordonnement causal qui sera appliquée dans la suite de ce document nécessite de disposer de relations structurelles au sens de Iwasaki et Simon. Les relations doivent représenter des phénomènes basiques répondant au principe de localité (cf. note de page 2-2).

Remarque 2-3 : Il est important de remarquer que les relations du MSR sont élémentaires et répondent au principe de localité.

Soit l'ensemble \mathbf{E} , de cardinal n_E , des relations. Dans la suite de ce document les applications sont notées avec le formalisme suivant : en indice figurent l'ensemble de départ et l'ensemble d'arrivée (par exemple $\mathbf{E} \rightarrow \mathbf{V}$). Nous définissons l'application $Var_{\mathbf{V} \rightarrow \mathbf{E}, \text{intervenant}}$ dont l'ensemble de départ est l'ensemble \mathbf{E} des relations et l'ensemble d'arrivée est l'ensemble \mathbf{V} des variables. Pour chaque équation e , l'ensemble $Var_{\mathbf{E} \rightarrow \mathbf{V}, \text{intervenant}}(e)$ désigne l'ensemble des variables intervenant dans e :

$$\forall e \in \mathbf{E}, Var_{\mathbf{E} \rightarrow \mathbf{V}, \text{intervenant}}(e) = \{V \in \mathbf{V} / V \text{ intervient dans } e\} \quad (2-5)$$

2.2.7 Composants, supports des relations

Chaque relation doit être associée à un ensemble de composants et de capteurs de Σ dont le bon comportement sous tend la relation. Ils constituent le support de la relation.

Un composant (ou un capteur) est associé à une relation si et seulement si cette relation détermine son (ou une partie de son) comportement. Pour déterminer le support d'une relation, nous définissons deux applications : *Comps* et *Support_{Relation}*.

L'ensemble de départ de l'application *Comps* est l'ensemble E des relations et l'ensemble d'arrivée est l'ensemble **COMPS** des composants. Cette application détermine les composants (autres que capteurs) associés à une relation. L'ensemble de départ de l'application *Support_{Relation}* est l'ensemble des relations et l'ensemble d'arrivée est **COMPS** \cup **CAPTEURS** :

$$\forall e \in E \text{ Comps}(e) = \{C \in \mathbf{COMPS} / C \text{ associé à } e\} \quad (2-6)$$

$$\forall e \in E \text{ Support}_{\text{Relation}}(e) = \text{Comps}(e) \cup \text{Capteurs}(Var_{E \rightarrow V} \text{ intervenant}(e)) \quad (2-7)$$

2.2.8 Conditions d'activation et configuration de Σ

Les étapes précédentes se réfèrent à une configuration précise et doivent être itérées pour chaque configuration. [Travé-Massuyès et Pons 1997] proposent une méthode systématique permettant de gérer simultanément tous les modes de fonctionnement.

2.2.9 Exemple de description d'un système

L'exemple Σ que nous proposons sur la Figure 2-7 illustre la méthode présentée dans cette section. Σ , qui n'a qu'une seule configuration, est constitué de conduites $P_{i=1..5}$ et de deux jonctions J_1 et J_2 . Les débits sont identifiés par un nom générique $F_{(I \text{ or } O), P_{i=1..5}}$ "I" désigne l'entrée, "O" désigne la sortie dans la conduite P_i . $F_{J_{i=1,2}}$ est le débit dans la jonction J_i . Les débits $F_{I,P_1}, F_{I,P_2}, F_{O,P_3}, F_{I,P_4}$ et F_{O,P_5} sont respectivement mesurés par les capteurs $F_{I,P_1}^S, F_{I,P_2}^S, F_{O,P_3}^S, F_{I,P_4}^S$ et F_{O,P_5}^S .

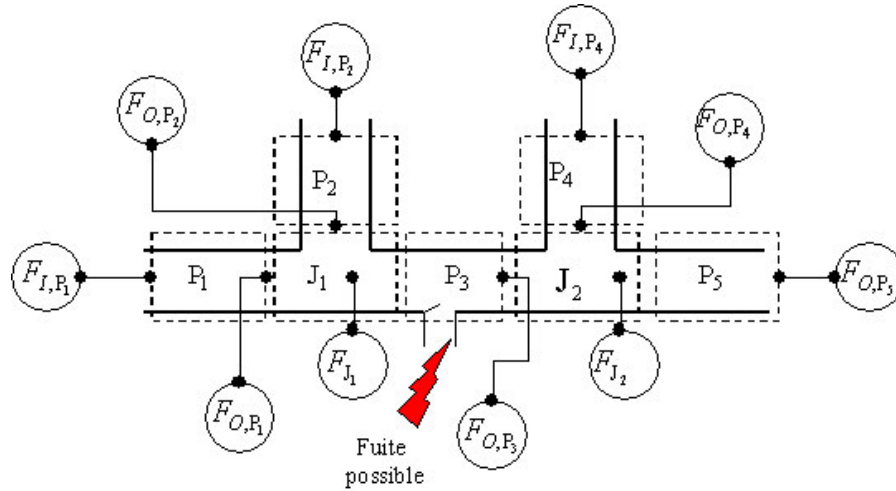


Figure 2-7: Exemple d'un réseau hydraulique

Les variables sont précisément positionnées. Certaines variables (i.e. F_{I,P_3}) ont été volontairement omises par souci de simplification. Cela n'enlève rien à la généralité de cet exemple. Les équations formelles et élémentaires (étant données les variables choisies) e_1, \dots, e_7 , sont utilisées pour décrire Σ .

Equation	Composant	Support	
$e_1 : F_{O,P_1} = e_1(F_{I,P_1})$	Conduite 1 (P_1)	$\{P_1, F_{I,P_1}^S\}$	(2-8)
$e_2 : F_{O,P_2} = e_2(F_{I,P_2})$	Conduite 2 (P_2)	$\{P_2, F_{I,P_2}^S\}$	(2-9)
$e_3 : F_{J_1} = e_3(F_{O,P_1}, F_{O,P_2})$	Jonction 1 (J_1)	$\{J_1\}$	(2-10)
$e_4 : F_{O,P_3} = e_4(F_{J_1})$	Conduite 3 (P_3)	$\{J_1, P_3, F_{O,P_3}^S\}$	(2-11)
$e_5 : F_{O,P_4} = e_5(F_{I,P_4})$	Conduite 4 (P_4)	$\{P_4, F_{I,P_4}^S\}$	(2-12)
$e_6 : F_{J_2} = e_6(F_{O,P_3}, F_{O,P_4})$	Jonction 2 (J_2)	$\{J_2, F_{O,P_3}^S\}$	(2-13)
$e_7 : F_{O,P_5} = e_7(F_{J_2})$	Conduite 5 (P_5)	$\{J_2, P_5, F_{O,P_5}^S\}$	(2-14)

Comme F_{I,P_1} influence F_{O,P_3} qui influence F_{O,P_5} , alors une perturbation sur F_{I,P_1} se propage sur F_{O,P_5} via F_{O,P_3} . Une fuite dans la conduite 3 (Figure 2-7) modifie F_{O,P_3} et par conséquent modifie F_{O,P_5} . Nous verrons dans la suite de ce document et tout particulièrement dans la section concernant la localisation des défauts que la notion de support d'une relation est cruciale pour le diagnostic. Nous montrons dans la section suivante, que selon les relations choisies pour décrire le système, les supports varient beaucoup.

2.2.10 Exemples de relations de redondance analytique

Le modèle précédemment décrit est le plus élémentaire²⁻⁵ qui puisse être établi considérant les composants et les variables. Dans la suite de ce paragraphe, des relations de redondance analytique sont construites entre les variables connues en combinant les équations $\{e_1, \dots, e_7\}$.

Une relation de redondance analytique RRA (cf. 1.5.4) dont le support est $\{P_1, P_2, P_3, J_1, F_{I,P_1}^S, F_{I,P_2}^S, F_{O,P_3}^S\}$, est déduite de $\{e_1, e_3, e_4\}$:

$$F_{O,P_3} - e_4(e_3(e_1(F_{I,P_1}), e_2(F_{I,P_2}))) = 0 \quad (2-15)$$

Une autre RRA de support $\{P_4, P_5, J_2, F_{O,P_3}^S, F_{I,P_4}^S, F_{O,P_5}^S\}$, est déduite de $\{e_5, e_6, e_7\}$:

$$F_{O,P_5} - e_7(e_6(F_{O,P_3}, e_5(F_{I,P_4}))) = 0 \quad (2-16)$$

L'intersection des supports des deux RRA (2-15) et (2-16) est F_{O,P_3}^S . Une troisième RRA de support $\{P_1, P_2, P_3, P_4, P_5, J_1, J_2, F_{I,P_1}^S, F_{I,P_2}^S, F_{I,P_4}^S, F_{O,P_5}^S\}$ est finalement déduite de $\{e_1, \dots, e_7\}$. Cette RRA correspond à une combinaison des relations (2-15) et (2-16) :

$$F_{O,P_5} - e_7(e_6(e_4(e_3(e_1(F_{I,P_1}), e_2(F_{I,P_2}))), e_5(F_{I,P_4}))) = 0 \quad (2-17)$$

Nous remarquons que le capteur F_{O,P_3}^S n'est pas inclus dans le support de (2-17). Cet exemple simple illustre qu'il est facile de construire des RRA. Il est aussi facile de construire une RRA par combinaison de deux autres RRA. Cependant nous observons que le support de cette nouvelle relation n'est pas égal à l'intersection des supports des relations qui ont été combinées. Nous verrons dans le chapitre 3.6 que le support de la RRA est un paramètre très important pour le diagnostic. Sa taille et les composants de celui-ci vont permettre d'obtenir un diagnostic plus ou moins précis. Cette exemple simple nous a montré qu'en effectuant des combinaisons de RRA, nous avons changé les supports des relations donc l'information que nous pourrions obtenir lors du diagnostic.

²⁻⁵ Les relations sont élémentaires : elles minimisent le nombre de composants dans le support de la relation.

Cet exemple simple illustre la difficulté de construction de RRA pertinentes pour le diagnostic. Plusieurs approches peuvent être choisies pour générer ces RRA à partir d'un système d'équations. [Krysander et Nyberg 2002] proposent une méthode de génération d'un grand nombre de RRA. Ensuite, les plus intéressantes, définies par leur «capacité de diagnostic», sont conservées. Dans le chapitre 3.3 nous proposons une méthode de génération de RRA.

Comme les relations sont élémentaires, alors le cardinal des variables intervenant dans chaque relation est faible, ce qui donne des supports petits pour le diagnostic.

2.2.11 Conclusion

Ce paragraphe rappelle les points importants des étapes (cf. Figure 2-8) d'obtention du MSR :

- Détermination du système physique Σ étudié : le choix de Σ dépend de la connaissance à disposition sur le procédé. Le MSR peut décrire seulement une partie du procédé.
- Détermination des composants qui constituent Σ : Σ est divisé en sous systèmes qui constituent ses composants. Si le diagnostic est orienté pour la maintenance alors le bon niveau de granularité est celui du remplacement de ces composants.
- Détermination des configurations de Σ : un modèle causal doit être obtenu pour chaque configuration de Σ .
- Identification des variables pertinentes du modèle à vocation de diagnostic : ces variables suffisent à décrire les phénomènes représentatifs de Σ selon un niveau de description guidé par le niveau de remplacement des composants.
- Expression de la structure des relations : des relations élémentaires, répondant au principe de localité doivent être établies.
- Détermination du support des relations : les composants qui sous tendent les relations sont associés à celles-ci.
- Détermination des conditions d'activation des relations de Σ .

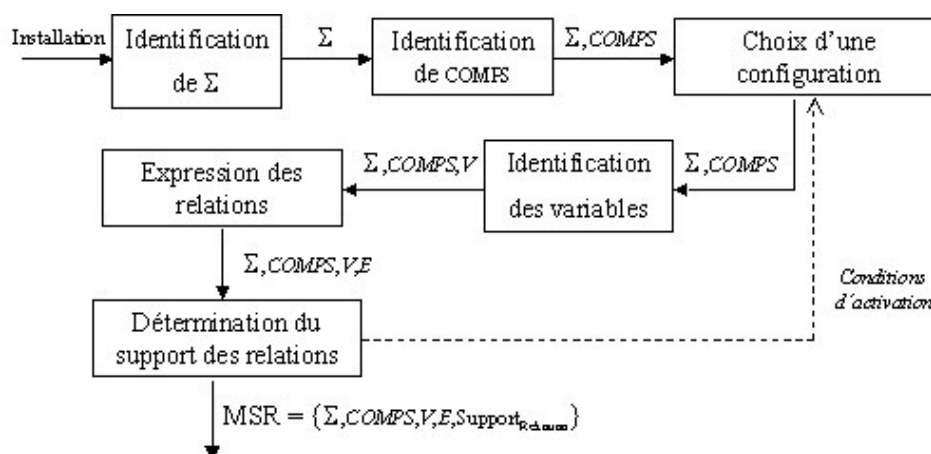


Figure 2-8: Méthode d'obtention du modèle structurel des relations

Nous avons appliqué cette méthode au pilote de FCC et la section 4.3.2 décrit les résultats obtenus.

Une fois les relations élémentaires écrites, choisies pertinentes pour la granularité des composants et décrivant une configuration spécifiée, alors l'ordonnancement causal du MSR est recherché. La section suivante présente la méthode d'obtention du modèle causal structurel par ordonnancement causal des relations du MSR.

2.3 Modèle causal structurel

2.3.1 Introduction

Cette section présente une méthode d'ordonnancement causal. Cette méthode est équivalente à la détermination de l'ordre de calcul des variables d'un système d'équations, depuis les variables exogènes, de proche en proche, vers chaque variable endogène souhaitée.

Plusieurs solutions, précédemment présentées, sont proposées pour définir ou extraire l'ordonnancement causal. Dans cette section, nous décrivons la méthode que nous avons décidé de suivre et qui est présentée dans les travaux de [Travé-Massuyès et Pons 1997]. Nous n'avons pas apporté d'amélioration à cette méthode mais nous proposons de la présenter dans un environnement plus général et industriel :

- nous avons d'abord décrit le MSR sur lequel s'applique cette méthode dans la section 2.2,
- nous l'illustrons sur un exemple et proposons des solutions pratiques pour les choix à effectuer pour son application dans des situations concrètes

Cette section présente :

- la manière dont Iwasaki et Simon appréhendent la notion de causalité,
- une méthode d'obtention de la causalité via la théorie des graphes,
- le bilan sur la connaissance contenue dans le modèle causal structurel.

2.3.2 Obtention de la causalité selon Iwasaki et Simon

Iwasaki et Simon déduisent l'ordonnancement causal d'une analyse purement structurelle des équations. Les équations sont structurelles : elle interprètent des mécanismes qui répondent au principe de localité par opposition à des équations qui ne représenteraient pas les mécanismes séparément. Le système d'équations doit contenir autant d'équations que de variables endogènes. Nous reviendrons sur ce point dans la section 2.3.3.1.

Cette section présente l'ordonnancement causal selon Iwasaki et Simon dans :

- les systèmes dynamiques,
- les systèmes statiques,
- les systèmes mixtes,
- un exemple de système statique.

2.3.2.1 La causalité dans les systèmes dynamiques

Les systèmes dynamiques peuvent se représenter par des équations différentielles. Dans ce manuscrit, les équations différentielles sont supposées être sous la forme canonique²⁻⁶. Cette hypothèse n'est pas restrictive car il est possible de transformer toute équation différentielle d'ordre «n» en «n» équations différentielles d'ordre 1.

Iwasaki et Simon distinguent deux types de causalité : la causalité intégrale et la causalité différentielle. La causalité intégrale s'adresse à une variable et à sa dérivée. Elle exprime que la valeur d'une variable $x_i(t)$ à l'instant t dépend de sa dérivée $\frac{dx_i}{dt}$ et de sa valeur à l'instant précédent $x_i(t-\Delta t)$. Cette notion sous-tend que la variable est distinguée de sa dérivée selon la Figure 2-9 :

$$\forall t, x_i(t) = x_i(t - \Delta t) + \frac{dx_i}{dt} \cdot \Delta t \quad (2-18)$$

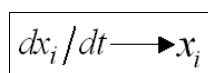


Figure 2-9 : Causalité intégrale

La causalité différentielle s'adresse à un groupe de variables $x_1, x_2, \dots, x_i, \dots, x_n$ intervenant dans une équation différentielle canonique telle que :

$$\frac{dx_i}{dt} = f(x_1, x_2, \dots, x_i, \dots, x_n) \quad (2-19)$$

L'interprétation causale naturelle exprime la dépendance de la dérivée par rapport aux autres variables intervenant dans l'équation selon la Figure 2-10.

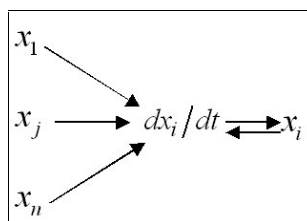


Figure 2-10 : Causalité différentielle et causalité intégrale

²⁻⁶ Une équation différentielle est dite sous forme canonique s'il n'y a qu'une seule dérivée dans l'équation et que cette dérivée est le seul terme apparaissant dans le membre gauche de l'équation.

Définition 2-1 : Un système dynamique de «n» équations différentielles du premier ordre à «n» variables est dit *complet* si tout sous-ensemble de «k» équations ($k \leq n$) contient au moins les dérivées de k variables.

L'ordonnancement causal d'un système complet sous la forme canonique est donné par les assertions suivantes :

- chaque variable dépend de sa dérivée selon la causalité intégrale,
- chaque dérivée dépend des autres variables intervenant dans l'équation canonique selon la causalité différentielle.

Les systèmes d'équations sont aussi composés de relations statiques. Ce paragraphe permet de déterminer la causalité dans les systèmes uniquement dynamiques. Le paragraphe suivant traite des systèmes uniquement statiques.

2.3.2.2 La causalité dans les systèmes statiques

Le système d'équations statiques ne doit pas être dégénéré, il doit contenir autant d'équations que de variables endogènes. La description de l'algorithme d'ordonnancement causal nécessite la définition préliminaire de deux notions :

Définition 2-2 : Un système statique (qualitatif) de «n» équations et «n» variables (endogènes) est dit *complet* si tout sous système de «k» équations contient au moins $k \leq n$ variables (endogènes).

Définition 2-3 : Etant donné un système complet \mathcal{S} , un sous système s de \mathcal{S} qui est aussi complet et qui ne contient pas de sous système complet est appelé sous système complet minimal (SCM).

L'ordre causal est défini selon l'algorithme récursif suivant :

Etant donné un système statique complet d'équations \mathcal{S} , soit \mathcal{S}_0 l'union de ses SCM, appelé d'ordre zéro. Comme \mathcal{S}_0 est complet, alors les variables intervenant dans \mathcal{S}_0 peuvent être déterminées en résolvant les équations de \mathcal{S}_0 . En substituant toutes les occurrences de ces variables dans $\mathcal{S} \setminus \mathcal{S}_0$ un nouveau système *dérivé d'ordre 1*, \mathcal{S}'_1 , est obtenu. Soit \mathcal{S}_1 , l'union des SCM de \mathcal{S}'_1 dits d'ordre 1. Cet algorithme est réitéré

jusqu'à ce que le dernier système complet dérivé ne contienne plus de sous système complet.

Nous définissons l'application $Var_{E \rightarrow V, intervenant, max, MCS}$ dont l'espace de départ est \mathbf{E} et l'espace d'arrivée \mathbf{V} . L'image $Var_{E \rightarrow V, intervenant, max, MCS}(e)$, de la relation e , est le sous ensemble de $Var_{E \rightarrow V, intervenant}(e)$ qui contient les variables qui appartiennent aux SCM d'ordre le plus élevé. Ces variables sont causalement dépendantes (selon Iwasaki et Simon) des variables du sous ensemble $Var_{E \rightarrow V, intervenant}(e) / Var_{E \rightarrow V, intervenant, max, MCS}(e)$:

$$\forall e \in \mathbf{E}, Var_{E \rightarrow V, intervenant, max, MCS}(e) = \left\{ V \in Var_{E \rightarrow V, impliquée}(e) \middle/ \begin{array}{l} V \text{ appartient au SCM} \\ \text{d'ordre le plus élevé} \end{array} \right\} \quad (2-20)$$

La détermination de l'ordre causal consiste donc à générer les SCM par ordre croissant. L'ordre causal finalement obtenu est unique.

Remarque 2-4 : Un SCM contenant plus d'une variable caractérise une boucle statique. Selon Iwasaki et Simon, toutes les variables d'un tel SCM sont mutuellement dépendantes et appartiennent au même niveau causal. Le problème des boucles de rétroaction est abordé différemment selon les auteurs.

Exemple :

Cet exemple illustre la recherche de l'ordonnancement causal du système \mathbf{S} défini par (2-21). Les variables V_1 , V_2 et V_5 qui interviennent seules dans les relations e_5 , e_6 et e_7 sont des variables exogènes. Cet exemple sera repris dans la suite de ce document. Le cas particulier des variables exogènes sera plus largement discuté.

$$\begin{aligned} Var_{E \rightarrow V, intervenant}(e_1) &= \{V_1, V_2, V_3, V_6\} \\ Var_{E \rightarrow V, intervenant}(e_2) &= \{V_3, V_4, V_5\} \\ Var_{E \rightarrow V, intervenant}(e_3) &= \{V_4, V_6\} \\ Var_{E \rightarrow V, intervenant}(e_4) &= \{V_4, V_6, V_7\} \\ Var_{E \rightarrow V, intervenant}(e_5) &= \{V_1\} \\ Var_{E \rightarrow V, intervenant}(e_6) &= \{V_2\} \\ Var_{E \rightarrow V, impliquées}(e_7) &= \{V_5\} \end{aligned} \quad (2-21)$$

Les SCM consécutifs de (2-21) sont présentés dans la suite de ce paragraphe. Ils permettent d'établir l'ordre causal.

\mathcal{S}_0 , l'union des SCM d'ordre 0 de (2-21) est $\{e_5, e_6, e_7\}$. Par conséquent, les variables V_1 , V_2 et V_5 sont substituées par leurs valeurs (notées v_1^* , v_2^* et v_5^*) pour toutes leurs occurrences dans $\{e_1, e_2, e_3, e_4\}$:

$$\begin{aligned} e_1 &: e_1(v_1^*, v_2^*, V_3, V_6) \\ e_2 &: e_2(V_3, V_4, v_5^*) \\ e_3 &: e_3(V_4, V_6) \\ e_4 &: e_4(V_4, V_6, V_7) \end{aligned} \tag{2-22}$$

Afin de déterminer \mathcal{S}_1 qui est l'union des SCM d'ordre 1, les sous systèmes complets de (2-22) doivent être identifiés. La Table 2 présente tous les sous systèmes extraits de (2-22).

$\mathcal{S}_{1,1}$ est un sous système composé de l'unique équation (e_1) et des 2 variables (V_3, V_6),

$\mathcal{S}_{1,2}$ est un sous système composé des 2 équations (e_1, e_2) et des 3 variables (V_3, V_4, V_6),

$\mathcal{S}_{1,3}$ est composé des 3 équations (e_1, e_2, e_3) et des 3 variables (V_3, V_4, V_6),

...

<u>Systèmes à 1 relation</u>	<u>Systèmes de 2 relations</u>	<u>Systèmes de 3 relations</u>
(4 possibilités)	(6 possibilités)	(4 possibilités)
$\mathcal{S}_{1,1} = (e_1)(V_3, V_6)$	$\mathcal{S}_{1,2} = (e_1, e_2)(V_3, V_4, V_6)$	$\mathcal{S}_{1,3} = (e_1, e_2, e_3)(V_3, V_4, V_6)$
$\mathcal{S}_{2,1} = (e_2)(V_3, V_4)$	$\mathcal{S}_{2,2} = (e_1, e_3)(V_3, V_4, V_6)$	$\mathcal{S}_{2,3} = (e_1, e_3, e_4)(V_3, V_4, V_6, V_7)$
$\mathcal{S}_{3,1} = (e_3)(V_4, V_6)$	$\mathcal{S}_{3,2} = (e_1, e_4)(V_3, V_4, V_6, V_7)$	$\mathcal{S}_{3,3} = (e_1, e_2, e_4)(V_3, V_4, V_6, V_7)$
$\mathcal{S}_{4,1} = (e_4)(V_4, V_6, V_7)$	$\mathcal{S}_{4,2} = (e_2, e_3)(V_3, V_4, V_6)$	$\mathcal{S}_{4,3} = (e_2, e_3, e_4)(V_3, V_4, V_6, V_7)$
	$\mathcal{S}_{5,2} = (e_3, e_4)(V_4, V_6, V_7)$	
	$\mathcal{S}_{6,2} = (e_2, e_4)(V_3, V_4, V_6, V_7)$	

Tableau 2-1 : Sous systèmes de (2-22)

L'unique SCM du Tableau 2-1 est $\mathcal{S}_1 = (e_1, e_2, e_3)(V_3, V_4, V_6)$. \mathcal{S}_1 est donc $\{e_1, e_2, e_3\}$. Comme \mathcal{S}_1 comporte trois variables V_3 , V_4 et V_6 , il constitue une boucle. Ces 3 variables sont substituées par leurs valeurs pour leurs occurrences dans $\{e_1, e_2, e_3\}$. Finalement, \mathcal{S}_2 est $\{e_4\}$.

A cette étape, l'ordonnancement causal du système (2-21) selon Iwasaki et Simon est exactement déterminé :

- $Var_{E \rightarrow V, \text{intervenant.max.MCS}}(e_1) = \{V_3, V_6\}$ informe que V_3 et V_6 dépendent de V_1 et V_2 ,
- $Var_{E \rightarrow V, \text{intervenant.max.MCS}}(e_2) = \{V_3, V_4\}$ informe que V_3 et V_4 dépendent de V_5 ,
- $Var_{E \rightarrow V, \text{intervenant.max.MCS}}(e_4) = \{V_7\}$ informe que V_7 dépend de V_4 et V_6 .

La relation e_3 apporte peu d'information en terme de causalité. En effet $Var_{E \rightarrow V, \text{intervenant}, \text{max}, \text{MCS}}(e_3) = \{V_4, V_6\}$ et $Var_{E \rightarrow V, \text{intervenant}}(e_3) = \{V_4, V_6\}$. Par conséquent V_4 et V_6 sont simplement mutuellement dépendantes. Les relations e_5 , e_6 et e_7 n'apportent pas d'information en terme de causalité.

La Figure 2-11 représente l'ordonnancement causal de l'exemple (2-21) selon Iwasaki et Simon. Dans cette représentation, les relations élémentaires de la boucle de \mathcal{S}_1 ne sont pas représentées.

Par exemple, la relation e_3 qui lie les variables V_4 et V_6 n'est pas exhibée.

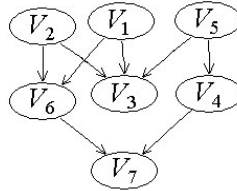


Figure 2-11 : Ordre causal de l'exemple (2-21) selon Iwasaki et Simon

2.3.2.3 La causalité dans les systèmes mixtes

Les structures mixtes sont composées d'équations différentielles et d'équations statiques. Ce sont les systèmes habituellement rencontrés.

L'ordre causal d'un système \mathcal{S} mixte est obtenu en construisant un nouveau système $\text{Inst}(\mathcal{S})$. $\text{Inst}(\mathcal{S})$ est un système qui inclut \mathcal{S} et tel que, pour chaque variable dérivée $\frac{dx_i}{dt}$ apparaissant dans \mathcal{S} , l'équation $x_i = x_i^*$ est artificiellement ajoutée à \mathcal{S} où x_i^* est une valeur constante.

L'ordonnancement causal des systèmes mixtes est obtenu avec les deux étapes suivantes :

- Etape n°1 : application de l'algorithme d'ordonnancement causal des systèmes statiques à $\text{Inst}(\mathcal{S})$.
- Etape n°2 : ajout des influences de la causalité intégrale de chaque dérivée sur sa primitive.

Cette méthode est illustrée par un exemple dans la suite de cette section.

Définition 2-4 : Un système S de n équations à n variables est *complet* si et seulement si :

- S est composé d'équations différentielles de premier ordre et d'équations statiques,
- $Inst(S)$ est un système statique complet (les dérivées des variables étant considérées comme des nouvelles variables endogènes distinctes).

Remarque 2-5 : Un système mixte complet admet un unique ordonnancement causal au sens d'Iwasaki et Simon.

Un système mixte complet dont le sous système statique ne contient pas de boucle de rétroaction admet un ordre causal, selon Iwasaki et Simon, directement utilisable par la méthode de diagnostic du Chapitre 3.

Si le sous système statique contient une boucle de rétroaction, il est nécessaire d'orienter cette boucle. La théorie des graphes présentée dans la section 2.3.3 propose une méthode d'obtention de toutes les solutions.

Exemple : Nous illustrons la méthode sur le système S suivant :

$$\begin{aligned} \dot{x} &= f(x, s) \\ \dot{s} &= g(x, s) \\ h &= h(x, s) \end{aligned} \tag{2-23}$$

Les trois étapes de la méthode sont illustrées sur la Figure 2-12 :

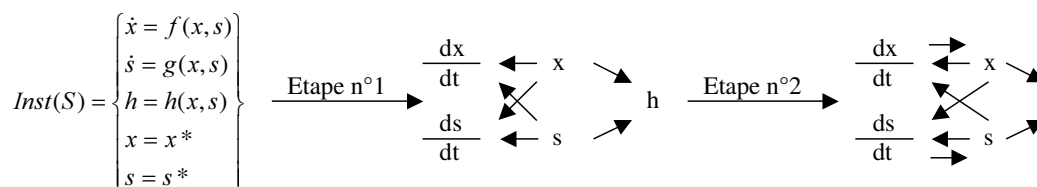


Figure 2-12 : La causalité des systèmes mixtes

2.3.3 Obtention de la causalité via la théorie des graphes

L'approche de Iwasaki et Simon est purement théorique. Elle a permis de définir formellement la notion de causalité. La méthode de diagnostic présentée dans le Chapitre 3 nécessite de déterminer un ordonnancement causal complet. Ce paragraphe présente un algorithme issu de la théorie des graphes qui permet d'obtenir l'ordonnancement causal d'une manière qui puisse être implémentée informatiquement.

[Porté *et al.* 1988] ont montré comment mettre en œuvre l'ordonnancement causal des systèmes statiques via la théorie des graphes.

Etant donné un système complet $\mathbf{S}=(\mathbf{E}, \mathbf{V})$ constitué d'un ensemble \mathbf{E} de « n » équations et d'un ensemble \mathbf{V} de « n » variables, alors le problème de l'ordonnancement causal se ramène à la recherche d'un couplage parfait [Gondran et Minoux 1979] dans le graphe biparti²⁻⁷ défini entre l'ensemble \mathbf{E} et l'ensemble \mathbf{V} .

Ce couplage parfait peut être obtenu, par exemple, par l'algorithme de Ford-Fulkerson [Ford et Fulkerson 1956]. Si le système \mathbf{S} est complet alors le graphe biparti admet un couplage parfait. Ce paragraphe présente les étapes nécessaires à l'élaboration de ce couplage parfait.

Cinq étapes ont été identifiées pour générer le MCS à partir des informations précédemment obtenues :

- Génération d'un graphe biparti préliminaire (qui peut correspondre à un système d'équations dégénéré),
- Génération du graphe biparti initial,
- Réalisation d'un couplage parfait sur le graphe biparti initial,
- Génération d'un graphe orienté déduit du couplage parfait,
- Génération du MCS.

²⁻⁷ Un graphe biparti est un graphe non orienté dans lequel les nœuds peuvent être divisés en deux ensembles. Aucun arc ne lie deux nœuds du même ensemble.

2.3.3.1 Graphe biparti préliminaire

La première étape consiste à générer un *graphe biparti* préliminaire. Le graphe biparti $\mathbf{G} = (\mathbf{V} \cup \mathbf{E}, \mathbf{ARCS})$ est défini par l'ensemble \mathbf{ARCS} des arcs non orientés entre les relations et les variables.

Nous définissons l'application *Exist* dont l'ensemble de départ est l'ensemble \mathbf{ARCS} et l'ensemble d'arrivée est l'ensemble {Vrai, Faux}. L'application *Exist* caractérise l'existence des arcs du MCS :

$$\begin{aligned} \text{Exist}(\text{arc}) = \text{vrai} &\Rightarrow \text{arc existe} \\ \text{Exist}(\text{arc}) = \text{faux} &\Rightarrow \text{arc n'existe pas} \end{aligned} \tag{2-24}$$

Soit l'application A de $(\mathbf{V}^* \mathbf{E})$ dans {Vrai, Faux}. La valeur de $A(V, e)$ est Vrai si et seulement si il existe un arc entre V et e (si et seulement $V \in \text{Var}_{E \rightarrow V, \text{intervenant}}(e)$) :

$$\begin{aligned} \forall e \in \mathbf{E}, \forall V \in \mathbf{V}, V \in \text{Var}_{E \rightarrow V, \text{intervenant}}(e) &\Rightarrow \text{Exist}(A(V, e)) = \text{vrai} \\ \forall e \in \mathbf{E}, \forall V \in \mathbf{V}, V \notin \text{Var}_{E \rightarrow V, \text{intervenant}}(e) &\Rightarrow \text{Exist}(A(V, e)) = \text{faux} \end{aligned} \tag{2-25}$$

Exemple : L'exemple illustré dans cette section diffère de l'exemple (2-21) par la définition des variables endogènes et exogènes. Considérons dans un premier temps que l'on ne dispose que de quatre relations pour décrire Σ . Il vient :

$$\mathbf{E} = \left\{ \begin{array}{l} e_1 : e_1(V_1, V_2, V_3, V_6) \\ e_2 : e_2(V_3, V_4, V_5) \\ e_3 : e_3(V_4, V_6) \\ e_4 : e_4(V_4, V_6, V_7) \end{array} \right\} \tag{2-26}$$

$$\mathbf{V} = \{V_1, V_2, V_3, V_4, V_5, V_6, V_7\} \tag{2-27}$$

$$\mathbf{V}_{\text{endo}} = \{V_3, V_4, V_5, V_6, V_7\} \tag{2-28}$$

$$\mathbf{V}_{\text{exo}} = \{V_1, V_2\} \tag{2-29}$$

Le graphe biparti correspondant est illustré par la Figure 2-13.

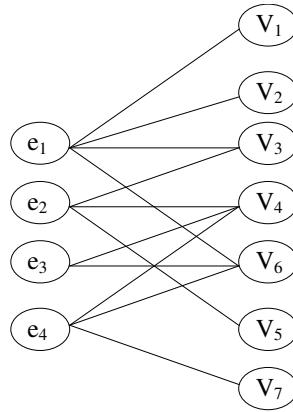


Figure 2-13: Graphe biparti préliminaire

La recherche de l'ordonnancement causal est équivalente à la détermination de l'ordre de calcul des variables d'un système d'équations. Elle est aussi équivalente à la construction d'une bijection, notée $Rel_{V \rightarrow E, \text{couplage}}$, de \mathbf{V} dans \mathbf{E} illustrée par la Figure 2-14. $Rel_{V \rightarrow E, \text{couplage}}(V)$ désigne la relation qui sera utilisée pour calculer la valeur de V . Soit l'application $(Rel_{V \rightarrow E, \text{couplage}})^{-1}$ de \mathbf{E} dans \mathbf{V} notée $Var_{E \rightarrow V, \text{couplage}}$. Cette application exprime quelle est la variable V qui est calculée à l'aide de l'équation e :

$$\forall e \in \mathbf{E}, Var_{E \rightarrow V, \text{couplage}}(e) = \{V \in \mathbf{V} / e = Rel_{V \rightarrow E, \text{couplage}}(V)\} \quad (2-30)$$

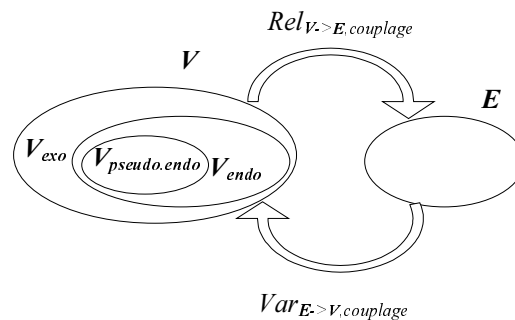


Figure 2-14: Illustration de la relation de couplage

Déterminer l'ordonnancement causal nécessite de disposer initialement d'un système d'équations complet. Cette contrainte nécessite, entre autres, de disposer d'autant de relations (n_E) que de variables (n_V). Si $n_E < n_V$ alors il est nécessaire de considérer que certaines variables endogènes sont pseudo-exogènes pour le MCS. Ces variables font partie de l'ensemble $V_{\text{pseudo.exo}}$. Dans ce document le cas $n_E > n_V$ n'est pas traité.

Pour chaque variable V pseudo-exogène, l'ensemble \mathbf{E} doit être artificiellement augmenté de la relation identité Id_V telle que :

$$\begin{aligned} \forall V \in \mathbf{V}_{exo}, Rel_{V \rightarrow E, couplage}(V) : Id_V \\ \forall V \in \mathbf{V}_{pseudo.exo}, Rel_{V \rightarrow E, couplage}(V) : Id_V \\ Id_V : V = v^* \quad (v^* \text{ est une valeur constante}) \end{aligned} \quad (2-31)$$

Les ensembles de variables endogènes $\mathbf{V}_{endo.causal}$ et exogènes $\mathbf{V}_{exo.causal}$ sont définis par :

$$\mathbf{V}_{endo.causal} = \mathbf{V}_{endo} \setminus \mathbf{V}_{pseudo.exo} \quad (2-32)$$

$$\mathbf{V}_{exo.causal} = \mathbf{V}_{exo} \cup \mathbf{V}_{pseudo.exo} \quad (2-33)$$

Exemple:

Dans l'exemple (2-26), $n_V = 7 \neq n_E = 4$. \mathbf{E} doit donc être augmenté de telle manière que $n_E = 7$ et que $\mathbf{S} = \mathbf{E} \cup \mathbf{V}$ soient complets. En fonction de ce choix, plusieurs systèmes peuvent être engendrés. Ils conduisent à des MCS différents. Ces choix n'ont pas d'influence sur la méthode qui va être présentée. Dans un premier temps, le couplage illustré par (2-34) est considéré : pour illustrer la méthode, V_5 est choisie arbitrairement pseudo-exogène. Les autres possibilités seront présentées dans le paragraphe 2.3.3.6 dans lequel des conseils pour le choix des variables pseudo-exogènes sont aussi proposés.

Nous remarquons que l'exemple (2-34) que nous avons construit dans cette section est identique à l'exemple (2-21). L'objectif de cette section est d'illustrer le rôle et l'importance des variables pseudo-exogènes.

$$\mathbf{E} = \left\{ \begin{array}{l} e_1 : e_1(V_1, V_2, V_3, V_6) \\ e_2 : e_2(V_3, V_4, V_5) \\ e_3 : e_3(V_4, V_6) \\ e_4 : e_4(V_4, V_6, V_7) \\ Id_{V_1} : V_1 = v_1^* \\ Id_{V_2} : V_2 = v_2^* \\ Id_{V_5} : V_5 = v_5^* \end{array} \right\} \quad (2-34)$$

$$\mathbf{V}_{exo.causal} = \{V_1, V_2, V_5\} \quad (2-35)$$

$$\mathbf{V}_{endo.causal} = \{V_3, V_4, V_6, V_7\} \quad (2-36)$$

Le graphe biparti de la Figure 2-15 est obtenu. L'ordonnancement causal va consister en premier lieu à établir un couplage parfait.

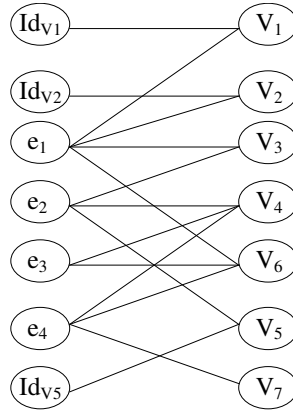


Figure 2-15: Graphe biparti \mathbf{G}

2.3.3.2 Couplage parfait du graphe biparti

Cette section présente les couplages parfaits possibles dans \mathbf{G} .

Définition 2-5: Le couplage parfait d'un graphe biparti est un ensemble d'arcs tel que chaque nœud est connecté à exactement un et unique arc.

Proposition 2-1 : Un graphe biparti associé à un système complet a un couplage parfait.

Soit l'application P de $(\mathbf{E}^* \mathbf{V})$ dans $\{\text{Vrai}, \text{Faux}\}$. Il existe un couplage parfait entre la relation e et la variable V si et seulement si $P(e, V) = \text{Vrai}$. La relation P définit le couplage parfait de \mathbf{G} :

$$\begin{aligned} \forall e \in \mathbf{E}, \forall V \in \mathbf{V}, \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{couplage}}(e) = V &\Leftrightarrow \text{Exist}(P(e, \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{couplage}}(e))) = \text{vrai} \\ \forall e \in \mathbf{E}, \forall V \in \mathbf{V}, \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{couplage}}(e) \neq V &\Leftrightarrow \text{Exist}(P(e, \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{couplage}}(e))) = \text{faux} \end{aligned} \quad (2-37)$$

Si le système $\mathbf{S} = \mathbf{E} \cup \mathbf{V}$ ne contient pas de boucle de rétroaction alors le couplage parfait est unique et correspond de manière univoque à l'ordonnancement selon Iwasaki et Simon [Porté *et al.* 1988] :

$$\forall e \in \mathbf{E}, \text{Card}(\text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{intervenant.max.MCS}}(e)) = 1 \Rightarrow \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{couplage}}(e) = \text{Var}_{\mathbf{E} \rightarrow \mathbf{V}, \text{intervenant.max.MCS}}(e) \quad (2-38)$$

Si le système $\mathbf{S} = \mathbf{E} \cup \mathbf{V}$ contient une boucle de rétroaction alors le couplage parfait n'est pas unique et chaque couplage correspond à une interprétation causale possible des boucles de rétroaction.

Exemple:

L'exemple (2-34) est identique à l'exemple (2-26) déjà traité dans le paragraphe 2.3.2.2.

La structure dérivée d'ordre zéro est $\mathcal{S}_0 = \{Id_{V_1}\}, \{Id_{V_2}\}, \{Id_{V_5}\}$, celle d'ordre 1 est $\mathcal{S}_1 = \{e_1, e_2, e_3\}$ et celle d'ordre 2 est $\mathcal{S}_2 = \{e_4\}$. A ce stade, seulement une partie des arcs du couplage parfait \mathbf{P} est connue. Ces arcs sont illustrés par la Figure 2-16.

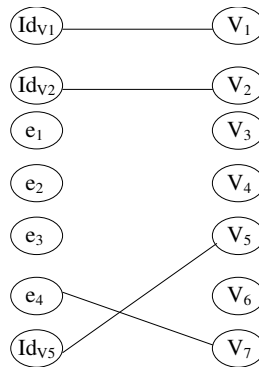


Figure 2-16: Graphe biparti, source du couplage parfait dans \mathbf{G} .

Lorsque l'on cherche à associer les variables restantes à une équation, on s'aperçoit qu'il y a plusieurs solutions :

- dans e_1 , V_3 et V_6 dépendent de V_1 et V_2 , donc e_1 est couplée à V_3 ou à V_6 ,
- dans e_2 , V_3 et V_4 dépendent de V_5 , par conséquent e_2 est couplée à V_3 ou V_4 ,
- concernant la relation e_3 , e_3 est couplée à V_4 ou à V_6 .

Ceci correspond exactement au sous système \mathcal{S}_1 que l'on avait trouvé par la méthode de Iwasaki et Simon.

Deux couplages parfaits sont donc envisageables :

- Couplage n°1 : e_1 est couplée à V_3 , e_2 est donc couplée à V_4 et e_3 est couplée à V_6 . Ce cas est illustré par la Figure 2-17.
- Couplage n°2 : e_1 est couplée à V_6 , e_2 est couplée à V_3 et e_3 est couplée V_4 .

Une autre recherche de couplage aboutit à l'échec : si e_1 est couplée à V_6 et e_2 est couplée à V_4 alors e_3 ne peut plus être couplée à aucune relation.

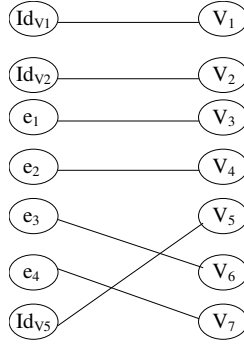


Figure 2-17: Le couplage parfait de \mathbf{G} , couplage n°1.

2.3.3.3 Graphes orientés

Le graphe orienté, dernière structure avant le MCS est dérivé du couplage parfait sur \mathbf{G} . Avant de présenter sa méthode d'obtention, quelques notions sont définies.

Les causes directes d'une variable V sont définies par l'application $Var_{V \rightarrow V, causes.directes}$ dont l'espace de départ est \mathbf{V} et l'espace d'arrivée est \mathbf{V} :

$$\forall V \in \mathbf{V}, Var_{V \rightarrow V, causes.directes}(V) = Var_{E \rightarrow V, intervenant}(Rel_{V \rightarrow E, couplage}(V) \setminus \{V\}) \quad (2-39)$$

Exemple:

$i = 1, 2, 5$

$$\begin{aligned} Var_{V \rightarrow V, causes.directes}(V_1) &= Var_{E \rightarrow V, intervenant}(Rel_{V \rightarrow E, couplage}(V_1) \setminus \{V_1\}) = Var_{E \rightarrow V, intervenant}(Id_{V_1}) \setminus \{V_1\} = \emptyset \\ Var_{V \rightarrow V, causes.directes}(V_3) &= Var_{E \rightarrow V, intervenant}(e_1) \setminus \{V_3\} = \{V_1, V_2, V_3, V_6\} \setminus \{V_3\} = \{V_1, V_2, V_6\} \\ Var_{V \rightarrow V, causes.directes}(V_4) &= Var_{E \rightarrow V, intervenant}(e_2) \setminus \{V_4\} = \{V_3, V_4, V_5\} \setminus \{V_4\} = \{V_3, V_5\} \\ Var_{V \rightarrow V, causes.directes}(V_6) &= Var_{E \rightarrow V, intervenant}(e_3) \setminus \{V_6\} = \{V_4, V_6\} \setminus \{V_6\} = \{V_4\} \\ Var_{V \rightarrow V, causes.directes}(V_7) &= Var_{E \rightarrow V, intervenant}(e_4) \setminus \{V_7\} = \{V_4, V_6, V_7\} \setminus \{V_7\} = \{V_4, V_6\} \end{aligned} \quad (2-40)$$

Soit le graphe orienté $\mathbf{G}' = (\mathbf{V} \cup \mathbf{E}, \mathbf{A}')$ où \mathbf{A}' est un ensemble d'arcs orientés tels que $a'(V, e)$ est orienté de \mathbf{V} vers \mathbf{E} et $a'(e, V)$ est orienté de \mathbf{E} vers \mathbf{V} :

$$\begin{aligned} \forall e \in \mathbf{E}, \forall V \in \mathbf{V}, A'(e, V) &= P(e, V) \\ \forall V \in \mathbf{V}, \forall e \in \mathbf{E}, \{V\} \in Var_{V \rightarrow V, causes.directes}(Var_{E \rightarrow V, couplage}(e)) &\Rightarrow Exist(A'(V, e)) = vrai \\ \forall V \in \mathbf{V}, \forall e \in \mathbf{E}, \{V\} \notin Var_{V \rightarrow V, causes.directes}(Var_{E \rightarrow V, couplage}(e)) &\Rightarrow Exist(A'(V, e)) = faux \end{aligned} \quad (2-41)$$

Le graphe orienté du couplage parfait de la Figure 2-17 est illustré par la Figure 2-18.

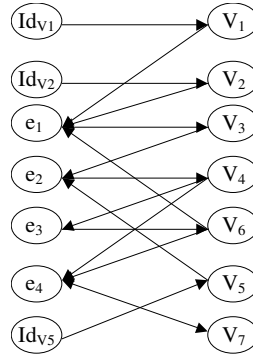


Figure 2-18: Graphe orienté G' du couplage n°1

2.3.3.4 Génération du modèle causal structurel

Le MCS est dérivé du graphe orienté G' . Tous les chemins dans G' liant les nœuds des variables via les nœuds des relations sont extraits. Les nœuds des relations sont ensuite éliminés. Le MCS représenté par le graphe causal $G_c = (E, I)$ contient des arcs orientés $I_{CMS}(V, Y)$ connectant les variables entre elles :

$$\begin{aligned} \forall V \in \mathcal{V}, \forall Y \in \mathcal{V}, \{Y\} \in \text{Var}_{\mathcal{V} \rightarrow \mathcal{V}, \text{causes.directes}}(V) &\Rightarrow \text{Exist}(I_{CMS}(V, Y)) = \text{vrai} \\ \forall V \in \mathcal{V}, \forall Y \in \mathcal{V}, \{Y\} \notin \text{Var}_{\mathcal{V} \rightarrow \mathcal{V}, \text{causes.directes}}(V) &\Rightarrow \text{Exist}(I_{CMS}(V, Y)) = \text{faux} \end{aligned} \tag{2-42}$$

Le graphe causal de l'exemple est présenté sur la Figure 2-19.

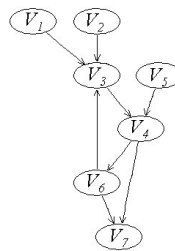


Figure 2-19: Graphe causal correspondant au couplage n°1

2.3.3.5 Nomenclature dans le graphe causal

Ce paragraphe présente quatre définitions auxquelles il sera fait référence dans la suite du document. Ces définitions sont les causes non directes, les causes, les causes connues et les conséquences d'une variable.

L'application *Connue* est aussi définie. Elle génère les variables connues d'un ensemble de variables (voir note de bas de page n°2-1).

L'application $Var_{V \rightarrow V, causes.backward}(V, n)$ définit les causes d'une variable V de profondeur n :

$$\begin{aligned} Var_{V \rightarrow V, causes.backward}(V, 1) &= Var_{V \rightarrow V, causes.directes}(V) \\ \forall V \in \mathbf{V}, \forall n \in \mathbf{N}^*, & \\ Var_{V \rightarrow V, causes.backward}(V, n) &= Var_{V \rightarrow V, causes.directes}(Var_{V \rightarrow V, causes.backward}(V, n-1)) \cup_{k=1 \dots n-1} Var_{V \rightarrow V, causes.backward}(V, k) \end{aligned} \quad (2-43)$$

Cette application est illustrée sur les exemples (2-49), (2-50), et (2-51).

L'ensemble de causes non directes de V est donné par l'application $Var_{V \rightarrow V, causes.-directes}$ de \mathbf{V} dans \mathbf{V} :

$$\forall V \in \mathbf{V}, Var_{V \rightarrow V, causes.-directes}(V) = \cup_{k=2 \dots \infty} Var_{V \rightarrow V, causes.backward}(V, k) \quad (2-44)$$

L'ensemble des causes de V est donné par l'application $Var_{V \rightarrow V, causes}$ de \mathbf{V} dans \mathbf{V} :

$$\forall V \in \mathbf{V}, Var_{V \rightarrow V, causes}(V) = Var_{V \rightarrow V, causes.-directes}(V) \cup Var_{V \rightarrow V, causes.directes}(V) \quad (2-45)$$

L'ensemble des conséquences ou enfants d'une variable V est l'ensemble des variables Y tel que V est une cause directe de Y . Il est donné par l'application $Var_{V \rightarrow V, conséquences.directes}$ de \mathbf{V} dans \mathbf{V} :

$$\forall V \in \mathbf{V}, Var_{V \rightarrow V, conséquences.directes}(V) = \{Y \in \mathbf{V} \setminus \{V\} / V \in Var_{V \rightarrow V, causes.directes}(Y)\} \quad (2-46)$$

Une variable du MCS est influencée par ses causes. Une variable exogène n'a pas de cause. Une influence entre une variable et sa cause directe est une influence directe.

$Connue(\mathbf{M})$ désigne le sous-ensemble des variables connues du sous-ensemble \mathbf{M} de \mathbf{V} :

$$\forall \mathbf{M} \subset \mathbf{V}, Connue(\mathbf{M}) = \{V \in \mathbf{M} \cap \text{OBS}\} \quad (2-47)$$

Il sera important de déterminer les causes connues d'une variable. Ces variables qui sont obtenues par l'application $Var_{E \rightarrow V, causes, connues}(e, n)$, seront aussi déterminées par l'algorithme de réduction (cf. Annexe A). L'ensemble des causes connues de la variable à une profondeur n est donné par l'application $Var_{V \rightarrow V, causes, connues}$ de $V^* \mathbb{N}$ dans \mathbf{N} :

$$\forall V \in \mathbf{V}, \forall n \in \mathbb{N} \quad Var_{V \rightarrow V, causes, connues}(V, n) = Connues(Var_{V \rightarrow V, causes, backward}(V, n)) \quad (2-48)$$

Exemple: Application de $Var_{E \rightarrow V, causes, backward}$ sur le couplage n°1 (Figure 2-19).

Les causes de profondeur 1 sont données par (2-43).

Les variables causes backward de profondeur 2 sont obtenues par :

$$\begin{aligned} & Var_{V \rightarrow V, causes, backward}(V_7, 2) \\ &= Var_{V \rightarrow V, causes, directes}(Var_{V \rightarrow V, causes, backward}(V_7, 1)) \setminus Var_{V \rightarrow V, causes, backward}(V_7, 1) \\ &= Var_{V \rightarrow V, causes, directes}(\{V_4, V_6\}) \setminus \{V_4, V_6\} \\ &= \{V_3, V_4, V_5\} \setminus \{V_4, V_6\} \\ &= \{V_3, V_5\} \end{aligned} \quad (2-49)$$

Les variables causes backward de profondeur 3 sont obtenues par :

$$\begin{aligned} & Var_{V \rightarrow V, causes, backward}(V_7, 3) \\ &= Var_{V \rightarrow V, causes, directes}(Var_{V \rightarrow V, causes, backward}(V_7, 2)) \setminus Var_{V \rightarrow V, causes, backward}(V_7, 1) \cup Var_{V \rightarrow V, causes, backward}(V_7, 2) \\ &= Var_{V \rightarrow V, causes, directes}(\{V_3, V_5\}) \setminus \{V_4, V_6, V_3, V_5\} \\ &= \{V_1, V_2\} \setminus \{V_4, V_6, V_3, V_5\} \\ &= \{V_1, V_2\} \end{aligned} \quad (2-50)$$

$$Var_{V \rightarrow V, causes, backward}(V_7, 4) = \emptyset \quad (2-51)$$

Les conséquences directes de V_4 sont :

$$\begin{aligned} & \{V_4\} \notin Var_{V \rightarrow V, causes, directes}(V_i) = \emptyset, i = 1, 2, 5 \\ & \{V_4\} \notin Var_{V \rightarrow V, causes, directes}(V_3) \\ & \{V_4\} \in Var_{V \rightarrow V, causes, directes}(V_6) \Rightarrow \{V_6\} \in Var_{V \rightarrow V, conséquences, directes}(V_4) \\ & \{V_4\} \in Var_{V \rightarrow V, causes, directes}(V_7) \Rightarrow \{V_7\} \in Var_{V \rightarrow V, conséquence, directes}(V_4) \end{aligned} \quad (2-52)$$

$$Var_{V \rightarrow V, \text{conséquences directes}}(V_4) = \{V_6, V_7\}$$

Appliquons $Var_{V \rightarrow V, \text{causes, connues}}$ au couplage n°1 (Figure 2-19).

$$\text{Supposons que : } Connues\{V_1, V_2, V_3, V_4, V_5, V_6, V_7\} = \{V_2, V_5, V_6\} \quad (2-53)$$

Il vient :

$$Var_{V \rightarrow V, \text{causes, connues}}(V_7, 1) = Connues(Var_{V \rightarrow V, \text{causes, backward}}(V_7, 1)) = Connues(\{V_4, V_6\}) = \{V_6\} \quad (2-54)$$

$$Var_{V \rightarrow V, \text{causes, connues}}(V_7, 2) = Connues(Var_{V \rightarrow V, \text{causes, backward}}(V_7, 2)) = Connues(\{V_3, V_5\}) = \{V_5\} \quad (2-55)$$

$$Var_{V \rightarrow V, \text{causes, connues}}(V_7, 3) = Connues(Var_{V \rightarrow V, \text{causes, backward}}(V_7, 3)) = Connues(\{V_1, V_2\}) = \{V_2\} \quad (2-56)$$

$$Var_{V \rightarrow V, \text{causes, connues}}(V_7, 4) = Connues(Var_{V \rightarrow V, \text{causes, backward}}(V_7, 4)) = \emptyset \quad (2-57)$$

$$Var_{V \rightarrow V, \text{causes, connues}}(V_7) = \{V_2, V_5, V_6\}$$

2.3.3.6 Autres modèles causaux engendrés

Nous remarquons que lors de l'application de la méthode d'ordonnement causal à (2-34), deux choix arbitraires doivent être faits par un expert du procédé. Le choix des variables pseudo exogènes et le choix de représentation de la boucle de rétroaction.

Exemple :

Toutes les possibilités d'ordre causal de (2-26) sont présentées dans la table suivante. Les couplages n°6 et n°7 contiennent une boucle entre V_4 et V_6 .

Couplage	Choix des variables pseudo-exogènes	Choix dans la boucle
1	V_5	V_3 est couplée à e_1
2	V_5	V_3 est couplée à e_2
3	V_3	pas de boucle
4	V_4	pas de boucle
5	V_6	pas de boucle
6	V_7	V_4 est couplée à e_4
7	V_7	V_4 est couplée à e_3

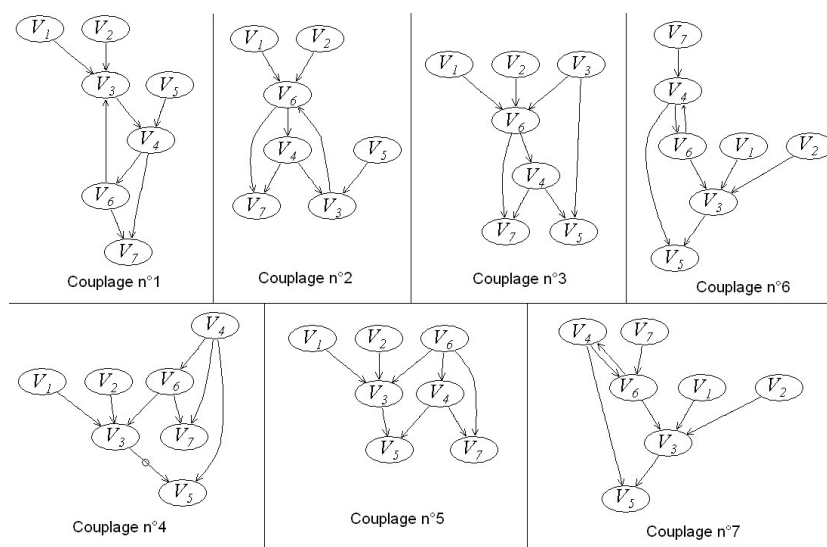


Figure 2-20: Graphes causaux possibles pour(2-26)

Nous constatons qu'il existe plusieurs interprétations causales d'un système d'équations. Nous donnons dans la section suivante quelques critères de choix.

2.3.4 Bilan sur la méthode d'obtention de la causalité

La méthode ne pouvant traiter que des systèmes de «n» équations à «n» inconnues, certaines variables endogènes doivent être considérées pseudo-exogènes. D'autre part, s'il y a présence d'une boucle de rétroaction, un choix de représentation de cette boucle doit être fait.

La présence de variables pseudo-exogènes est la conséquence d'un manque d'équations pour décrire le procédé. Nous proposons de choisir les variables pseudo-exogènes selon les critères suivants :

- les différentes dynamiques des variables. Il est préférable de choisir exogènes les variables qui ont une dynamique d'évolution lente.
- la confiance accordée aux capteurs. Il est préférable d'avoir confiance dans les mesures des variables pseudo-exogènes.
- la sensibilité de la relation. Les variables les moins sensibles de la relation devraient être considérées pseudo-exogènes.

Nous remarquons que le choix des variables pseudo-exogènes engendre la présence ou l'absence de boucles de rétroaction. Nous proposons que ce choix soit guidé par l'expert du procédé selon son interprétation causale préférée de ces boucles.

La méthode générale de l'obtention de l'ordonnement causal d'un système mixte complet d'équations est la suivante :

- nous avons élaboré dans la section 2.2 une méthode de génération du MSR,
- nous proposons des critères de choix des variables pseudo-exogènes,
- nous nous appuyons sur les travaux de [Porté *et al.* 1988] pour déterminer l'ordonnement causal dans la section 2.3.3,
- nous proposons des choix de représentation des boucles de rétroaction.

2.3.5 Information contenue dans le modèle causal structurel

Le MCS se représente par un graphe causal (\mathbf{G}_C) mais contient plus d'information que celui-ci. Cette section récapitule la connaissance contenue dans le MCS.

Dans les paragraphes précédents, le modèle était présenté sous la forme de relations pouvant lier entre elles plusieurs variables. Dans la suite de ce document, la notion d'influence est utilisée. L'ensemble des influences du MCS est \mathbf{I}_{MCS} . Une influence du MCS (représentée par un arc dans le \mathbf{G}_C) de la variable V sur la variable Y est notée $I_{MCS}(V, Y)$. Une influence désigne la contribution de la variable V au comportement de la variable Y . Elle a les attributs suivants :

- son existence $Exist(I_{MCS}(V, Y))$ (cf. (2-37)),
- son support $Support_{Influence}(I_{MCS}(V, Y))$,
- la contribution sous tendue ($Contribution(I_{MCS}(V, Y))$, (cette connaissance est facultative).

Les deux paragraphes suivants présentent le support et les contributions des influences.

2.3.5.1 Association des influences avec des composants

Chaque influence du MCS est associée à un ensemble de composants physiques. Un composant physique est associé à une influence si et seulement si elle détermine l'ensemble ou une partie de son comportement. Cette notion est équivalente à celle présentée dans le paragraphe 2.2.7. Elle lie des influences et non des relations aux composants physiques. Nous définissons l'application $Support_{Influence}$ de \mathbf{I}_{MCS} dans \mathbf{COMPS} qui détermine l'ensemble des composants associés à une influence :

$$\forall (V, Y) \in \mathcal{V}^2, \text{Support}_{\text{Influence}}(I_{\text{MCS}}(V, Y)) = \text{Comps}(Rel_{V \rightarrow E, \text{couplage}}(Y)) \cup \text{Capteurs}(\{V, Y\}) \quad (2-58)$$

2.3.5.2 Contributions quantitatives du modèle causal structurel

A ce stade de la méthode, deux approches sont envisageables pour quantifier le MCS. La quantification peut être effectuée avant ou après l'opération de réduction. Pour chaque variable V , la relation e couplée à V doit être transformée, par des manipulations algébriques des équations, de telle manière que V (à gauche) soit une fonction F_V de ses causes directes (à droite).

$$\forall V \in \mathcal{V}, V = F_V(\text{Var}_{\mathcal{V} \rightarrow \mathcal{V}, \text{causes directes}}(V)) \quad (2-59)$$

Si la fonction F_V est linéaire alors la contribution $\text{Contribution}(I_{\text{CMS}}(V, Y))$ est immédiatement déduite de F_V . Sinon elle doit être linéarisée autour du point de fonctionnement pour s'écrire sous la forme :

$$\forall V \in \mathcal{V}, V = \sum_{Y \in \text{Var}_{\mathcal{E} \rightarrow \mathcal{V}, \text{causes directes}}(V)} \text{Contribution}(I_{\text{CMS}}(Y, V))Y \quad (2-60)$$

Ces manipulations sont illustrées sur un exemple du pilote de FCC dans [Heim *et al.* 2003b].

2.3.5.3 Paramètres du modèle

A cette étape de la méthode, les paramètres physiques issus de la connaissance du procédé ou de la littérature peuvent être utilisés pour quantifier le modèle. Nous verrons dans la suite de ce document que les paramètres peuvent aussi être obtenus par une procédure d'identification des paramètres. Cette procédure ne sera cependant appliquée que après avoir effectué les opérations de réduction et d'approximation présentées dans les sections suivantes.

2.3.6 Systèmes à plusieurs modes de fonctionnement

La majorité des systèmes ont plusieurs modes opératoires, ce qui leur confère un caractère hybride. Ils sont le siège de processus continus, mais ils incluent par exemple des contrôleurs de nature discrète.

Il est nécessaire de reprendre la méthode d'ordonnement causal que nous avons développée pour chaque configuration du système physique. Lorsqu'il y a plusieurs configurations, il devient inextricable d'obtenir l'ordonnement causal, pour chaque configuration, de manière manuelle.

Une solution consiste à utiliser un modèle équationnel modal dans le sens où des conditions sont associées à ses équations. Ces conditions sont associées aux influences du modèle causal structurel.

[Travé-Massuyès et Pons 1997] proposent une extension de la méthode d'ordonnement causal pour de tels systèmes. Lors du changement de mode, la structure causale correspondant au nouveau mode est engendrée de manière incrémentale, c'est-à-dire en partant de celle dans le mode précédent plutôt que d'être calculée à nouveau entièrement.

Nous n'avons pas implémenté cette méthode dans le cadre de l'application sur le pilote de FCC et nous n'avons considéré qu'un seul mode de fonctionnement. Une difficulté déjà prévisible est qu'une telle méthode ne semble applicable que si le modèle et ses paramètres sont parfaitement connus dans toutes les configurations, ce qui n'entrait pas dans le cadre de la thèse.

2.3.7 Conclusion

Le MCS s'obtient à partir du MSR en appliquant un algorithme d'ordonnement causal (cf. Figure 2-9). L'approche d'Iwasaki et Simon formalise et caractérise l'ordonnement causal mais sans donner de méthode opérationnelle (algorithme) pour l'obtenir. [Travé-Massuyès et Pons 1997] en dérivent un algorithme qui consiste à chercher les SCM. [Porté *et al.* 1988] proposent une méthode opérationnelle pour obtenir l'ordonnement causal d'Iwasaki et Simon :

- la première étape consiste à trouver un couplage parfait pour pouvoir ensuite utiliser l'algorithme de Ford et Fulkerson (d'autres algorithmes peuvent être utilisés);
- à partir de ce couplage parfait, un graphe orienté est généré selon la méthode de la section 2.3.3.3;

- les composantes fortement connexes (qui correspondent aux boucles de rétroaction) sont ensuite recherchées dans ce graphe orienté. Les nœuds de ces composantes sont agrégés pour produire le graphe orienté résultant qui est identique à celui obtenu par l'ordonnancement causal d'Iwasaki et Simon.

Si le système n'a pas de boucle, le couplage parfait est unique et la dernière étape ne trouvera pas de composantes fortement connexes. Ainsi, l'ordonnancement causal s'obtient directement. Si le système a des boucles, alors plusieurs couplages parfaits existent et les graphes orientés sont différents (mais ne diffèrent que par les orientations à l'intérieur des boucles). Cependant, les composantes fortement connexes obtenues vont agréger exactement les mêmes variables. Si la méthode est arrêtée à la seconde étape, toutes les interprétations causales des boucles sont obtenues. L'expert en choisit une car il a souvent une interprétation causale préférée de ces boucles.

Cette méthode a été appliquée au pilote de FCC (cf. Chapitre 4). Il a été nécessaire de définir un certain nombre de variables pseudo-exogènes et donc des choix arbitraires ont été effectués. Ces choix ont été faits par un expert du procédé selon les critères proposés dans le paragraphe 2.3.3.6.

Deux autres opérations sont ensuite appliquées sur le MCS pour l'obtention d'un modèle causal utilisable pour le diagnostic d'un procédé industriel : la réduction et l'approximation. Ces étapes sont décrites dans la section 2.4.

2.4 Modèle causal approché

2.4.1 Introduction

Dans cette section, les opérations de réduction puis d'approximation sont appliquées sur le MCS. Seules les influences entre les variables connues sont pertinentes pour le diagnostic, mais le MCS contient des variables connues et inconnues. Par conséquent, l'objectif de l'opération de réduction est d'appréhender les influences entre les variables connues en éliminant les variables inconnues intermédiaires. Pour effectuer cette opération, nous nous appuyons sur les travaux de

[Travé-Massuyès *et al.* 2001] auxquels nous apportons une amélioration qui est décrite dans cette section.

L'algorithme de réduction produit le modèle causal réduit (MCR). Le MCR peut contenir :

- des relations difficiles à estimer,
- des relations qui sont négligeables vis à vis de l'objectif de diagnostic,
- des variables exogènes inconnues,
- des variables constantes dans le temps.

Prendre en compte ces relations ou les influences de ces variables peut, dans un premier temps, être fastidieux étant donné l'objectif. Nous avons élaboré l'opération d'approximation qui permet de négliger ces phénomènes tout en conservant la connaissance qui garantira des diagnostics corrects. Cette opération est appliquée au MCR et produit le modèle causal approché (MCA). Cette opération est aussi obligatoirement appliquée aux variables exogènes non mesurées. En effet dans ce cas, il est impossible de quantifier les influences de ces variables.

2.4.2 Réduction

Cette section illustre l'opération de réduction sur l'exemple de la Figure 2-21. Les variables V_i , V_j , V_k et V_z s'influencent directement et la variable V_z n'est pas mesurée. L'algorithme de réduction va consister à extraire V_z tout en conservant les influences entre les autres variables.

Une opération préliminaire doit être effectuée avant de faire l'opération de réduction à proprement dite. Nous avons introduit cette opération qui n'était pas proposée dans [Travé-Massuyès *et al.* 2001]. Elle s'applique dans le cas où la variable inconnue à éliminer (V_z sur la Figure 2-21) agit sur elle même via une autre variable (V_i sur la Figure 2-21). Dans ce cas :

- L'influence de V_z sur V_i est éliminée et pour toutes les variables V_k qui influencent V_z , différentes de V_i , une influence de V_k sur V_i est créée. Si cette opération n'est pas effectuée au préalable alors la réduction qui est décrite à la suite n'aboutit pas à une solution.

L'algorithme de réduction à proprement dit est ensuite appliqué :

- Pour toutes les causes directes V_k et V_i de V_z et pour toutes les conséquences directes V_j de V_z , deux nouvelles influences sont créées de V_k et V_i sur V_j . Les relations associées aux influences de $I_{\text{MCR}}(V_k, V_j)$ et $I_{\text{MCR}}(V_i, V_j)$ sont obtenues aisément en effectuant des compositions de relations. Les influences de V_k et V_i sur V_z sont éliminées.

L'annexe A détaille l'algorithme de réduction.

2.4.3 Approximation

Si l'influence entre la variable V_k et la variable V_j doit être négligée alors l'opération d'approximation est appliquée sur l'influence entre la variable V_k et la variable V_j : l'influence de la variable V_k sur la variable V_j est éliminée et les composants de l'ensemble $\text{Support}_{\text{Influence}}(I_{\text{MCR}}(V_k, V_j))$ sont ajoutés à l'ensemble $\text{Support}_{\text{Influence}}(\text{Var}_{V \rightarrow V, \text{causes directes}}(\text{variable}), V_j)$ des autres causes directes de V_j .

Nous avons élaboré cette opération pour pouvoir éliminer (négliger) une influence tout en gardant l'identité des composants qui étaient sur le support de cette influence. Nous verrons que cette information sera utilisée pour la localisation des défauts dans le chapitre 3.

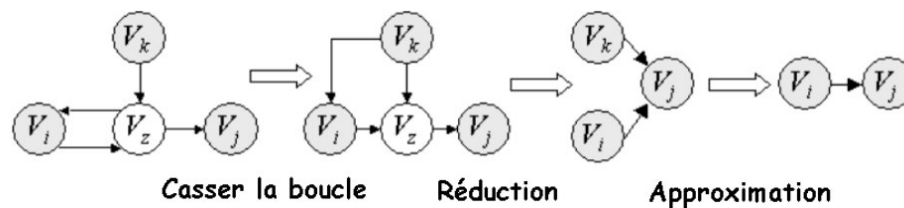


Figure 2-21 : Opérations de réduction et d'approximation

2.4.4 Conclusion

La Figure 2-22 résume les deux étapes de réduction et d'approximation. Ces étapes permettent respectivement d'exhumer du MCS les relations entre les variables mesurées puis de négliger certains phénomènes tout en garantissant une modélisation correcte.

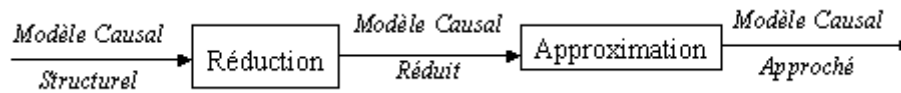


Figure 2-22: Obtention du modèle causal approché

2.5 Quantification du modèle

2.5.1 Introduction

Le modèle peut être linéaire ou non, il peut aussi avoir une forme particulière (réseaux de neurones, observateurs,...).

Le modèle dont nous disposons sur le pilote de FCC est un modèle de connaissance non linéaire dont nous ne connaissons pas tous les paramètres. Nous avons choisi de linéariser ce modèle autour d'un point de fonctionnement. La linéarisation permet d'obtenir de manière analytique la forme du modèle puis une identification des paramètres du modèle permet d'estimer les valeurs de ces paramètres. [Heim *et al.* 2003b] détaille cette manipulation. Nous avons choisi cette démarche car :

- Il n'existe pas, en général, de modèle couvrant toute la plage de fonctionnement d'un procédé,
- Cette approximation permet d'identifier assez aisément les paramètres du modèle.

Le domaine de validité de cette représentation est cependant limité et n'est justifié que si le point de fonctionnement ne change pas souvent.

La section 2.5.2 présente la représentation choisie. Deux approches sont ensuite envisageables pour quantifier les paramètres du modèle :

- L'approche classique est souvent stochastique : elle permet d'obtenir les paramètres sous forme de nombres réels (\mathbb{R}) (par opposition aux intervalles).
- L'approche par identification ensembliste est présentée dans la section 2.5.3.

Nous avons consacré l'Annexe B à l'identification ensembliste.

2.5.2 Représentation du modèle

Soit Y une variable endogène du MCA, les «n» causes directes de Y étant les variables notées $U_1, U_2, \dots, U_j, \dots, U_n$. Soit F_j , la fonction de transfert portée par l'arc entre U_j et Y (cf. Figure 2-23) :

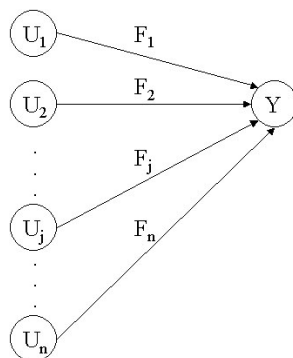


Figure 2-23 : Illustration de la fonction de transfert

La fonction de transfert F_j est définie par son numérateur B_j et son dénominateur A_j (z est l'opérateur décalage) :

$$F_j(z) = \frac{B_j(z)}{A_j(z)}, j = 1 \dots n \quad (2-61)$$

$$Y = \sum_{j=1 \dots n} F_j U_j(z)$$

Soit X une variable quelconque du MCA. La mesure de X à l'instant z_k est notée $X^m(z_k)$, et son point de fonctionnement est noté X° . L'écart de la mesure au point de fonctionnement est noté :

$$\Delta X^m(z_k) = X^m(z_k) - X^\circ \quad (2-62)$$

2.5.3 Identification des paramètres du modèle

Dans l'annexe B nous décrivons les deux méthodes d'identification que nous avons choisies.

La première méthode, stochastique, est très brièvement décrite car elle est classique.

Nous verrons dans la suite de ce manuscrit que nous avons besoin de connaître des intervalles contenant les paramètres. Par conséquent, nous avons testé une méthode d'identification ensembliste permettant d'obtenir ces intervalles. Cette

méthode s'appuie sur les travaux de [Leseq, Barraud et Tran-Dinh 2003]. Des ellipsoïdes sont utilisées pour définir une région qui contient certainement les valeurs des paramètres. Puis, pour obtenir les paramètres sous forme d'intervalles, il est nécessaire de projeter l'ellipsoïde sur les axes. Les paramètres ainsi projetés sont très pessimistes pour le diagnostic. Par ailleurs, la taille de l'ellipsoïde varie beaucoup avec les entrées, donc les expérimentations pour l'identification demandent à être soigneusement planifiées.

2.5.4 Conclusion

Nous avons fait le choix de linéariser le modèle autour du point de fonctionnement. En pratique, il nous aurait été très difficile et très long d'obtenir un modèle non linéaire.

Il nous a été ensuite nécessaire de quantifier les paramètres du modèle. Nous avons utilisé des méthodes classiques et éprouvées qui donnent des résultats satisfaisants.

Parallèlement, nous avons testé une méthode d'identification ensembliste dont les résultats n'ont pas pu être exploités (par la méthode de détection ensembliste) car nous avons induit un pessimisme lors de la projection des paramètres sur les axes. D'autre part cette méthode nécessite de disposer d'enregistrements dans lesquels les entrées sont bien excitées et nous disposons pour le pilote de FCC d'enregistrements de mauvaise qualité.

2.6 Conclusion

Dans ce chapitre, nous avons présenté une méthode d'ordonnancement causal. Le modèle causal obtenu décrit les influences entre les variables d'un système physique en comportement normal.

Nous avons jugé nécessaire d'élaborer une méthode préliminaire de génération du modèle structurel des relations, afin de bien situer et poser le problème. Sept étapes préliminaires permettent de générer le MSR (cf. Figure 2-8). Le MSR est un système composé d'équations élémentaires, répondant au principe de localité, et contient des variables connues et non connues. Le support de ses relations est identifié.

Nous appliquons ensuite un algorithme d'ordonnement causal au MSR. Cet algorithme permet de générer le modèle causal structurel. Le MSR est divisé en un sous-système statique et un sous-système dynamique. La causalité des systèmes dynamiques est imposée par les notions de causalité intégrale et de causalité différentielle (cf. 2.3.2.1). La causalité des systèmes statiques est obtenue, par exemple, par l'algorithme de Ford et Fulkerson.

Le MCS contient les mêmes informations que le MSR, cependant les influences entre les variables sont orientées.

Deux opérations dites de réduction puis d'approximation sont ensuite appliquées sur le MCS (cf. Figure 2-22).

L'opération de réduction appliquée au MCS permet d'extraire les influences entre les variables connues : le modèle causal réduit est obtenu.

Nous avons finalement élaboré l'opération d'approximation qui est appliquée au MCR. Elle permet de générer le modèle causal approché : cette opération permet de négliger des influences tout en conservant la connaissance qui garantit une modélisation puis un diagnostic corrects (le MCA est utilisé comme support du module de diagnostic). Dans notre application, les influences du MCA sont associées à des fonctions de transfert classiques. Leurs paramètres sont estimés ou peuvent être obtenus via une connaissance physique (données de la littérature, paramètres connus ou mesurés, etc.). Le MCA est donc un modèle dynamique de bon comportement du système. Il propose des approximations et explique les influences entre les variables. Cette méthode a été appliquée à un pilote de FCC (cf. Chapitre 4).

Comme le MCA permet de prédire le comportement normal des variables et d'expliquer les phénomènes de propagation des défauts, il constitue un outil pertinent pour le diagnostic à base de modèles. Des résidus spécifiques peuvent être extraits du MCA. Des alarmes sont ensuite obtenues en comparant les sorties du modèle aux mesures. Les influences du MCA sont associées à des composants physiques (qui constituent leur support). Cette connaissance et les alarmes permettent de suspecter les composants susceptibles d'être défectueux. Le Chapitre 3 présente le MCA au sein d'un module de diagnostic.

Chapitre 3

Méthode de diagnostic

Le Chapitre 2 a présenté une méthode permettant de construire le modèle causal approché (MCA) d'un procédé industriel continu quelconque. Le Chapitre 3 présente une méthode de diagnostic s'appuyant sur ce MCA : l'enchaînement de trois techniques complémentaires de génération d'alarmes, de localisation et d'identification de défauts, permet d'obtenir, à partir des observations, un message qui identifie la panne sur un composant et propose des actions à entreprendre.

L'approche de redondance analytique est utilisée pour générer des alarmes. Des relations de redondance analytique (RRA, cf. 1.5.3.2) spécifiques sont extraites du modèle causal approché. Les RRA sont vérifiées en comparant les sorties des relations aux mesures. Le modèle étant souvent imprécis et incertain, deux approches différentes, permettant de prendre en compte ces incertitudes, ont été explorées et comparées.

Les relations du MCA constituent des RRA. Connaissant les supports de ces relations, la localisation consiste ensuite à établir les diagnostics à partir des alarmes. Nous nous plaçons dans le cadre de diagnostic basé sur la cohérence tel que défini par [Reiter 1987]. Un diagnostic [Reiter 1987] est un ensemble de composants qui doivent tous se comporter anormalement pour expliquer les observations.

Une méthode à base de règles est finalement utilisée pour l'identification des défauts. Chaque composant physique est associé à une base de règles qui constitue un modèle semi-qualitatif de mauvais fonctionnement du composant. Les bases de règles associées aux composants qui font partie d'un diagnostic sont activées. Si les symptômes associés à une défaillance du composant sont observés via un traitement du signal des mesures (reconnaissance de formes, détection de sauts, ...), alors un message qui identifie la panne sur un composant est transmis à l'opérateur.

3.1 Introduction

Le MCA quantifié permet de générer deux références pour chacune des variables endogènes : une référence dite «globale» et une référence dite «locale». La référence globale correspond au régime nominal de fonctionnement : elle est évaluée à partir des valeurs des consignes et des perturbations connues. La référence locale informe sur le comportement attendu étant données les valeurs mesurées³⁻¹ des variables causes directes. Pour chaque variable endogène du MCA, les valeurs des deux références sont comparées à la mesure. La section 3.2 présente les références qui sont les sorties du MCA.

Une alarme dite «globale» est obtenue en comparant la mesure à la référence globale. Une alarme dite «locale» est obtenue en comparant la mesure à la référence locale. La section 3.3 présente ces deux alarmes. Ces comparaisons ne peuvent être effectuées telles quelles car le modèle est souvent incertain.

La plupart des simulateurs considèrent avoir une connaissance totalement déterministe du système. En réalité, les systèmes complexes sont soumis à des incertitudes (la structure du modèle n'est pas connue) ou au mieux à des imprécisions (la structure du modèle est connue mais ses paramètres sont imprécis) [Bonarini et Bontempi 1994]. Lorsque le système est bien connu, le modèle est très souvent simplifié pour être approprié à une application industrielle. Les paramètres du modèle peuvent évoluer au cours du temps de façon imprévisible mais toujours acceptable.

³⁻¹ Le terme « mesure » correspond à un abus de langage. En effet, il ne s'agit pas de mesures mais de valeurs transmises par le SCADA (qui comprennent les mesures, les consignes, les actions).

Pour toutes ces raisons, il est nécessaire de proposer des techniques permettant de prendre en compte les incertitudes et les imprécisions.

La section 3.4 présente qualitativement deux techniques permettant de prendre en compte ces incertitudes. La première consiste à utiliser un simulateur classique (dont les paramètres sont des réels²) et à représenter sous la forme d'une valeur réelle la sortie du modèle, l'imprécision étant prise en compte par un raisonnement sur les résidus fondé sur la logique floue [Evsukoff 1998]. La seconde approche consiste à représenter la sortie du modèle sous la forme d'intervalles (les paramètres du modèle sont des intervalles). Un certain nombre d'auteurs ont déjà exploré les techniques ensemblistes pour le diagnostic, en particulier [Adrot 2000], [Loiez et Taillibert 1997], [Ploix et Follot 2001], [Janati, Adrot et Ragot 2001], [Taillibert 1998], [Fagarasan, Ploix et Gentil 2001]. Nous avons, de notre côté, étudié l'apport des intervalles modaux proposés par [Armengol 1999], qui est un partenaire dans le projet CHEM. La section 3.5 propose une comparaison des méthodes de génération d'alarmes sur des données simulées. Des valeurs sont proposées pour les paramètres de détection et les méthodes sont comparées selon des critères précis.

Chaque influence du MCA est associée à un ensemble de composants (capteurs inclus) qui constituent son support. À partir de ces associations, et des alarmes précédemment générées, un diagnostic est établi. Nous nous sommes inspirés du domaine du diagnostic à base de modèles de l'IA pour l'élaboration des diagnostics [Reiter 1987]. La section 3.6 présente la méthode de localisation des défauts.

Lorsqu'un composant est contenu dans un diagnostic³⁻³, un système à base de règles est activé (système expert). Il analyse automatiquement la connaissance experte sur ce composant [Heim, Cauvin et Gentil 2001]. Si les symptômes associés à une défaillance sont observés (par traitement du signal) suivant des tests logiques, alors un message qui identifie la panne sur le composant est généré. Ce message, qui apparaît sur l'interface opérateur, contient une phrase en langage naturel qui décrit la défaillance supposée de ce composant, propose des actions à entreprendre pour valider

²⁻² Par opposition à un paramètre donné sous la forme d'un intervalle.

³⁻³ Ici le terme diagnostic est employé dans le cadre du diagnostic basé sur la cohérence [Reiter 1987]

l'information et pour rétablir le bon fonctionnement (en s'appuyant sur la connaissance experte). La section 3.7 présente la méthode d'identification des défauts.

3.2 Génération des références

3.2.1 Introduction

Rappelons que le modèle causal approché (MCA) est un modèle dynamique et que toutes ses variables sont connues. Par conséquent, un grand nombre de relations de redondance analytique peut en être extrait (par exemple, la relation liant une variable à toutes ses causes à une profondeur donnée) [Fagarasan, Ploix et Gentil 2001]. Cette section propose d'utiliser, pour chaque variable, deux références (implicitement deux RRA) particulières et pertinentes pour le diagnostic. La première est la référence globale (évaluée à partir des valeurs des consignes et des perturbations connues). Cette référence permet de déterminer si la variable est cohérente avec les consignes du procédé. La seconde est la référence locale (évaluée à partir des mesures des variables qui influencent directement la variable). Elle permet de déterminer si la variable est cohérente vis-à-vis de ses causes directes. Les notions de référence locale et de référence globale sont décrites dans [Evsukoff 1998].

3.2.2 Références globales

3.2.2.1 Génération des références globales

Soit une variable Y . La valeur de la référence globale de Y à l'instant z_k est notée $Y^{globale,BO}(z_k)$. La différence entre $Y^{globale,BO}(z_k)$ et le point de fonctionnement de Y (noté Y°) est notée $\Delta Y^{global,BO}(z_k)$, selon :

$$\Delta Y^{global,BO}(z_k) = Y^{globale,BO}(z_k) - Y^\circ \quad (3-1)$$

La valeur $\Delta Y^{global,BO}(z_k + 1)$ est calculée à partir des valeurs des références globales de ses causes directes $\Delta U_j^{globale,BO}(z_k)$ (qui sont elles-mêmes calculées à partir des valeurs des références globales de leurs causes directes jusqu'à ce que l'on calcule une référence globale à partir des variables exogènes).

Cet écart est obtenu, pour l'exemple de la Figure 2-23, selon :

$$\Delta Y^{globale,BO}(z_k + 1) = (1 - A_j(z)) \cdot \Delta Y^{globale,BO}(z_k) + \sum_{j=1}^{j=n} B_j(z) \cdot \Delta U_j^{globale,BO}(z_k) \quad (3-2)$$

La sortie mesurée du processus n'est pas réinjectée dans le modèle (cf. Figure 3-1). Il s'agit d'une simulation en **boucle ouverte**. La sortie du modèle est calculée depuis l'instant initial avec les valeurs des entrées et des sorties précédentes du modèle.

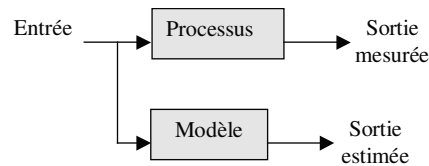


Figure 3-1 : Simulation sans ré-injection des sorties mesurées dans le modèle

3.2.2.2 Ordre de calcul des références globales

Comme la valeur de la référence globale d'un nœud dépend des valeurs des références globales de ses causes directes, il apparaît que l'ordre de calcul des références globales n'est pas arbitraire et doit être spécifié.

Un algorithme (cf. Annexe C) permet de déterminer cet ordre. Les valeurs des nœuds sont ensuite calculées selon le nombre croissant de causes

3.2.2.3 Modèle utilisé pour générer les références globales

Nous avons choisi d'utiliser un modèle déterministe³⁻⁴ pour générer les alarmes globales. Un modèle ensembliste (cf. note en bas de page) pourrait être aussi envisagé pour générer les références globales. Ce cas n'a cependant pas été étudié car les références globales se propageant de variables en variables, un modèle ensembliste générerait des intervalles beaucoup trop larges pour les variables ayant beaucoup d'antécédents.

³⁻⁴ Le modèle déterministe est un modèle classique dans lequel les paramètres sont des réels par opposition au modèle ensembliste où les paramètres et la sortie sont des intervalles.

3.2.2.4 Bilan sur la référence globale

La référence globale donne le comportement attendu de la variable étant données les valeurs des variables exogènes (ou pseudo-exogènes) du MCA. La référence globale est calculée en propageant les références globales de proche en proche depuis les entrées du MCA jusqu'à la variable calculée.

Une autre référence dite locale est définie pour chaque variable du MCA : elle informe sur le comportement local de la variable.

3.2.3 Génération des références locales

Pour la forme du modèle de la référence locale, nous avons envisagé plusieurs possibilités. Nous proposons dans la suite de ce document des critères (modèle déterministe ou modèle ensembliste, confiance dans le modèle, détection robuste ou sensible) pour choisir la façon d'utiliser le modèle.

Soit une variable Y . La valeur de la référence locale de Y à l'instant z_k , notée $Y^{locale,BO}(z_k)$, est calculée à partir des mesures U_j^m des causes directes de Y . La différence entre $Y^{locale,BO}(z_k)$ et le point de fonctionnement Y° de Y est noté $\Delta Y^{locale,BO}(z_k)$ selon :

$$\Delta Y^{local,BO}(z_k) = Y^{loaale,BO}(z_k) - Y^\circ \quad (3-3)$$

La valeur $\Delta Y^{local,BO}(z_k + 1)$ peut être calculée, pour l'exemple de la Figure 2-21, en boucle ouverte selon le schéma bloc de la Figure 3-1 :

$$\Delta Y^{locale,BO}(z_k + 1) = (1 - A_j(z)) \cdot \Delta Y^{locale,BO}(z_k) + \sum_{j=1}^{j=n} B_j(z) \cdot \Delta U_j^m(z_k) \quad (3-4)$$

La référence locale peut aussi être obtenue en **boucle fermée** : la mesure de sortie est alors réinjectée selon le schéma de la Figure 3-2. Ce schéma s'apparente à un observateur très simple (prédiction de la sortie à un pas à partir des mesures passées).

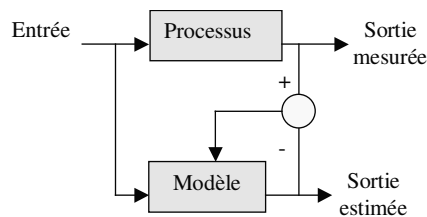


Figure 3-2 : Simulation avec ré-injection des sorties mesurées dans le modèle

Ce document propose d'obtenir la référence locale à l'instant z_k , notée $Y^{locale,Nk}(z_k)$ et son écart par rapport au point de fonctionnement noté $\Delta Y^{locale,Nk}(z_k)$, par une simulation sur un horizon glissant de taille est $N_k \in \mathbb{N}^*$. A chaque valeur de N_k correspond une valeur de $Y^{locale,Nk}(z_k)$ à l'instant z_k . Tout d'abord nous notons :

$$\Delta Y^{locale,Nk}(z_k) = Y^{locale,Nk}(z_k) - Y^o \quad (3-5)$$

Pour calculer la sortie du modèle $\Delta Y^{locale,Nk}(z_k + N_k)$ à l'instant $z_k + N_k$, nous définissons les intermédiaires de calcul $\Delta Y_1^{locale,B.F}(z_k + 1)$, $\Delta Y_2^{locale,B.F}(z_k + 2)$, ..., $\Delta Y_{N_k-1}^{locale,B.F}(z_k + N_k - 1)$. Pour l'exemple de la Figure 2-21 nous obtenons la sortie du modèle $\Delta Y^{locale,Nk}(z_k + N_k)$ par la formule récurrente :

$$\begin{aligned} \Delta Y_1^{locale,B.F}(z_k + 1) &= (1 - A_j(z)) \cdot \Delta Y^m(z_k) + \sum_{j=1}^{j=n} B_j(z) \cdot \Delta U_j^m(z_k) \\ \Delta Y_2^{locale,B.F}(z_k + 2) &= (1 - A_j(z)) \cdot \Delta Y_1^{locale,B.F}(z_k + 1) + \sum_{j=1}^{j=n} B_j(z) \cdot \Delta U_j^m(z_k + 1) \\ \Delta Y^{locale,Nk}(z_k + N_k) &= (1 - A_j(z)) \cdot \Delta Y_{N_k-1}^{locale,B.F}(z_k + N_k - 1) + \sum_{j=1}^{j=n} B_j \cdot \Delta U_j^m(z_k + N_k - 1) \end{aligned} \quad (3-6)$$

La Figure 3-1 illustre la formule (3-6). L'initialisation $\Delta Y^m(z_k)$ est la même pour le modèle et le processus.

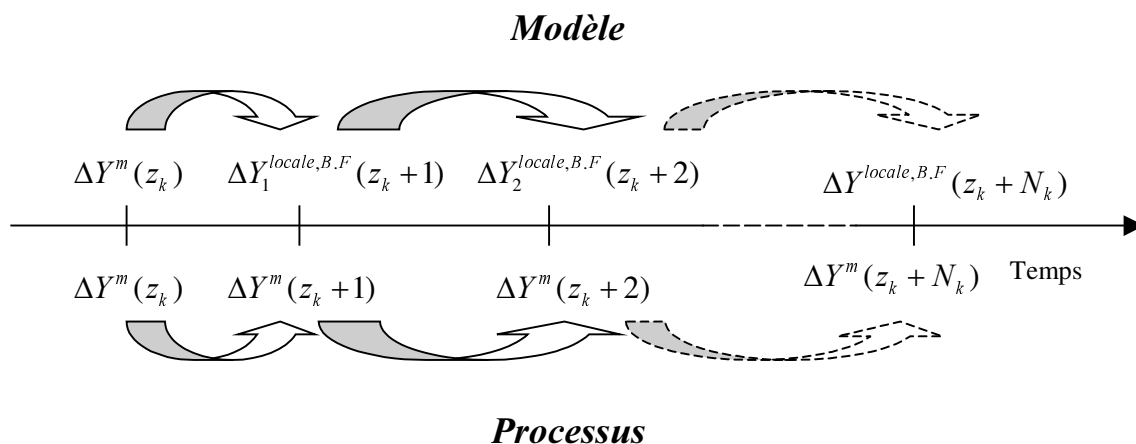


Figure 3-3 : Simulation sur une fenêtre de taille N_k

Le cas $N_k = 1$ correspond à une prédiction à un pas. Si la confiance dans le modèle est importante, alors la référence locale est calculée en boucle ouverte. Par exemple, une boucle de régulation peut être représentée par une consigne «C» qui agit sur une mesure «M». Le modèle étant une fonction de transfert F d'ordre 1 [Heim, Cauvin et Gentil 2000a]. La confiance dans un tel modèle est importante car le gain de 1 est parfaitement connu et le temps de réponse est en pratique souvent reproductible [[Heim, Cauvin et Gentil 2000a].

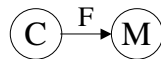


Figure 3-4 : Graphe causal d'une représentation simplifiée d'une régulation

Si le modèle est très imprécis (par exemple une vanne avec du jeu), il est nécessaire de corriger continuellement sa sortie : on a intérêt à utiliser un algorithme en boucle fermée. En pratique on fixe la taille de la fenêtre pour chaque variable (par exemple 10). Nous verrons qu'il est aussi possible d'envisager des stratégies traitant plusieurs tailles de fenêtres différentes (par exemple 1, 2 ,5 ,10) pour une seule variable.

Remarque 3-1 : La ré-injection de la mesure en boucle fermée est utilisée lorsque la confiance dans le modèle est faible. L'inconvénient d'une telle manipulation est que le modèle peut suivre une dérive lente engendrée par une défaillance. La méthode de détection n'est alors plus sensible à ce type de défauts.

3.2.4 Conclusion

Deux références sont générées par le MCA, la référence globale et la référence locale. La référence globale est calculée en boucle ouverte à partir des valeurs des variables exogènes du MCA. La notion de boucle fermée n'a pas été envisagée pour le calcul des références globales car celles-ci doivent rester indépendantes des mesures des variables endogènes. La référence locale d'une variable est calculée à partir des valeurs de ses causes directes, en boucle ouverte ou fermée. Une troisième technique consistant à effectuer une simulation en boucle fermée sur des fenêtres glissantes est proposée (cf. relation (3-6)). Ces références, qui sont comparées aux mesures des variables, permettent de générer des alarmes. Les sections suivantes présentent les alarmes générées.

3.3 Génération des alarmes

3.3.1 Introduction

La section 3.3 présente les alarmes générées en comparant pour chaque variable la référence globale et la mesure puis en comparant la référence locale et la mesure. Si le modèle est déterministe alors l'indicateur de défaut est un résidu. Si le modèle est ensembliste alors l'indicateur de défaut est l'appartenance de la mesure à un intervalle qui est la sortie du modèle.

3.3.2 Les alarmes globales

Pour le modèle déterministe, le résidu global $\rho_X(z_k)$ de la variable endogène X est défini à l'instant z_k comme la différence entre la mesure de X et sa référence globale par :

$$\rho_X(z_k) : \Delta X^m(z_k) - \Delta X^{\text{globale}, B.O}(z_k) \quad (3-7)$$

Un seuil bas $X_{\text{SB,global}}$ d'alarme et un seuil haut $X_{\text{SH,global}}$ d'alarme sont utilisés pour caractériser le résidu global. Si ce dernier est dans l'intervalle $[X_{\text{SB,global}}, X_{\text{SH,global}}]$, alors l'alarme globale vaut 0 sinon l'alarme globale vaut 1. Nous proposons dans la suite de ce document des valeurs pour les seuils $X_{\text{SB,global}}$ et $X_{\text{SH,global}}$.

Remarque 3-2 : Une défaillance peut ne pas se manifester ou se manifester de manière intermittente (par exemple lorsque le système est excité). Un résidu (le résidu global en particulier) faible ne préjuge donc absolument pas de l'absence de défaut.

Le support de la RRA $\rho_X(z_k) = 0$ est grand³⁻⁵. Si le résidu global $\rho_X(z_k)$ est grand, alors il est difficile de déterminer quel composant est en cause car tous les composants de son support sont candidats. Les références globales apportent donc peu d'information sur la localisation des défauts.

³⁻⁵ Il contient tous les composants des supports des influences directes de X et des supports (capteurs exclus) des influences de toutes les causes endogènes de X , ainsi que le capteur des variables exogènes prédécesseurs de X .

Par contre, l'information apportée par les références globales est pertinente pour le diagnostic pour les deux raisons suivantes :

- Les alarmes globales permettent d'expliquer la propagation des alarmes par rapport à un comportement global attendu par l'opérateur. Elles sont donc très importantes pour la conduite des procédés.
- Une dérive lente peut, par exemple, éloigner une variable de sa consigne. Les alarmes globales sont sensibles à ce type de défauts. Il est possible que les alarmes locales ne soient, pour des raisons numériques (cf. *Remarque 3-1*), pas sensibles à ces défauts.

Cette technique présente une lacune. Il serait envisageable qu'un jeu de consignes, fixées par l'opérateur, entraîne le processus dans un mode de fonctionnement dangereux. Les alarmes que nous proposons ne permettent pas de détecter ces défaillances humaines. Des alarmes à seuils classiques permettraient de prévenir une telle situation.

3.3.3 Génération des alarmes locales

Pour le modèle déterministe, le résidu local $\lambda_X(z_k)$ de la variable X est défini à l'instant t comme la différence entre la variation de la mesure de X et sa référence locale selon:

$$\lambda_X(z_k) : \Delta X^m(z_k) - \Delta X^{locale}(z_k) \quad (3-8)$$

Un seuil bas $X_{SB,local}$ et un seuil haut $X_{SH,local}$ d'alarmes sont aussi utilisés pour caractériser le résidu local⁶.

³⁻⁶ Nous proposons dans la suite de ce document des valeurs pour ces seuils.

Pour le modèle ensembliste, une autre grandeur est utilisée : l'intervalle de sortie du modèle. Si l'intersection de la mesure (qui est un intervalle)⁷ et de cet intervalle est vide, alors un défaut est détecté et l'alarme vaut 1, sinon elle vaut 0. La Figure 3-5 illustre le principe :

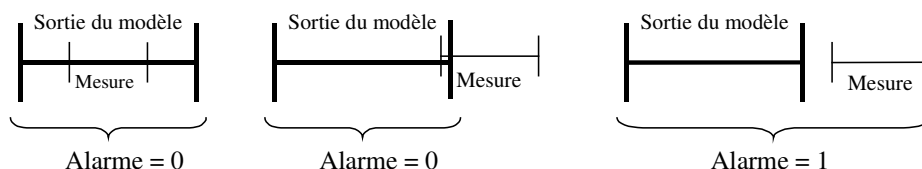


Figure 3-5 : Indicateurs de défauts de la méthode ensembliste

Le support de la RRA $\lambda_X = 0$ est minimal car les relations du MSR sont élémentaires (cf. note de bas de page 2-2). Autrement dit, si $\lambda_X = 0$ n'est pas observé alors un minimum de composants sont candidats. Les références locales permettent donc d'obtenir des informations précises sur la localisation des défauts et donc de produire un diagnostic pertinent.

Prise en compte de la dynamique du défaut

Pour appréhender des défauts de dynamiques différentes, la référence locale est calculée sur des fenêtres glissantes de tailles différentes⁸ (cf. formule (3-6)). Cette technique agit directement sur la manière dont la sortie du modèle est calculée et elle ne s'appuie pas sur des paramètres à fixer (seuils de détection, etc.). Elle n'est donc pas inhérente à une technique de détection : elle peut être appliquée à un modèle déterministe ou ensembliste.

Une grande taille de fenêtre permet d'augmenter la sensibilité de la méthode de détection : des défauts de dynamiques plus lentes (dérives lentes) peuvent être détectés. Cependant, la méthode devient plus sensible aux erreurs de modélisation.

⁷ Nous verrons dans la suite de ce document comment représenter la mesure sous la forme d'un intervalle.

⁸ Prédiction sur un horizon glissant.

Remarque 3-3 : Si la taille de la fenêtre augmente, alors l'incertitude sur la sortie du modèle augmente et la méthode de détection est moins sensible. Une solution permettant de prendre en compte ces phénomènes antagonistes consiste à utiliser des tailles de fenêtres différentes.

Nous verrons que l'algorithme de détection basé sur le modèle ensembliste est particulièrement adapté pour traiter, de manière simple, les fenêtres de tailles différentes. Par contre dans le cas du modèle déterministe, il est nécessaire de déterminer les seuils pour chaque taille de fenêtre. Par conséquent, nous avons décidé de traiter des fenêtres de tailles différentes dans le cas du modèle ensembliste mais nous n'avons traité qu'une fenêtre de taille 1 (i.e. prédiction à un pas) dans le cas du modèle déterministe.

3.3.4 Conclusion

Deux types d'alarmes sont extraits du MCA. L'alarme globale indique si le comportement globalement attendu est normal. L'alarme locale permet de localiser quelles sont les variables qui se comportent localement anormalement (ou normalement). Les indicateurs de défauts sont différents selon la forme du modèle.

3.4 Techniques de détection

3.4.1 Introduction

Cette section présente deux techniques de détection permettant de prendre en compte les incertitudes du modèle et des mesures pour la décision de détection. Nous choisissons deux approches pour représenter les paramètres du modèle.

La première approche consiste à utiliser un modèle déterministe, (dont les paramètres sont des réels)⁹ et à représenter sous la forme d'une valeur réelle sa sortie. Dans ce cas, l'imprécision est prise en compte lors de l'interprétation des résidus (différence entre la sortie du modèle et la mesure). Nous nous sommes appuyés sur les travaux de [Evsukoff 1998] pour l'interprétation par le raisonnement par logique floue

³⁻⁹ Nous verrons dans la suite de ce document que les paramètres peuvent être représentés par des intervalles.

des résidus. La logique floue, présentée dans la section 3.4.3, nécessite d'interpréter simultanément les résidus et leur variation.

La seconde approche consiste à représenter la sortie du modèle sous la forme d'un intervalle (les paramètres du modèle sont des intervalles), ce modèle sera dit «ensembliste». Nous avons choisi d'utiliser les intervalles modaux pour générer les intervalles, sorties du modèle. Nous nous sommes appuyés sur le travaux de [Armengol 1999] qui propose des stratégies permettant de les dédier à la détection. La section 3.4.4 présente les travaux de [Armengol 1999], [Goldsztejn 2000].

Avec ces deux techniques, la décision repose sur des paramètres de réglage permettant d'augmenter sa sensibilité ou sa robustesse³⁻¹⁰.

3.4.2 Logique binaire

Définissons tout d'abord, selon [Evsukoff 1998], les notions de défaut amont et de défaut local, d'une variable endogène X quelconque :

- Un défaut est dit amont si X est cohérente avec les mesures de ses causes directes mais n'est pas cohérente avec les variables exogènes du MCA.
- Un défaut est dit local si la mesure de X n'est pas cohérente avec ses causes directes.

Nous remarquons alors que :

- les résidus globaux ρ_X sont des indicateurs de défauts amont,
- les résidus locaux λ_X sont des indicateurs de défauts locaux.

Par conséquent, il vient la table de signature des défauts (fautes simples) du Tableau 3-1. Le signe «✓» signifie que le résidu est un indicateur pour le défaut.

	Défaut Amont	Local
ρ_X	✓	
λ_X		✓

Tableau 3-1 : Table de signature des défauts

³⁻¹⁰ Nous donnerons le résultat de nos réflexions pour le choix de ces paramètres dans la suite de cette section.

L'assertion suivante est extraite du Tableau 3-1 :

«Si λ_X est grand alors il y a un défaut local : le composant défaillant est dans le support des influences de X ».

Sous hypothèse d'exonération³⁻¹¹, l'assertion suivante est vérifiée :

«Si ρ_X est grand, et λ_X est petit, alors il y a un défaut amont. Le composant défaillant est dans le support des influences des causes non directes de X ».

Nous proposons dans cette section une interprétation qualitative (résidu grand ou petit) des résidus. Nous verrons dans la suite de ce document qu'il est nécessaire de fixer des seuils pour effectuer cette décision binaire.

Pour prendre la décision, [Evsukoff 1998] interprète d'abord la valeur du résidu global puis ensuite la valeur du résidu local. Si le résidu global est petit, alors il accorde moins d'importance à l'interprétation de la valeur résidu local³⁻¹². Dans la plupart des cas, si le résidu local est grand, alors le résidu global l'est aussi.

De manière générale, le résidu global est très souvent grand lorsque le résidu local est grand. Néanmoins, dans le cas pratique du pilote de FCC nous avons observé des situations où le résidu global est petit alors que le résidu local est grand. Ce phénomène est lié au système de régulation dont le rôle est de compenser les perturbations ou les défaillances pour maintenir les variables régulées à leurs consignes. Si la mesure d'une variable régulée est bien maintenue à sa consigne, alors son résidu global reste faible, mais cela ne préjuge en rien la valeur du résidu local. Nous avons, par conséquent, dissocié l'interprétation du résidu global et du résidu local.

Les résidus ne peuvent être exploités tels quels car les imprécisions des mesures et des modèles seraient à l'origine de détections intempestives de défauts. Une méthode basée sur le raisonnement approché est proposée dans la section 3.4.3. Elle permet de rendre plus graduelle la description de ces résidus. Une méthode ensembliste est présentée dans la section 3.4.4.

³⁻¹¹ L'hypothèse d'exonération suppose qu'une défaillance se manifeste toujours au niveau du résidu.

³⁻¹² Il pondère le résultat du raisonnement sur le résidu local par le résultat du raisonnement sur le résidu global.

3.4.3 Logique floue

La logique floue est une théorie mathématique qui, dans son application pratique, permet par exemple de tenir compte des incertitudes et permet une fusion des informations [Klir et Yuan 1995]. En logique binaire, l'appartenance d'un objet à un ensemble est vraie ou ne l'est pas. La logique floue est plus « souple » que la logique binaire car l'appartenance d'un objet à un ensemble flou est graduelle. La mise en œuvre de la logique floue se fait en trois étapes la plupart du temps. La fuzzification, le raisonnement par logique floue et la défuzzification. Ces trois étapes peuvent se représenter par le schéma suivant :

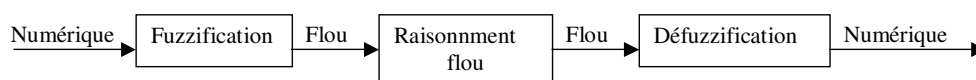


Figure 3-6 : Mise en œuvre de la logique floue sur des grandeurs numériques

La fuzzification consiste à rendre qualitatives des grandeurs quantitatives en définissant des fonctions d'appartenance à des sous ensembles flous. Autrement dit, elle permet d'associer des termes linguistiques à des valeurs numériques. Des opérateurs sur les ensembles flous sont définis. Ils doivent répondre à certaines conditions (monotonie, commutativité, associativité, etc.) qui leur confèrent des propriétés. L'utilisation des opérateurs se traduit par des opérations simples entre les valeurs prises par les fonctions d'appartenance. La défuzzification permet de repasser à une interprétation numérique de l'ensemble flou de la conclusion. Il existe plusieurs méthodes de défuzzification (centre de gravité, hauteur, utilisation de bornes, etc.).

Fuzzification

Le principe de raisonnement sur les résidus est exactement identique pour les résidus globaux et pour les résidus locaux. Cette section en présente les grands principes. Les algorithmes sont détaillés dans [Heim, Cauvin et Gentil 2000b], et [Evsukoff 1998] propose une description théorique des concepts.

L'ensemble descriptif $L = \{NN, N, Z, P, PP\}$ (cf. Figure 3-7) est utilisé pour décrire les valeurs des résidus. Les variables linguistiques sont interprétées de la manière

suivante : (NN pour négatif et très éloigné de 0, N pour négatif, P pour positif, PP pour positif et très éloigné de zéro et Z pour zéro). Si le résidu :

- est positif et inférieur au Seuil de Pré-Alarme (noté SPA Figure 3-7), alors sa valeur est considérée comme nulle étant données les incertitudes du modèle et des mesures,
- dépasse le Seuil de Réglage de sensibilité de l'Alarme (SRA), alors sa valeur est légèrement anormale,
- dépasse le Seuil d'Alarme (SA), alors sa valeur est inacceptable.

Les fonctions d'appartenance ont été choisies assez classiques, trapézoïdales ou triangulaires, de façon à ne nécessiter que trois paramètres de réglage SPA, SRA et SA pour chaque variable.

Les appartenances du résidu $r(t)$, exemple de la Figure 3-7, sont $0/NN$, $0/N$, $0,5/Z$, $0,5/P$ et $0/PP$.

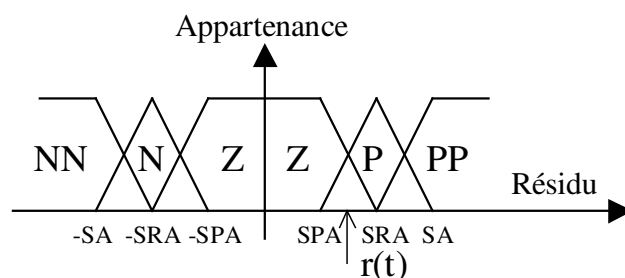


Figure 3-7 : Fuzzification des résidus

Trois notions sémantiques auraient pu être choisies pour déterminer l'état d'un résidu. Le choix de cinq notions permet d'obtenir une plus grande sensibilité de l'analyse. D'autres variables linguistiques auraient pu être utilisées entre Z et PP pour prendre en compte différentes nuances de la valeur des résidus. Il est préférable d'avoir quelques variables linguistiques qui se recouvrent bien plutôt que d'avoir beaucoup de variables linguistiques qui se recouvrent mal.

La méthode proposée par [Evsukoff 1998] interprète le résidu et sa variation. Les variations des résidus sont donc aussi fuzzifiées. La variation du résidu global est :

$$\Delta\rho_X(z_k+1) = \rho_X(z_k+1) - \rho_X(z_k) \quad (3-9)$$

La variation du résidu local est :

$$\Delta\lambda_X(z_k+1) = \lambda_X(z_k+1) - \lambda_X(z_k) \quad (3-10)$$

L'ensemble descriptif $\dot{L} = \{N^*, Z^*, P^*\}$ (cf. Figure 3-8) est utilisé pour décrire la variation des résidus. Il ne comprend que trois notions sémantiques car il est difficile de faire la différence entre une grandeur qui augmente lentement et une grandeur qui augmente vite. Si la variation du résidu :

- est positive et inférieure au Seuil de Pré-Alarme (SPA* sur la Figure 3-8), alors le résidu varie normalement,
- est supérieure au Seuil d'Alarme (SA*), alors le résidu augmente de manière anormale.

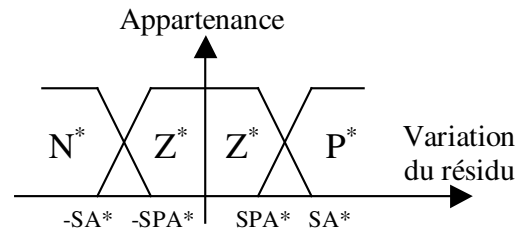


Figure 3-8 : Fuzzification des variations des résidus

Inférence (raisonnement par logique floue)

Les raisonnements traités dans cette étude sont du type :

«Si le résidu est négatif et sa variation est négative, alors la situation est grave».

Cette assertion se traduit pour le résidu global $\rho_X(t)$:

Si $\rho_X(t)$ est N et $\dot{\rho}_X$ est N*, alors «la situation est grave». (3-11)

Cette conclusion est partielle car elle se rapporte au couple (N,N*). Quinze conclusions partielles sont obtenues par combinaison des variables des ensembles L et \dot{L} . La conclusion, «la situation est grave» se traduit par une appartenance à un sous ensemble flou. S'il s'agit du résidu global, alors cet ensemble est construit sur le

référentiel symbolique $LCG = \{OK, AL\}$. La notion sémantique OK signifie qu'aucun défaut n'est détecté et AL qu'un défaut est détecté. Ainsi la relation (3-11) devient pour le résidu global :

$$\text{Si } \rho_X(t) \text{ est } N \text{ et } \rho_X(t+1) - \rho_X(t) \text{ est } N^*, \text{ alors } C_{N,N^*} = a_{NN,N^*,OK}/OK \text{ et } a_{N,N^*,AL}/AL \quad (3-12)$$

C_{NN,N^*} est la conclusion, encore appelée qualification symbolique de l'état de la variable analysée pour le couple (N, N^*) . $a_{NN,N^*,OK}$ et $a_{NN,N^*,AL}$ sont les degrés d'appartenance de la conclusion C_{NN,N^*} aux sous ensembles flous de LCG. Cette notation traduit le degré d'appartenance de la conclusion à la plage de fonctionnement normal OK et à la plage de fonctionnement anormal AL. Les appartenances des conclusions à l'ensemble LCG sont récapitulées dans le Tableau 3-2.

	N^*	Z^*	P^*
NN	0/OK et 1/AL	0/OK et 1/AL	0,2/OK et 0,8/AL
N	0/OK et 1/AL	0,4/OK et 0,6/AL	0,6/OK et 0,4/AL
Z	0,8/OK et 0,2/AL	1/OK et 0/AL	0,8/OK et 0,2/AL
P	0,6/OK et 0,4/AL	0,4/OK et 0,6/AL	0/OK et 1/AL
PP	0,2/OK et 0,8/AL	0/OK et 1/AL	0/OK et 1/AL

Tableau 3-2 : Degrés d'appartenance des conclusions à l'ensemble $LCG = \{OK, AL\}$

La fuzzification a permis de traduire numériquement l'assertion « $\rho_X(t)$ est N » de la relation (3-12). Le raisonnement par logique floue consiste ensuite à traduire le lien logique «et» et l'inférence «si ... alors» par des opérations numériques. Un choix judicieux de normes permet d'effectuer le raisonnement par logique floue via des calculs matriciels simples [Evsukoff 1998]. Le résultat est finalement une appartenance aux plages de fonctionnement OK et AL, (notées $\mu(OK)$ et $\mu(AL)$) obtenue en effectuant une agrégation des quinze conclusions partielles.

Le même raisonnement est effectué sur le résidu local. La conclusion est construite sur le référentiel symbolique $LCL = \{AM, LO\}$. La notion sémantique AM (amont) signifie que la référence locale est cohérente avec la mesure et la notion sémantique LO (local) signifie que la référence locale n'est pas cohérente avec la mesure.

Décision

La prise de décision qui attribue la valeur des alarmes est booléenne. Aucune défuzzification n'est effectuée. La description par logique floue est transformée en prise de décision binaire selon les règles suivantes :

- Si $\mu(\text{AL}) = 1$ alors l'alarme globale de X vaut 1,
- Si $\mu(\text{LO}) = 1$ alors l'alarme locale de X vaut 1.

[Evsukoff 1998] propose des grandeurs intermédiaires qui permettent de rendre plus ou moins robuste la méthode de détection. Nous n'avons pas utilisé ces possibilités dans le cadre du pilote de FCC car les défauts sont souvent brutaux.

Remarque 3-4 : La logique floue est parfois considérée, à tort, comme un filtre du signal. La méthode qui a été précédemment présentée s'applique à un signal peu bruité. Un filtre classique doit préalablement être utilisé pour atténuer le bruit sur les signaux. Un signal très bruité entraîne des fonctions d'appartenances elles mêmes bruitées.

Les paragraphes suivants présentent deux techniques permettant de régler la sensibilité de la méthode de détection à base de logique floue. Ces techniques ne sont cependant pas inhérentes à cette méthode et peuvent s'appliquer à tout modèle déterministe.

Prise en compte de la persistance des défauts

Une solution permettant de prendre en compte la persistance du défaut consiste à analyser non pas la valeur instantanée du résidu mais ses « n_{res} » dernières valeurs. Une grande valeur pour « n_{res} » réduit le taux de fausses alarmes mais retarde la décision. La méthode permet d'intégrer simplement cette technique par des agrégations de matrices [Evsukoff 1998]. La même démarche consistant à travailler sur une fenêtre « n_{res^*} » pour la variation du résidu est implémentée dans le système ASCO.

Prise en compte de la dynamique du défaut

La prise en compte de la variation du résidu (cf. Figure 3-8) permet d'appréhender la dynamique du défaut et d'anticiper la décision. Par exemple, si le résidu est grand et augmente vite (dynamique rapide du défaut), alors la méthode considère que la situation est plus grave que s'il est grand et reste constant ou diminue.

Une autre solution consisterait à calculer la sortie du modèle sur un horizon glissant (cf. formule (3-6)). Néanmoins, seule une prédiction à un pas a été implémentée.

Cette section a présenté l'interprétation des résidus selon [Evsukoff 1998] que nous avons choisie de suivre.

3.4.4 Méthode ensembliste

3.4.4.1 Introduction

Les paramètres du modèle sont souvent connus avec une précision finie qui peut être représentée par des intervalles. Le modèle doit tenir compte de cette imprécision. Un problème classique en calcul ensembliste est de trouver l'ensemble des valeurs prises par une fonction réelle continue lorsque ses variables appartiennent à des ensembles fixés. L'évaluation exacte de l'image de l'extension naturelle est souvent trop complexe pour être faite. Soit f une fonction réelle continue définie sur une partie A de \mathbb{R}^n . L'image de son extension naturelle est l'ensemble des valeurs prises par f lorsque les variables parcourent l'ensemble A :

$$W(f, A) = \{f(x), x \in A\} \tag{3-13}$$

L'arithmétique classique des intervalles permet de définir une extension de la fonction continue, définie uniquement sur la restriction des intervalles de \mathbb{R}^n . Les calculs sur les réels sont remplacés par des calculs sur les bornes des intervalles.

Les opérations arithmétiques classiques d'addition «+» et de soustraction «-» sont par exemple définies par les formules suivantes ($a_1 < b_1$ et $a_2 < b_2$) :

$$[a_1; b_1] + [a_2; b_2] = [a_1 + a_2; b_1 + b_2] \quad (3-14)$$

$$[a_1; b_1] - [a_2; b_2] = [a_1 - b_2; b_1 - a_2] \quad (3-15)$$

Par exemple, l'image de la fonction $f(X, Y) = X - Y$ avec $X = [0; 1]$ et $Y = [0; 1]$ est exactement l'intervalle $[-1; 1]$. La simplification induite par l'arithmétique classique des intervalles introduit un pessimisme dans le calcul en raison de la décomposition du problème en sous problèmes indépendants. Cette décomposition élimine les relations entre les différentes occurrences des variables dans l'expression de la fonction. Par exemple, l'image de la fonction $f(X) = X - X$ (deux occurrences de X) avec $X = [0; 1]$ est $[-1; 1]$ alors que le résultat exact est évidemment zéro. L'arithmétique classique des intervalles permet uniquement d'obtenir une enveloppe extérieure qui contient le résultat exact de l'extension naturelle. Les intervalles modaux décrits dans cette section, ont été mis au point par E.Gardenyes et M.A.Sainz [Gardeñes *et al.* 2001] ; Ils permettent de calculer une enveloppe intérieure et une enveloppe extérieure. Ces enveloppes s'avèrent pertinentes pour la détection de défauts [Armengol 1999].

L'arithmétique classique des intervalles présente une seconde lacune : aucune quantification explicite n'est associée aux intervalles (la quantification $\forall x$ est différente de la quantification $\exists x$).

Un exemple de l'utilité des quantificateurs est le suivant : l'équation $a + x = b$ avec les contraintes $a \in [-1; 1]$ et $b \in [5; 8]$ admet plusieurs solutions selon le choix de la quantification associée aux variables :

- $X = [6; 7]$ si la quantification est $\forall a \in [-1; 1] \quad \forall x \in X, \quad \exists b \in [5; 8], \quad a + x = b,$
- $X = [6; 7]$ si la quantification est $\forall b \in [5; 8] \quad \exists a \in [-1; 1], \quad \exists x \in X, \quad a + x = b,$
- $X = [4; 9]$ si la quantification est $\forall a \in [-1; 1] \quad \exists x \in X, \quad \exists b \in [5; 8], \quad a + x = b,$
- $X = [4; 9]$ si la quantification est $\exists x \in X \quad \exists b \in [5; 8], \quad \exists a \in [-1; 1], \quad a + x = b.$

La théorie des intervalles modaux, en associant un quantificateur aux intervalles, met en place un langage qui permet la modélisation de la quantification.

3.4.4.2 Brève présentation de la théorie des intervalles modaux

Nous avons choisi de nous appuyer sur les travaux de [Armengol 1999] pour la génération des intervalles modaux. De plus [Armengol 1999] les applique au diagnostic et propose des stratégies de détection particulièrement pertinentes.

La théorie des intervalles modaux est complexe au premier abord et ce document ne présente que ses grands principes, l'objectif principal étant d'en étudier l'apport pour la détection de défauts. La théorie des intervalles modaux permet de définir deux fonctions sémantiques généralement notées f^* et f^{**} , qui sont respectivement liées aux enveloppes extérieures et intérieures. Cette section présente des définitions et le principe d'obtention des deux fonctions sémantiques.

Définition des intervalles modaux

Un intervalle classique est caractérisé par l'ensemble des nombres réels compris entre ses bornes. Les intervalles modaux bénéficient d'une représentation plus riche. Ils sont constitués d'un support, qui correspond à un intervalle classique, et d'une modalité qui est un type de quantification (\forall ou \exists). L'ensemble des intervalles³⁻¹³ classiques étant noté $I(\mathbb{R})$, l'ensemble des intervalles modaux noté $I^*(\mathbb{R})$ est :

$$I^*(\mathbb{R}) = I(\mathbb{R}) \times \{\forall, \exists\} \quad (3-16)$$

Pour un intervalle modal $A \in I^*(\mathbb{R})$, l'intervalle classique qui lui est associé est appelé le support de A , noté $\text{Set}(A)$. Le quantificateur associé à sa modalité est noté $\text{Mod}(A)$. La définition de $\text{Set}(A)$ est nécessaire car les bornes des intervalles modaux ne sont pas nécessairement ordonnées de manière croissante (exemple : $\text{Set}([5;-3]) = [-3;5]$). Dans la suite de ce document, $[a,b]$ ($a \leq b$) désigne un intervalle classique et $[a,b]'$ un intervalle modal.

³⁻¹³ Les quantificateurs \forall et \exists sont aussi souvent notés U (universel) et E (existential).

Prédicats sur l'ensemble des intervalles modaux

Des prédicats sont définis sur l'ensemble des intervalles modaux. Un prédicat réel est une expression logique définie sur un ensemble qui peut être vraie ou fausse (exemple : $\sin(x)=0, x \in \mathbb{R}$). L'ensemble des prédicats réels est noté $\text{Pred}(\mathbb{R})$. Les prédicats, définis sur l'ensemble des intervalles modaux, permettent de différencier les intervalles modaux existentiels et universels. Ils permettent aussi de définir des opérations sur les intervalles modaux. L'ensemble des prédicats associés à un intervalle modal A' est noté $\text{Pred}(A')$.

L'ensemble $\text{Pred}(A', \exists)$ des prédicats associés à un intervalle existentiel A' contient tous les prédicats vrais pour au moins un réel de $\text{Set}(A')$ (où 1 représente vrai):

$$\text{Pred}(A', \exists) = \{P \in \text{Pred}(\mathbb{R}) / \exists x \in \text{Set}(A'), P(x) = 1\} \quad (3-17)$$

L'ensemble $\text{Pred}(A', \forall)$ des prédicats associés à un intervalle universel A' contient tous les prédicats vrais pour tous les réels de $\text{Set}(A')$:

$$\text{Pred}(A', \forall) = \{P \in \text{Pred}(\mathbb{R}) / \forall x \in \text{Set}(A'), P(x) = 1\} \quad (3-18)$$

Par exemple, sur la *Figure 3-9* :

- $P_1(x) \in \text{Pred}([-1;1]', \exists)$, $P_2(x) \in \text{Pred}([-1;1]', \exists)$ et $P_3(x) \notin \text{Pred}([-1;1]', \exists)$
- $P_1(x) \in \text{Pred}([-1;1]', \forall)$, $P_2(x) \notin \text{Pred}([-1;1]', \forall)$ et $P_3(x) \notin \text{Pred}([-1;1]', \forall)$

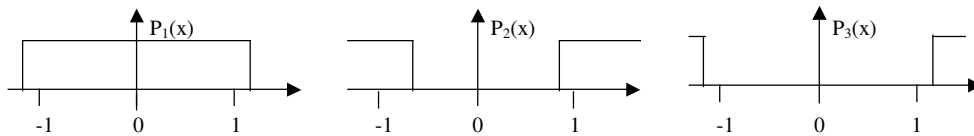


Figure 3-9 : Prédicats associés aux intervalles modaux

Inclusion modale

L'inclusion entre intervalles modaux est définie par l'inclusion ensembliste sur l'ensemble des prédicats :

$$\begin{aligned} A \in I^*(\mathbb{R}), B \in I^*(\mathbb{R}) \\ A \subset B \Leftrightarrow \text{Pred}(A) \subset \text{Pred}(B) \Leftrightarrow \forall P \in \text{Pred}(\mathbb{R}), P \in A \Rightarrow P \in B \end{aligned} \quad (3-19)$$

Des propriétés et des théorèmes sont attribués à l'inclusion modale. L'inclusion modale permet d'apporter un ordre partiel sur l'ensemble des intervalles modaux.

Union et intersection d'intervalles modaux

L'ordre partiel apporté aux intervalles modaux permet de définir la borne supérieure et la borne inférieure de deux intervalles modaux A et B qui correspondent respectivement à l'union et à l'intersection. Les définitions formelles sont les mêmes que pour les intervalles classiques, ce sont des définitions générales de bornes inférieures et supérieures :

$$\begin{aligned}
 & A \in I^*(\mathbb{R}), B \in I^*(\mathbb{R}) \\
 & A \cup B = C \Leftrightarrow (\forall X \in I^*(\mathbb{R}), (A \subseteq X \text{ et } B \subseteq X) \Leftrightarrow C \subseteq X) \\
 & A \cap B = C \Leftrightarrow (\forall X \in I^*(\mathbb{R}), (X \subseteq A \text{ et } X \subseteq B) \Leftrightarrow X \subseteq C)
 \end{aligned} \tag{3-20}$$

L'union de A et B est le plus petit intervalle qui contient (au sens de l'inclusion modale) A et B. L'intersection de A et B est le plus grand intervalle contenu (au sens de l'inclusion modale) dans A et B. L'expression classique des opérations d'union et d'intersection se prolonge aux intervalles modaux. Soient $[a_1, b_1]$ et $[a_2, b_2]$ deux intervalles modaux quelconques :

$$[a_1, b_1] \cup [a_2, b_2] = [\min\{a_1, a_2\}, \max\{b_1, b_2\}] \tag{3-21}$$

$$[a_1, b_1] \cap [a_2, b_2] = [\max\{a_1, a_2\}, \min\{b_1, b_2\}] \tag{3-22}$$

Les opérateurs union et intersection ne sont pas commutables entre eux. L'union est un quantificateur dit universel et l'intersection est un quantificateur dit existentiel.

Il est intéressant de constater que l'on peut écrire l'extension naturelle à l'aide d'une réunion d'intervalles. On retrouvera cette utilisation d'intervalles dégénérés pour définir les extensions sémantiques des fonctions aux intervalles modaux.

$$W(f, A) = \cup_{x \in A} [f(x); f(x)] \tag{3-23}$$

Fonctions sémantiques

Chaque fonction f réelle continue de \mathbb{R}^n dans \mathbb{R} permet de construire deux fonctions sémantiques f^* et f^{**} qui sont des extensions de f aux intervalles modaux. f^* et f^{**} sont à valeur de $I^*(\mathbb{R}^n)$ vers $I^*(\mathbb{R})$. Les fonctions sémantiques f^* et f^{**} sont définies sur $I^*(\mathbb{R}^3)$ par :

- $\text{Mod}(X) = \exists$ et $\text{Mod}(Y) = \exists \Rightarrow f^*(X,Y)=f^{**}(X,Y)= \bigcup_{x \in X', y \in Y'} [f(x,y), f(x,y)]$,
- $\text{Mod}(X) = \forall$ et $\text{Mod}(Y) = \forall \Rightarrow f^*(X,Y)=f^{**}(X,Y)= \bigcap_{x \in X', y \in Y'} [f(x,y), f(x,y)]$
- $\text{Mod}(X) = \exists$ et $\text{Mod}(Y) = \forall \Rightarrow f^*(X,Y)= \bigcup_{x \in X'} \bigcap_{y \in Y'} [f(x,y), f(x,y)]$,
- $\text{Mod}(X) = \forall$ et $\text{Mod}(Y) = \exists \Rightarrow f^*(X,Y)= \bigcup_{y \in Y'} \bigcap_{x \in X'} [f(x,y), f(x,y)]$,
- $\text{Mod}(X) = \exists$ et $\text{Mod}(Y) = \forall \Rightarrow f^{**}(X,Y)= \bigcap_{y \in Y'} \bigcup_{x \in X'} [f(x,y), f(x,y)]$,
- $\text{Mod}(X) = \forall$ et $\text{Mod}(Y) = \exists \Rightarrow f^{**}(X,Y)= \bigcap_{x \in X'} \bigcup_{y \in Y'} [f(x,y), f(x,y)]$,

Si les intervalles X et Y sont de même modalité, alors $f^*(X,Y)=f^{**}(X,Y)$. L'intersection est toujours appliquée aux arguments universels et l'union aux arguments existentiels. Pour la fonction sémantique f^* , l'intersection est toujours effectuée en premier et pour la sémantique f^{**} en deuxième.

Les deux extensions f^* et f^{**} sont liées aux enveloppes extérieures et intérieures (au sens de l'inclusion modale) :

$$\forall X \in I^*(\mathbb{R}), f^{**}(X) \subset f^*(X) \quad (3-24)$$

Extension rationnelle

L'extension rationnelle modale $fR(X)$, qui correspond à l'image exacte, est définie exactement de la même manière que l'extension rationnelle classique, en utilisant les opérateurs modaux. Son utilisation est beaucoup plus complexe que dans la théorie classique principalement à cause de l'existence des deux extensions sémantiques qui sont différentes dans le cas général. L'étude de l'optimalité consiste à rechercher les conditions suffisantes pour lesquelles l'extension rationnelle coïncide avec les extensions sémantiques. Par exemple, les fonctions uni-incidentes vérifient $f^{**}(X) \subset fR(X) \subset f^*(X)$ quel que soit l'intervalle modal X considéré. Dans ces conditions, l'optimalité est équivalente à la condition $f^{**}(X)=f^*(X)$.

Dans le cas général, des conditions sur f doivent être vérifiées pour étudier la modalité :

- la monotonie des fonctions considérées,
- les modalités des variables impliquées.

Un processus d'optimisation Branch and Bound est appliqué. Il permet de séparer l'espace étudié en parties possédant des propriétés (monotonie de f , modalités des arguments, etc.) sur lesquelles des théorèmes s'appliquent. Ce processus permet d'obtenir un intervalle surévalué de plus en plus petit et un intervalle sous évalué de plus en plus grand. Ils convergent vers l'extension naturelle $fR(X)$.

Conclusion

Les intervalles modaux permettent de générer deux intervalles (ou enveloppes) :

- un intervalle surévalué caractérisé par $*$ qui contient l'intervalle exact,
- un intervalle sous-évalué caractérisé par f^{**} qui est contenu dans l'intervalle exact.

Un critère d'optimisation (par exemple distance entre ces deux intervalles) permet d'itérer les calculs et de faire converger les deux intervalles vers l'intervalle exact. Ceci est illustré par la figure suivante :

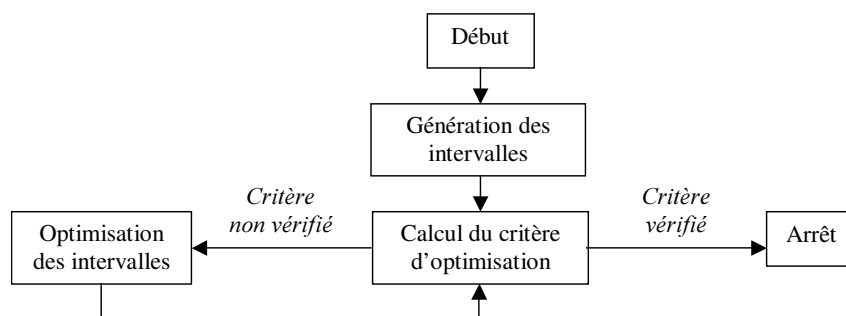


Figure 3-10 : Principe d'optimisation des intervalles

3.4.4.3 Détection de défauts

Les intervalles, sorties du modèle, sont calculés sur des horizons glissants selon la formule (3-6). Ceci permet de prendre en compte des défauts de dynamique différente. La taille des fenêtres (valeur de N_k) est automatiquement choisie selon une stratégie qui est décrite dans la suite de ce paragraphe. Etant donné une taille de fenêtre fixée, la stratégie de détection est basée sur les deux propriétés suivantes :

- L'intervalle sous-évalué est contenu dans l'intervalle exact, donc toutes les valeurs situées dans cet intervalle sont des valeurs possibles de la sortie du modèle connaissant les incertitudes.
- L'intervalle sur-évalué contient l'intervalle exact, donc toutes les valeurs situées à l'extérieur de cet intervalle constituent des valeurs qui ne peuvent pas être prises par la sortie du modèle connaissant les incertitudes.

L'algorithme de détection présenté ci-après est illustré par la Figure 3-11.

```
Début
Initialiser la taille de la fenêtre
Tant Que (il reste du temps de calcul) ou (aucun défaut n'est détecté)
  Générer des intervalles sur un horizon glissant de taille n
  Tant Que (la mesure est entre l'intervalle sur-évalué et l'intervalle sous-évalué)
    Optimiser les intervalles
  Fin Tant Que
  Si la mesure est dans l'intervalle sous-évalué alors augmenter n sinon défaut
Fin Tant Que
Fin
```

Algorithme 3-1 : Stratégie de détection ensembliste

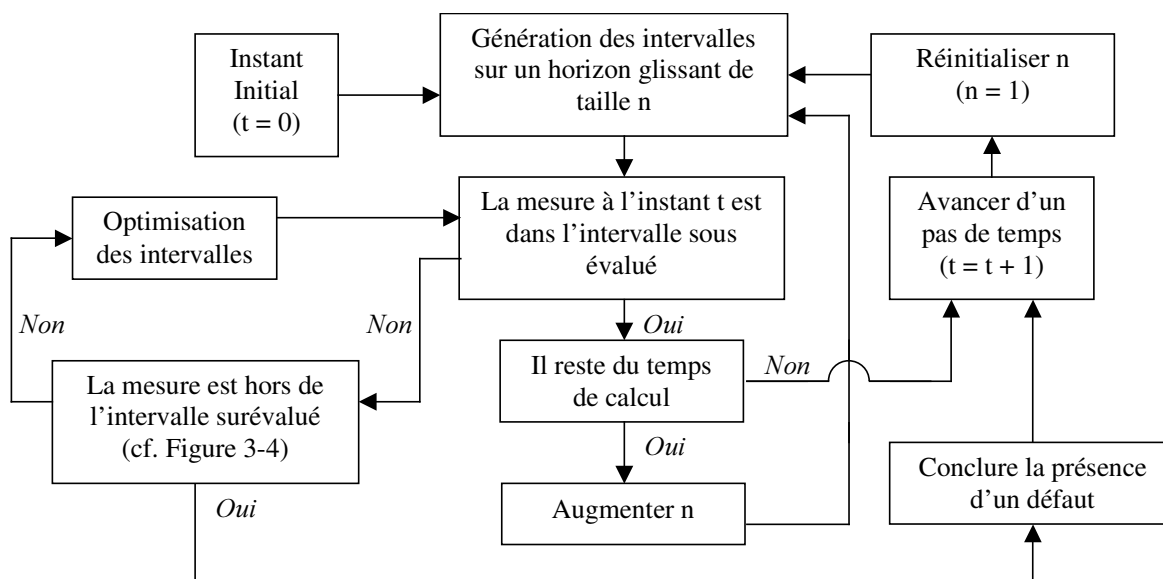


Figure 3-11 : Stratégie pour la détection de défauts

Prise en compte de la persistance du défaut

La prise en compte de la persistance du défaut consiste à conclure que l'alarme vaut 1 si et seulement si une série de $n_{\text{déf}}$ alarmes successives est observée.

Prise en compte de la dynamique du défaut

La prise en compte de la dynamique du défaut est inhérente à l'algorithme de détection. Plusieurs tailles de fenêtres sont étudiées (par exemple 1 puis 5 puis 15). Si le temps est suffisant, alors la fenêtre de taille 1 est d'abord traitée (défauts de dynamique rapide), puis les fenêtres de tailles croissantes sont successivement traitées pour envisager les défauts de dynamiques progressivement plus lentes. Cette propriété doit cependant être nuancée (cf. *Remarque 3-3*).

3.4.4.4 Conclusion

La prise en compte des incertitudes du modèle et des mesures peut être effectuée par un modèle ensembliste. Les calculs ensemblistes ne sont appliqués qu'à la référence locale. [Armengol 1999] propose d'utiliser la théorie des intervalles modaux pour effectuer les calculs ensemblistes. La stratégie originale de cette méthode est d'itérer pour optimiser la sortie du modèle si les résultats de la détection le requièrent. Nous notons qu'une hypothèse forte de la théorie des intervalles modaux consiste à considérer que chaque paramètre du modèle est assurément dans un

intervalle et qu'elle ne varie pas dans le temps. D'autres approches consistent à considérer que les paramètres varient de manière aléatoire dans un intervalle au cours du temps.

Les seuils de détection sont mutuellement dépendants des paramètres du modèle et des mesures et sont donc liés entre eux de manière automatique par le modèle ensembliste. Par opposition, si les seuils sont fixes, il peuvent être très pessimistes pour une variable X et beaucoup moins pour une autre variable Y, et par conséquent la variable Y peut être en alarme bien avant la variable X.

3.4.5 Conclusion

Cette section a présenté les grands principes de deux techniques de détection à base de modèles de bon fonctionnement.

L'une s'appuie sur un modèle déterministe et sur un raisonnement par logique floue, l'autre s'appuie sur un modèle ensembliste et une stratégie particulière de détection. Le raisonnement par logique floue permet d'interpréter le résidu global et le résidu local. La méthode ensembliste génère la sortie du modèle tout en l'interprétant. Elle n'a été appliquée que pour la référence locale. L'approche ensembliste propose d'utiliser des modèles sur des horizon glissants pour la prise en compte de la dynamique des défauts. Cette possibilité n'a pas été implémentée et constitue une perspective. L'approche par logique floue propose d'interpréter la variation du résidu pour anticiper les défaillances.

Dans les deux cas, les paramètres des algorithmes de détection peuvent être modifiés pour rendre les techniques plus sensibles ou plus robustes. La section suivante indique comment fixer ces paramètres et compare ces techniques sur des défauts types simulés.

3.5 Comparaison des techniques de détection

3.5.1 Introduction

Dans cette section, nous effectuons une comparaison des deux techniques de détection que nous venons de décrire. Nous proposons d'évaluer les apports et différences respectives. Nous avons établi une liste de critères qui apportent des informations pour cette comparaison.

La section 3.5.2 propose une comparaison qualitative et la section 3.5.3 propose une comparaison quantitative sur des données simulées.

3.5.2 Comparaison qualitative des techniques

3.5.2.1 Introduction

Nous avons remarqué que peu d'auteurs proposent des méthodes permettant de fixer les valeurs des paramètres des algorithmes de détection. En particulier, [Armengol 1999] ne propose pas de méthode pour déterminer les intervalles sur les paramètres du modèle ensembliste. Nous avons donc jugé important de proposer dans cette section des méthodes permettant de fixer ces paramètres :

- Dans la section 3.5.2.2, nous comparons, selon la méthode envisagée, les efforts nécessaires pour quantifier les paramètres du modèle. Ces paramètres sont obtenus hors ligne sur des données avec des méthodes d'identification.
- Dans la section 3.5.2.3, nous proposons des valeurs pour les paramètres de détection.

3.5.2.2 Effort de calage des paramètres des modèles

Identification à partir de méthodes stochastiques

Les méthodes stochastiques d'identification de paramètres sont éprouvées et faciles à mettre en œuvre. Elles donnent, sous certaines hypothèses, accès aux incertitudes sur les paramètres sous la forme d'écart-types. On trouve de nombreux outils de mise en œuvre, tels qu'une boîte à outil de Matlab.

Identification à partir de méthodes ensemblistes

Les méthodes d'identification ensemblistes (cf. Annexe B) ne sont pas utilisées de manière courante. Par conséquent, leur utilisation nécessite une expertise et l'utilisation d'outils non commerciaux.

Cette contrainte nous a incité, en pratique, à proposer des intervalles sur les paramètres via une approche stochastique, même si cette approche ne peut se justifier théoriquement. Par conséquent, parallèlement à l'identification ensembliste qui donne un paramètre θ sous la forme $[\theta_{\min}; \theta_{\max}]$, l'intervalle de la formule (3-25) sera proposé (σ_θ étant l'écart type sur le paramètre θ obtenu lors de l'identification stochastique). Des considérations statistiques³⁻¹⁴ assurent qu'il y a une probabilité de 0.0063 % qu'un échantillon soit en dehors de cet intervalle.

$$[\theta - 4\sigma_\theta; \theta + 4\sigma_\theta] \quad (3-25)$$

3.5.2.3 Effort pour le calage des paramètres de détection

Trois méthodes de détection sont distinguées dans la suite de ce manuscrit :

- une approche de référence,
- l'approche par logique floue,
- l'approche ensembliste.

Paramètres de détection des algorithmes de l'approche de référence

Une approche simple sera utilisée comme référence dans la suite de ce document. Elle consiste à générer une alarme par simple dépassement de seuil. Nous proposons de fixer ce seuil à partir de données d'apprentissage comme quatre fois l'écart type entre la mesure et la sortie du modèle. L'écart type σ_{mes_mod} , calculé sur N points entre la sortie du modèle $Y_{modèle}$ et la sortie mesurée $Y_{mesurée}$, est donné par :

$$\sigma_{mes_mod} = \sqrt{\sum_{k=1}^{k=N} \frac{(Y_{mesure}(k) - Y_{modèle}(k))^2}{N-1}} \quad (3-26)$$

Cette grandeur est en générale fournie par les algorithmes d'identification classiques.

¹⁴ Sous hypothèse d'un résidu d'identification aléatoire suivant une loi normale.

Paramètres de détection des algorithmes de l'approche par logique floue

Nous proposons de fixer les paramètres de détection de l'approche par logique floue selon des considérations statistiques. Ces paramètres SPA, SRA, SA, SPA* et SA* sont présentés sur la Figure 3-7 et la Figure 3-8.

SPA Seuil de pré-alarme

Le paramètre SPA est fixé comme l'écart type $\sigma_{\text{mes_mod}}$ défini par (3-26). Cet écart type prend en compte le bruit de mesure et l'erreur de modèle.

SRA Seuil de réglage pour la sensibilité du détecteur de la méthode par logique floue

Le paramètre SRA est fixé à $\frac{1}{2}(\text{SPA} + \text{SA})$. Au delà de cette valeur, la valeur du résidu est de moins en moins acceptable.

SA Seuil d'alarme

Le paramètre SA est fixé à $4 * \sigma_{\text{mes_mod}}$. Si le résidu dépasse cette valeur il est considéré anormal. Des considérations statistiques¹⁵ assurent que, au-delà de cette valeur, il y a 0.0063 % de chance d'avoir une fausse alarme.

SPA* Seuil de pré-alarme pour la variation du résidu

Ce seuil est fixé à $2 * \text{SPA}$. En deçà de cette valeur, le résidu augmente ou diminue de manière acceptable. Des données de mauvais fonctionnement peuvent être utilisées pour caler ce paramètre.

SA* Seuil d'alarme pour la variation du résidu

Ce seuil est fixé à $2 * \text{SA}$. Au delà de cette valeur, le résidu augmente ou diminue de manière anormale.

Il est à noter que, avec la méthode que nous proposons, un seul paramètre, $\sigma_{\text{mes_mod}}$, est nécessaire pour définir les 8 ensembles flous nécessaires au raisonnement et que ce paramètre est fourni naturellement lors de l'identification du modèle.

¹⁵ Sous hypothèse d'un résidu d'identification aléatoire suivant une loi normale centrée.

Remarque 3-5 : Nous remarquons qu'il est nécessaire de distinguer le seuil SPA pour le résidu local et le seuil SPA pour le résidu global. Dans la suite de ce document, nous notons respectivement $SPA_{L,X}$ et $SPA_{G,X}$ le seuil d'alarme du résidu local et le seuil d'alarme du résidu global de la variable X.

Il est ensuite nécessaire de fixer les paramètres « n_{res} » et « n_{res*} » qui permettent de régler la sensibilité de la méthode à la persistance des défauts (cf. 3.4.3). Nous avons fixé ces valeurs de manière arbitraire à 2 et 3. Elles peuvent ensuite être modifiées par un expert avec le raisonnement suivant : si on doit détecter des défauts de dynamique lente, on a intérêt à augmenter la taille de la fenêtre et pour des défauts de dynamique rapide, une fenêtre courte est suffisante. On notera qu'une fenêtre longue retarde la décision mais la rend plus robuste.

Paramètres de détection des algorithmes de l'approche ensembliste

Les deux paramètres à fixer dans le cadre de l'approche ensembliste sont :

- les incertitudes σ_{mes} sur chaque mesure;
- les tailles N_k des fenêtres glissantes pour la simulation sur un horizon glissant (cf. formule (3-6)). Des fenêtres de tailles de 3, 5, 10 et 15 périodes d'échantillonnage sont proposées par défaut.

L'incertitude sur la mesure est une caractéristique du capteur qui est donnée par le constructeur. Nous proposons, en son absence, de représenter la mesure par un intervalle qui est déterminé par l'écart type σ_{mes} du bruit de la mesure :

$$[measure - 4.\sigma_{mes} ; measure + 4.\sigma_{mes}] \quad (3-27)$$

L'écart type σ_{mes} est obtenu lorsque la mesure Y_{mesure} est stable au point de fonctionnement Y° :

$$\sigma_{mes} = \sqrt{\sum_{k=1}^{k=N} \frac{(Y_{mesure}(k) - Y^\circ)^2}{N-1}} \quad (3-28)$$

On a toujours $\sigma_{mes} < \sigma_{mes_mod}$ obtenu par la formule (3-26).

3.5.2.4 Conclusion

Les paramètres d'un modèle classique s'obtiennent facilement car les techniques d'identification sont bien répandues et maîtrisées. Il est plus difficile d'obtenir les paramètres du modèle ensembliste. L'étude engagée a montré que ces techniques sont prometteuses mais ne sont pas encore tout à fait au point. En pratique, des considérations statistiques permettent de définir des intervalles sur les paramètres, même si cette approche ne se justifie pas théoriquement.

Il peut être indispensable de disposer de données de validation (différentes des données d'identification) pour quantifier le paramètre $\sigma_{\text{mes_mod}}$ de détection de l'approche par logique floue (qui permet de fixer les bornes SPA, SRA... des classes floues).

Dans le cadre de l'approche ensembliste il n'y a pas de paramètres de détection comparables et seuls les intervalles sur les paramètres du modèle et la mesure doivent être déterminés.

Des paramètres sont proposés pour régler la sensibilité des méthodes mais ils peuvent être modifiés par l'expert en fonction des défauts à détecter.

Les méthodes peuvent être appliquées sur des fenêtres de longueur variable qui dépendent de la dynamique du défaut et donc de l'expertise sur le procédé.

3.5.3 Comparaison des techniques sur des données simulées

3.5.3.1 Introduction

Nous avons choisi de générer des données simulées par la fonction de transfert du premier ordre de la Figure 3-12. Nous avons ajouté :

- des incertitudes, δ_{a_1} , δ_{b_0} sur le pôle a_1 et le coefficient du numérateur b_0 ,
- des défauts multiplicatif, e_{a_1} et e_{b_0} sur a_1 et b_0 ,
- une incertitude δ_Y , un défaut e_Y et un bruit b à la sortie Y .

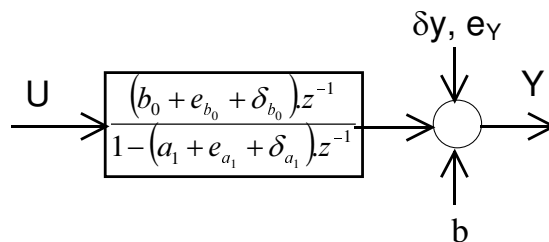


Figure 3-12 : Fonction de transfert illustrant les techniques de détection

Nous n'avons pas ajouté de bruit sur l'entrée. Pour les tests, la période d'échantillonnage est de 5 secondes avec les valeurs $a_1 = 0,95$ et $b_0 = 0,05$ et la durée des scénarios est de 2000 secondes.

L'entrée du système, dans tous les cas suivants est constituée de deux créneaux unitaires de $t = 20$ sec. à $t = 515$ sec. puis de $t = 1020$ sec. à $t = 1515$ sec.

Représentation des incertitudes du modèle

Nous avons ajouté des incertitudes sur les paramètres et sur la sortie, car en pratique, en fonctionnement normal, ils évoluent dans le temps. Une sinusoïde est utilisée pour représenter leur variation normale¹⁶. Une dérive lente aurait pu être choisie mais il est supposé que les paramètres évoluent sans cesse dans une plage définie. Ces évolutions sont supposées de dynamique lente donc des créneaux, qui sont brutaux, n'ont pas été envisagés. Les figures 3-12, 3-13 et 3-14 illustrent respectivement en bleu les incertitudes sur b_0 (δ_{s,b_0}), a_1 (δ_{s,a_1}) et Y ($\delta_{s,y}$). L'incertitude δ_{s,b_0} est contenue dans l'intervalle $[-0,025; 0,025]$, l'incertitude δ_{s,a_1} dans l'intervalle $[-0,0048; 0,0048]$ et l'incertitude $\delta_{s,y}$ dans l'intervalle $[-0,05; 0,05]$.

Nous remarquerons qu'une hypothèse très importante de la méthode de détection de l'approche ensembliste est que, en fonctionnement normal, les paramètres sont dans un intervalle qui est connu et qu'ils ne varient pas dans le temps³⁻¹⁷. Cette hypothèse rend les calculs d'optimisation sur un horizon glissant, via les intervalles modaux, particulièrement complexes. Par contre, si les paramètres sont considérés variables dans le temps, le problème d'optimisation est plus simple.

³⁻¹⁶ Nous aurions aussi pu choisir une somme de sinusoïdes pour représenter ces incertitudes.

³⁻¹⁷ D'autres approchent considèrent que les paramètres évoluent, dans le temps, dans un intervalle.

Certaines expériences ont été faites avec ces incertitudes mises à 0, ce qui correspond à une expérimentation avec des paramètres constants. Nous allons tout de même tester cette méthode sur des données avec des paramètres qui varient dans le temps pour étudier la sensibilité de la méthode à ces variations.

Représentation des défauts

Les défauts sont représentés par des échelons, des rampes ou des sinusoides.

- Des échelons (vert), qui simulent des défauts constants sont ajoutés au temps $t = 1000$ sec. sur b_0 (e_{p,b_0}), sur a_1 (e_{p,a_1}) et sur Y ($e_{p,Y}$). Ils sont respectivement illustrés en vert sur les figures 3-12, 3-13 et 3-14.
- Des rampes (bleu) qui simulent des défauts de dynamique lente sont ajoutés sur b_0 (e_{r,b_0}), sur a_1 (e_{r,a_1}) et sur Y ($e_{r,Y}$). Elles sont respectivement illustrées en bleu clair sur les figures 3-12, 3-13 et 3-14.
- Des sinusoides, qui simulent des défauts intermittents sont ajoutés (rose) sur b_0 (e_{s,b_0}), a_1 (e_{s,a_1}) et Y ($e_{s,Y}$). Ces sinusoides sont respectivement illustrées en rose sur les figures 3-12, 3-13 et 3-14. L'amplitude de ces sinusoides double au temps $t = 1000$ secondes. Nous remarquerons qu'un défaut représenté sous la forme d'une sinusoïde est intermittent (par exemple un tuyau qui se bouche et se débouche). Nous n'avons pas observé ce type de défaut sur le pilote de FCC. Ce choix de défaut nous permettra d'illustrer qu'une défaillance se manifester de manière intermittente.

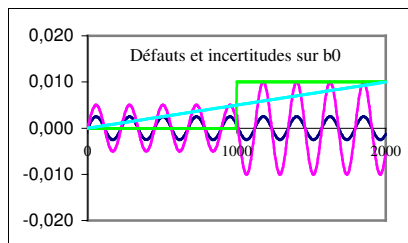


Figure 3-13 : Incertitudes et défauts sur b_0

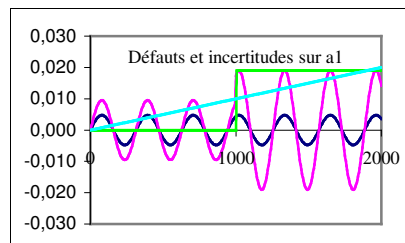


Figure 3-14 : Incertitudes et défauts sur a_1

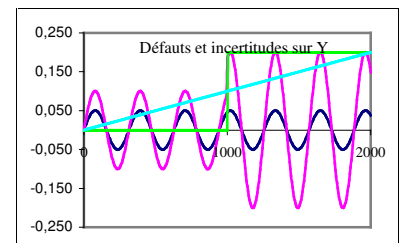


Figure 3-15 : Incertitudes et défauts sur Y

3.5.3.2 Création de données de bon et de mauvais fonctionnement

Nous avons construit différents types de données de bon fonctionnement (pour lesquelles les défauts sont nuls) :

- le modèle est idéal,
- la sortie est incertaine et les paramètres sont constants,
- la sortie du modèle est parfaitement connue et les paramètres du modèle sont incertains.

Nous avons ajouté, à la sortie, deux types b de bruit d'énergie faible (132 dB):

- un bruit uniforme (de borne 0,05), cf. Figure 3-16,
- un bruit blanc gaussien, cf. Figure 3-17.

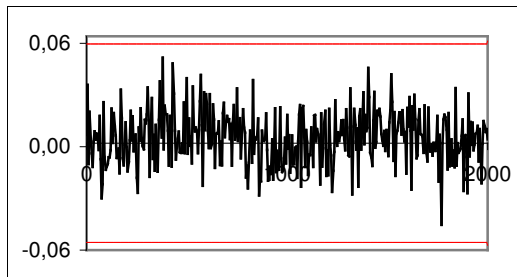


Figure 3-16 : Bruit borné

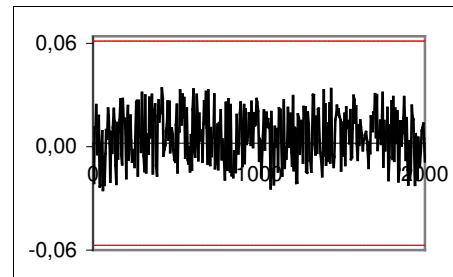


Figure 3-17 : Bruit blanc

- L'écart type du bruit uniforme est $\sigma_{\text{mes}} = 0,0144$. Par conséquent l'intervalle $[-4*0,0144 ; 4*0,0144]$ (cf. formule (3-28)), représenté en rouge sur la Figure 3-16, se justifie dans ce cas car il contient exactement le bruit.
- L'écart type du bruit blanc est $\sigma_{\text{mes}} = 0,0148$, et son amplitude maximale est 0,0593. Par conséquent l'intervalle $[-4*0,0148 ; 4*0,0148]$ (cf. formule (3-28)), se justifie dans ce cas et il contient exactement le bruit.

Le rapport signal d'un bruit b dont la moyenne sur N points est \bar{b} est calculé selon la formule (3-28) pour une mesure dont la valeur non bruitée $Y_{\text{non_bruitée}}$ a pour moyenne la valeur $\bar{Y}_{\text{non_bruitée}}$ sur ces N points. Sur le pilote de FCC nous avons un bruit très faible, de l'ordre de 70 dB.

$$\frac{S}{B} = 20 * \log_{10} \left(\frac{\sum_{k=0}^{k=N} (Y_{non_bruiée}(t_0 + k.Td) - \bar{Y}_{non_bruiée})^2}{\sum_{k=0}^{k=N} (b(t_0 + k.Td) - \bar{b})^2} \right) \quad (3-29)$$

D'autres essais ont été effectués avec des niveaux de bruit plus fort (30 dB), mais ils n'apportent pas d'information supplémentaire par rapport aux essais rapportés ci-dessous. Il est bien évident que le bruit dégrade fortement les performances des méthodes de détection. D'autre part les systèmes d'acquisition de données filtrent en général le signal.

3.5.3.3 Test sur des données de bon fonctionnement

Description des tests sur des données de bon fonctionnement

Nous proposons quatre situations de détection selon le choix du modèle et le choix de la méthode de détection.

- Le modèle correspond à une prédiction à un pas (formule (3-6)) et la détection est effectuée par l'algorithme de référence (cf. section 3.5.2.3). Cette situation est notée BF (Boucle Fermée).

Le paramètre de détection σ_{mes_mod} (cf. formule (3-26)) est obtenu, à partir de données d'identification, en comparant la mesure à la sortie du modèle en prédiction à un pas.

- Le modèle correspond à une simulation en boucle ouverte (formule (3-4)) et la détection est effectuée par l'algorithme de référence. Cette structure est notée BO (Boucle Ouverte).

Le paramètre de détection σ_{mes_mod} (cf. formule (3-26)) est obtenu, à partir des données d'identification, en comparant la mesure à la sortie du modèle en boucle ouverte.

- Le modèle correspond à une simulation en boucle ouverte et l'algorithme de détection est basé sur l'approche par logique floue. Cette structure est notée Flou BO.

Le paramètre de détection est σ_{mes_mod} (cf. formule (3-26)) est obtenu, à partir des données d'identification en comparant la mesure à la sortie du modèle en boucle ouverte. Les paramètres de réglage des classes floues sont obtenus selon la méthode décrite en 3.5.2.3. Les paramètres « n_{res} » et « n_{res*} » sont respectivement fixés à 2 et à 3.

- Le modèle est ensembliste et ses paramètres sont fixés par une identification stochastique. La méthode de détection correspond à la méthode ensembliste précédemment décrite. Cette situation est notée Inter. Les intervalles sur les paramètres sont obtenus, à partir des données d'identification, selon la formule (3-25). L'intervalle sur la sortie est donné par la formule (3-28). Les fenêtres de tailles de 1, 5, 10 et 15 périodes d'échantillonnage sont utilisées.

Le Tableau 3-3 énumère les résultats obtenus :

- La première colonne attribue pour chaque cas un numéro de référence de 1 à 9,
- Les trois premières lignes correspondent au cas sans bruit, les trois suivantes au cas avec un bruit borné et les trois dernières au cas avec un bruit blanc. La 2nd colonne rappelle le type de bruit.
- Pour chacun de ces trois cas de bruit, un modèle sans incertitude, un modèle avec une incertitude sur la sortie et un modèle avec des incertitudes sur les paramètres sont envisagés. Les colonnes 3, 4 et 5 informent sur les incertitudes.
- Les colonnes 6 et 8 donnent respectivement les paramètres b_0 et a_1 identifiés par la méthode du modèle avec un critère quadratique et la méthode d'optimisation de type descente de Gauss-Newton (algorithme OE de Matlab).
- Les colonnes 7 et 9 contiennent les écarts types σ_{b_0} et σ_{a_1} sur les paramètres b_0 et a_1 obtenus via la méthode d'identification précédente.
- La colonne 10 contient l'écart type σ_{BF} entre la sortie du modèle et la sortie mesurée¹⁸, en prédiction à un pas, cet écart type est utilisé pour fixer le seuil d'alarme de la méthode BF.
- La colonne 11 contient l'écart type σ_{BO} entre la sortie du modèle en boucle ouverte et la sortie mesurée, cet écart type est utilisé pour fixer le seuil d'alarme de la méthode BO.
- Les colonnes 12, 13, 14 et 15 contiennent le nombre de fausses détections observées sur ces données par les méthodes BF, BO, Inter et Flou BO.

¹⁸ Cet écart type est obtenu sur les données d'identification. En pratique, des données supplémentaires de validations peuvent être utilisées pour l'évaluer.

Cas	Qualité du modèle				Résultats de l'identification						Fausses detection			
	Bruit	Sortie	a_1	b_0	b_0	σ_{b_0}	a_1	σ_{a_1}	σ_{BF}	σ_{BO}	BF	BO	Inter.	Flou BO
1	-	0	0	0	0,05000	0,00E+00	0,9500	0,00E+00	3,44E-09	2,45E-09	0	0	395	0
2	-	$\delta_{S,Y}$	0	0	0,04944	1,66E-03	0,9505	1,89E-03	3,85E-03	3,55E-02	0	0	98	0
3	-	0	$\delta_{S,a1}$	$\delta_{S,b0}$	0,05042	6,51E-04	0,9492	6,87E-04	2,12E-03	1,99E-02	0	0	78	0
4	Borné	0	0	0	0,05022	1,58E-04	0,9498	1,68E-04	1,96E-02	1,44E-02	0	0	0	0
5	Borné	$\delta_{S,Y}$	0	0	0,04958	1,20E-03	0,9503	1,24E-03	2,00E-02	3,90E-02	0	0	0	0
6	Borné	0	$\delta_{S,a1}$	$\delta_{S,b0}$	0,05085	8,06E-04	0,9489	8,39E-04	1,97E-02	2,48E-02	0	0	0	0
7	Blanc	0	0	0	0,05009	1,43E-04	0,9499	1,52E-04	1,96E-02	1,43E-02	0	0	0	0
8	Blanc	$\delta_{S,Y}$	0	0	0,04955	1,03E-03	0,9504	1,07E-03	2,00E-02	3,90E-02	0	0	0	0
9	Blanc	0	$\delta_{S,a1}$	$\delta_{S,b0}$	0,05071	1,25E-03	0,9490	1,34E-03	1,97E-02	2,46E-02	0	0	0	0

Tableau 3-3 : Comparaison de techniques, données de bon fonctionnement

Limites des hypothèses de la méthode Inter.

Des alarmes sont observées pour la structure Inter dans les cas 1, 2 et 3.

- Dans le cas n°1, 395 alarmes sont générées par la méthode Inter. Cela provient d'un bruit numérique¹⁹ et n'est pas lié à la méthode de détection.
- Dans les cas 2 et 3, les paramètres du modèle varient dans le temps, alors que la méthode par intervalles suppose qu'ils sont constants. Ceci est la source d'alarmes. La méthode Inter. est donc, dans ce cas, particulièrement sensible à ces variations.

Dans les cas 5, 6, 8 et 9 les paramètres du modèle sont variables dans le temps alors que la méthode par intervalle n'est pas adaptée à cette situation. Nous avons tout de même effectué les tests et nous n'observons pas d'alarme.

Sur les données simulées nous avons imposé que le paramètre b_0 varie dans l'intervalle $[0,0475; 0,0525]$ (cf. *Figure 3-13*), et les intervalles $[b_0 - 4 * \sigma_{b_0}; b_0 + 4 * \sigma_{b_0}]$ obtenus selon la formule (3-25) sont :

- dans les cas 3 : $[0,04781; 0,05303]$,
- dans les cas 6 : $[0,04763; 0,05407]$,
- dans les cas 9 : $[0,04570; 0,05572]$.

Dans le cas 3, l'intervalle $[b_0 - 4 * \sigma_{b_0}; b_0 + 4 * \sigma_{b_0}]$ ne contient pas $[0,0475; 0,0525]$, ce qui montre que l'intervalle obtenu via l'identification classique peut être sous évalué.

¹⁹ Les intervalles sont réduits à des réels. La sortie du modèle et la mesure sont des valeurs numériques toujours différentes et leur comparaison engendre des alarmes.

Sur les données simulées (signaux discrets) nous avons imposé que le paramètre a_1 varie dans l'intervalle $[0,9452; 0,9548]$ (cf. *Figure 3-13*), et les intervalles $[a_1 - 4\sigma_1; a_1 + 4\sigma_1]$ obtenus selon la formule (3-25) sont :

- dans les cas 3 : $[0,9465; 0,9519]$
- dans les cas 6 : $[0,9455; 0,9523]$,
- dans les cas 9 : $[0,9436; 0,9544]$.

Dans les cas 3, 6 et 9 l'intervalle $[a_1 - 4\sigma_1; a_1 + 4\sigma_1]$ est sous évalué.

Nous utilisons néanmoins, dans ce manuscrit, la formule (3-25) pour paramétrer les intervalles du modèle Inter.

Limites des hypothèses de la méthode Flou BO

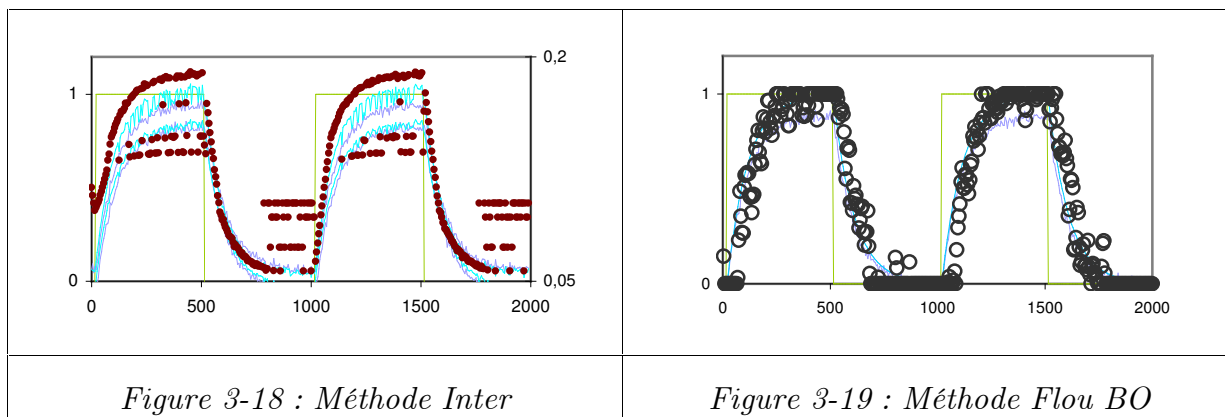
Nous avons effectué un test qui permet de mettre en valeur les limites des hypothèses de la méthode flou BO.

Ce test est effectué pour le modèle du cas n°6. Des données simulant le bon fonctionnement sont créées avec $b_0 = 0,0477$ et $a_1 = 0,9456$. Ces paramètres sont très proches des bornes que nous avons imposées en fonctionnement normal.

- La méthode BF a été testée et aucune fausse détection n'a été observée.
- La méthode BO a été testée et 71 fausses détections ont été observées.
- La méthode Inter a été testée et aucune fausse détection n'a été observée. La Figure 3-18 illustre en bleu gris la mesure. Les enveloppes internes et externes sont en bleu clair²⁻²⁰. L'axe de droite correspond à la différence entre la borne supérieure et la borne inférieure de l'enveloppe. Nous remarquons que cette différence augmente avec l'excitation du modèle.
- La méthode Flou BO a été testée et 65 fausses détections ont été observées. La Figure 3-19 illustre en bleu gris la mesure. La sortie du modèle en boucle ouverte est en bleu clair. Les appartenances à la classe alarme $\mu(\text{AL})$ sont représentées par des ronds noirs. On remarque que des fausses détections sont observées lorsque le système est excité. Dans le cadre de la

²⁻²⁰ Ces deux enveloppes sont confondues dans le cas particulier des fonctions de transfert du premier ordre.

méthode Flou BO, l'incertitude le seuil de détection SA sur le résidu est constant et il n'est pas fonction de l'excitation de l'entrée.



Cet exemple, qui correspond à un cas extrême, montre une limite de la méthode Flou BO. Dans la section suivante, nous testons les méthodes sur des données de mauvais fonctionnement.

3.5.3.4 Comparaison des techniques pour des données avec défaut

Nous avons testé les modèles identifiés dans les cas 1 à 9 sur des données de mauvais fonctionnement. Ces tests, qui sont recensés dans le Tableau 3-4, permettent de comparer la sensibilité des méthodes. Pour chacun des 9 cas, nous avons envisagé plusieurs situations :

- le défaut $e_{s,Y}$ est ajouté à Y (ligne 1 pour le cas n°1, ligne 3 pour le cas n°2...),
- les défauts e_{s,b_0} et e_{s,a_1} sont respectivement ajoutés à b_0 et a_1 (ligne 2 pour le cas n°1, ligne 4 pour le cas n°2 ...),
- le défaut $e_{p,Y}$ est ajouté à Y dans le cas n°4 (ligne 9),
- les défauts e_{p,b_0} et e_{p,a_1} sont ajoutés dans le cas n°4 (ligne 10),
- le défaut $e_{r,Y}$ est ajouté à Y dans le cas n°4 (ligne 11),
- les défauts e_{r,a_1} et e_{r,b_0} sont ajoutés dans le cas n°4 (ligne 12).

La méthode Inter n'a bien entendu pas été testée dans le cas n°1. La méthode Flou BO n'a pas été testée sur les cas 5 à 9 : nous verrons que les premiers cas sont suffisant pour illustrer son apport.

La méthode BO est la plus performante (sauf lignes 3 à 6) en terme de nombre de détection, cependant les autres méthodes apportent des raffinements (prise en compte de la dynamique ou de la persistance des défauts), qui les rendent plus sensibles ou retardent la décision et préviennent les fausses détections. La comparaison des techniques ne peut donc se résumer à la comparaison du nombre de détections.

Ligne	Cas	Défauts			Détections			
		Y	b_0	a_1	BF	BO	Inter.	Flou BO
1	1	$e_{S,Y}$	0	0	400	400		211
2	1	0	$e_{S,b0}$	$e_{S,a1}$	396	396		202
3	2	$e_{S,Y}$	0	0	95	107	155	56
4	2	0	$e_{S,b0}$	$e_{S,a1}$	48	33	61	17
5	3	$e_{S,Y}$	0	0	240	231	281	122
6	3	0	$e_{S,b0}$	$e_{S,a1}$	118	100	171	57
7	4	$e_{S,Y}$	0	0	0	288	130	268
8	4	0	$e_{S,b0}$	$e_{S,a1}$	0	157	60	146
9	4	$e_{P,Y}$	0	0	1	200	62	-
10	4	0	$e_{P,b0}$	$e_{P,a1}$	3	190	138	-
11	4	$e_{R,Y}$	0	0	0	295	31	-
12	4	0	$e_{R,b0}$	$e_{R,a1}$	0	275	110	-
13	5	$e_{S,Y}$	0	0	0	92	0	-
14	5	0	$e_{S,b0}$	$e_{S,a1}$	0	30	17	-
15	6	$e_{S,Y}$	0	0	0	178	89	-
16	6	0	$e_{S,b0}$	$e_{S,a1}$	0	71	33	-
17	7	$e_{S,Y}$	0	0	1	288	138	-
18	7	0	$e_{S,b0}$	$e_{S,a1}$	0	153	64	-
19	8	$e_{S,Y}$	0	0	1	91	66	-
20	8	0	$e_{S,b0}$	$e_{S,a1}$	0	30	18	-
21	9	$e_{S,Y}$	0	0	2	177	60	-
22	9	0	$e_{S,b0}$	$e_{S,a1}$	0	74	11	-

Tableau 3-4 : Récapitulatif des tests

Nous notons que ces tests ont été chacun effectués sur un seul essai. Il serait judicieux de les effectuer sur un grand nombre d'essais et de considérer la moyenne des alarmes recensées par exemple.

Défauts de forme échelon

Nous étudions dans un premier temps les défauts de forme échelon qui correspondent aux lignes 9 et 10.

Dans le cadre de la méthode BO, et pour le défaut échelon additif sur la sortie (ligne 9), nous ne dénombrons aucun manque à la détection (200 alarmes). Par contre pour les défauts échelons multiplicatifs (ligne 10), nous dénombrons 10 manques à la

détections (190 alarmes). Nous verrons dans la suite de cette section que les défauts multiplicatifs ne se manifestent que lorsque l'entrée est excitée. C'est ce phénomène qui explique les 10 manques à la détection.

Dans le cadre de la méthode Inter et pour le défaut échelon additif sur la sortie (ligne 9), nous dénombrons 138 manques à la détection (62 alarmes). Pour les défauts multiplicatifs échelons (ligne 10), nous dénombrons 62 manques à la détection (138 alarmes). Cette méthode est donc moins sensible que la précédente dans cette situation. Contrairement au cas précédent, nous obtenons plus d'alarmes (138) pour les défauts multiplicatifs, que pour le défaut additif (62).

La méthode BF est très peu sensible dans les situations des lignes 9 et 10.

Défauts rampe

Nous étudions dans un second temps les défauts de forme rampe qui correspondent aux lignes 11 et 12.

Dans le cadre de la méthode BO, et pour le défaut rampe additif sur la sortie (ligne 11), nous dénombrons 105 manques à la détection (295 alarmes). Pour les défauts multiplicatifs, nous dénombrons 125 manques à la détections (275 alarmes). La rampe est progressive dans le temps et, par conséquent, les manques à la détection ont lieu au début scénario. Au début du scénario il est difficile de faire la différence entre un défaut et une variation acceptable du paramètre.

La méthode Inter est beaucoup moins sensible à ces défauts (ligne 11 et 12) que la méthode BO. La méthode BF ne détecte aucun défaut. Les modèles en boucle fermée ont tendance à suivre les défauts de dynamique lente.

Défauts oscillants

Nous remarquons que pour les lignes 3, 4, 5 et 6, la méthode Inter génère plus d'alarmes que les autres méthodes. Elle est, dans cette situation, plus performante que les autres méthodes. Nous avons préalablement remarqué que dans les cas des modèles cas 2 et cas 3, nous avons des alarmes pour les faibles variations des

paramètres (Tableau 3-3). Ces tests confirment que, dans cette situation, la méthode est très sensible aux variations des paramètres et aux incertitudes et défauts sur la sortie.

Dans la suite de cette section, nous proposons d'illustrer les apports des techniques sur les 5 cas encadrés du Tableau 3-4 que nous détaillons.

La première situation que nous détaillons (cf. Figure 3-20) correspond au cas n°4, avec un défaut additif et la méthode de détection BO : **288** alarmes sont dénombrées. L'entrée (courbe verte) est composée de deux créneaux successifs excitant le système, la sortie mesurée (courbe bleu-gris) évolue autour de la sortie du modèle (courbe bleu clair). L'échelle de ces données est sur l'axe de gauche. L'axe de droite correspond aux valeurs prises par les alarmes (points noirs). Le défaut additif étant oscillants, les alarmes sont intermittentes.

La seconde situation que nous détaillons (cf. Figure 3-21) correspond au cas n°4 avec les défauts multiplicatifs et la méthode BO : **157** détections sont dénombrées. Durant la période pendant laquelle le système n'est plus excité (entre 650 et 1250 secondes) les défauts ne sont plus détectés. Ceci illustre que les défauts multiplicatifs ne se manifestent que lorsque le système est excité.

La troisième situation (cf. Figure 3-22) correspond au cas n°3 avec les défauts multiplicatifs et la méthode BF : **118** détections sont dénombrées alors que la méthode BO donne 100 détections. Dans les cas 4 à 9, avec présence de bruit la méthode BO détecte des défauts alors que la méthode BF ne détecte rien. La méthode BF permet de faire des corrections sur le modèle mais est moins sensible. Cette particularité peut être considérée comme un avantage ou un inconvénient selon le résultat souhaité.

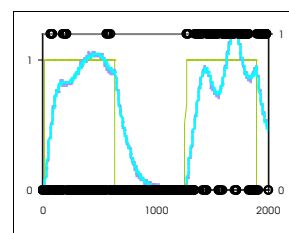
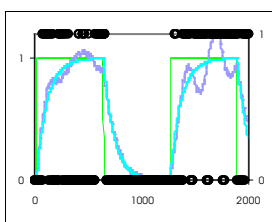
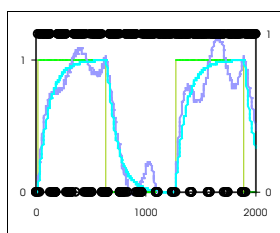


Figure 3-20 : Cas n°4, BO Figure 3-21 : Cas n°3, BO Figure 3-22 : Cas n°3, BF

La quatrième situation (cf. Figure 3-23) correspond au cas n°3 avec les défauts multiplicatifs et la méthode BO Flou. L'appartenance à la classe floue alarme ($\mu(\text{AL})$) est représentée par les ronds noirs : **57** alarmes ($\mu(\text{AL})=1$) sont dénombrées alors que la méthode BO donne 100 alarmes. Le nombre d'alarmes est plus faible que dans la méthode BO, cependant, dans la période de 700 à 900 secondes, $\mu(\text{AL})$ s'éloigne significativement de 0 alors que dans la même période, pour le cas comparable de la Figure 3-21, aucun défaut n'est détecté. Cet exemple illustre bien l'apport de l'évolution graduelle de $\mu(\text{AL})$ qui permet de savoir que la situation devient anormale alors que l'alarme ne vaut pas encore 1. Compter uniquement les situations $\mu(\text{AL}) = 1$ ne correspond pas vraiment à l'esprit de la méthode, qui est supposée attirer l'attention de l'opérateur sur des défauts naissants. Il est par exemple intéressant de remarquer que $\mu(\text{AL})$ oscille comme le défaut et ce phénomène peut être exploité par un dégradé de couleur sur l'interface pour l'opérateur.

La cinquième situation (cf. Figure 3-24) correspond au cas n°4 avec la méthode Inter. et le défaut additif : **130** alarmes sont dénombrées alors que la méthode BO donne 288 alarmes. La courbe bleu clair représente les enveloppes internes et externes³⁻²¹. La mesure²² est en bleu gris. Cette méthode est moins sensible que la méthode BO et en particulier dans la première partie (0 à 1000 secondes) lorsque l'amplitude du défaut est plus faible (cf. Figure 3-14). Par contre elle garantit le résultat (pas de fausses détections)²³.

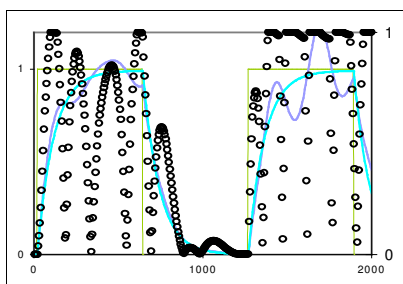


Figure 3-23 : Cas 3, BO Flou

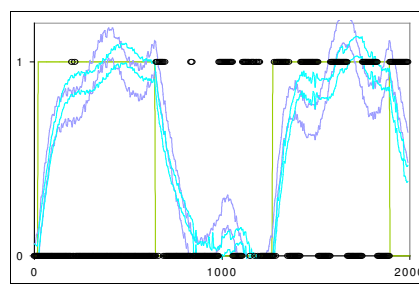


Figure 3-24 : Cas 4, méthode Int.

³⁻²¹ Ces deux enveloppes sont confondues dans le cas particulier des fonctions de transfert du premier ordre.

³⁻²² Notons que la mesure est représentée par une enveloppe.

²³ Sous hypothèse que les paramètres restent constants et que des intervalles les contenant sont bien connus.

3.5.3.5 Conclusion

Dans cette section, nous avons comparé sur des données simulées quatre techniques de détection. Chaque technique présente des inconvénients et des avantages en fonction du résultat souhaité. Nous aboutissons aux conclusions suivantes :

- La méthode BO est la plus performante (en terme d'alarmes) sur les données bruitées, par contre elle ne permet pas de diagnostiquer des systèmes instables.
- La méthode de référence BF donne de bon résultats, elle présente l'avantage (ou l'inconvénient) de faire une correction continue et peut être utilisée pour des systèmes instables.
- La méthode floue BO apporte des raffinements permettant de prendre en compte la persistance et la dynamique du défaut. Par contre ceci entraîne des retards à la détection ($\mu(AL)=1$). Elle présente l'avantage original de savoir si la situation est plus ou moins anormale contrairement à une décision binaire.
- La méthode Inter. présente l'avantage de garantir le résultat et permet de générer des bornes évolutives en fonction de l'excitation du système contrairement aux autres méthodes pour lesquelles les seuils sont fixes.

3.5.4 Complémentarité des techniques

Dans cette section, nous présentons une méthode de détection que nous avons élaborée et qui permet de combiner les avantages de l'approche ensembliste (les seuils sont fonctions de l'excitation de l'entrée) et les avantages de l'approche par logique floue (évolution graduelle de l'appartenance à la classe alarme).

Le noyau (Z) des résidus utilisé par les algorithmes flous (cf. Figure 3-7) n'a pas la même sémantique que l'enveloppe interne qui est générée par les intervalles. Une mesure située dans le noyau Z est considérée non anormale car elle est compatible avec l'erreur de mesure et l'erreur de modèle. Une mesure située hors du noyau Z est considérée (plus ou moins) anormale. On ne peut pas décider de l'état normal ou anormal d'une mesure située dans ou hors de l'enveloppe interne.

La notion d'enveloppe interne de l'approche par intervalles modaux n'a donc pas d'équivalent dans l'approche par logique floue. Par contre, l'enveloppe externe a une

sémantique plus proche des sémantiques (NN ou PP) utilisées par les algorithmes à base de logique floue. Une mesure située hors de l'enveloppe externe est assurément anormale. Une mesure située dans l'enveloppe externe peut être anormale ou normale. C'est cette notion que nous utilisons par la suite quand nous parlons d'enveloppe.

Nous proposons dans ce manuscrit de générer deux enveloppes externes puis d'étudier, via un raisonnement à base de logique floue, la position de la mesure par rapport à ces enveloppes. La première enveloppe externe définit les bornes du noyau Z et la seconde enveloppe définit les bornes du noyau PP et NN par symétrie.

La première enveloppe est générée à partir des écarts types sur les paramètres $[\theta \pm \sigma_\theta]$ et des incertitudes sur les mesures (cf. (3-28)). Elle permet de définir les bornes du noyau Z. Les bornes de cette enveloppe sont notées min1 et max1. La mesure est représentée par l'intervalle $[min-mes; max-mes]$. Si l'intersection entre cette enveloppe et la mesure est vide, alors la variable ne peut pas être considérée normale.

La seconde enveloppe est générée à partir des écarts types sur les paramètres $[\theta \pm 4.\sigma_\theta]$ et des incertitudes sur les mesures (cf. (3-27)). Elle permet de définir les bornes des noyaux PP et NN. Les bornes de cette enveloppe sont notées min2 et max2. Si l'intersection entre cette enveloppe et la mesure est nulle, alors on a une alarme (cf. Figure 3-5).

La Figure 3-25 illustre les enveloppes que nous avons définies et que nous allons utiliser dans la suite de cette section.

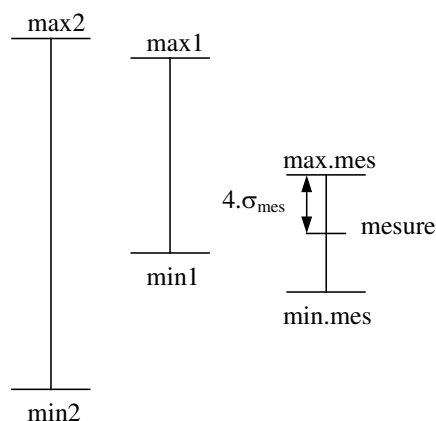


Figure 3-25 : Enveloppes utilisées pour la combinaison des techniques

Le principe de la méthode que nous avons élaborée consiste à comparer, via un raisonnement flou, la position de la mesure (qui est un intervalle) avec celles des enveloppes. Or, dans la méthode par logique floue de la section 3.4.3, nous avons comparé deux réels (une mesure et une référence) et pas deux intervalles.

Pour pouvoir comparer la position de ces intervalles, nous effectuons un changement de variables et définissons un nouveau **résidu**. Ce résidu est défini comme la différence entre le milieu de la mesure et le milieu de la seconde enveloppe :

$$résidu = \frac{max.mes + min.mes}{2} - \frac{max2 + min2}{2} \quad (3-30)$$

Nous définissons alors les nouveaux seuils des classes floues puis nous illustrons sur des exemples de valeurs du résidu :

$$SPA = 4.\sigma_{mes} + \frac{max1 - min1}{2}$$

$$SRA = \frac{1}{2}(SPA + SA)$$

$$SA = 4.\sigma_{mes} + \frac{max2 - min2}{2}$$

$$SPA^* = 2.SPA$$

$$SA^* = 2.SA$$

Les valeurs caractéristiques prises par ce résidu sont (cf. Figure 3-26) :

- **Cas n°1 : résidu=0** $\Rightarrow min.mes = \frac{max2 + min2}{2} - 4.\sigma_{mes}$

Ce qui signifie que la mesure est centrée dans la 2nd enveloppe.

- **Cas n°2 : résidu=SPA** $\Rightarrow min.mes = \frac{max2 + min2}{2} + \frac{max1 - min1}{2}$

Ce qui signifie que la mesure est juste au dessus d'une première enveloppe qui serait centrée dans la seconde enveloppe³⁻²⁴.

- **Cas n°3 : résidu=SA** $\Rightarrow min.mes = max2$

Ce qui signifie que la mesure est juste au dessus de la 2nd enveloppe.

³⁻²⁴ La première enveloppe n'est pas nécessairement centrée dans la seconde enveloppe.

La Figure 3-26 illustre les trois cas précédents. La mesure est en noir, la première enveloppe est en vert et la seconde enveloppe est en rouge.

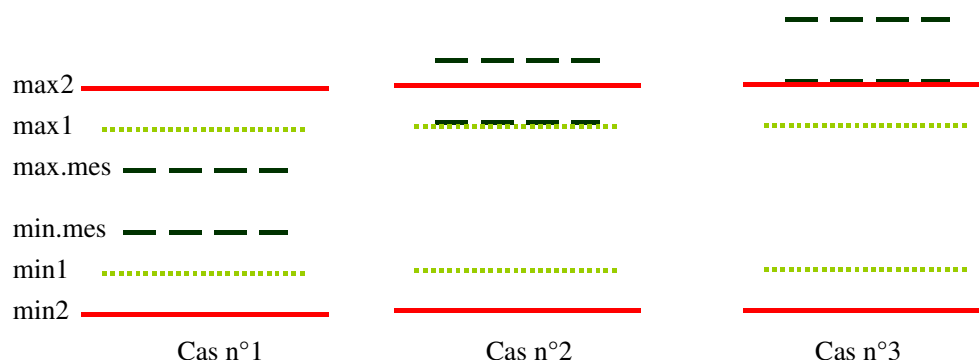


Figure 3-26 : Illustration des valeurs caractéristiques prises par le nouveau résidu

Les résultats que nous avons obtenus avec cette technique sont illustrés sur la Figure 3-27 et la Figure 3-28. La Figure 3-27 représente la première enveloppe (en rose) et la seconde enveloppe (en bleu clair). La Figure 3-28 représente la seconde enveloppe (en bleu clair) et la mesure (en bleu gris) et les appartenances à la classe alarme $\mu(\text{AL})$ (ronds noirs). 130 alarmes sont dénombrées avec la méthode Inter seule (cf. Figure 3-24) et avec la méthode qui vient d'être présentée. Cela est normal car la seconde enveloppe est exactement la même pour les deux méthodes. Par contre, cette méthode présente l'avantage de proposer une évolution graduelle de l'alarme grâce à l'interprétation par le raisonnement par logique floue des positions des deux enveloppes.

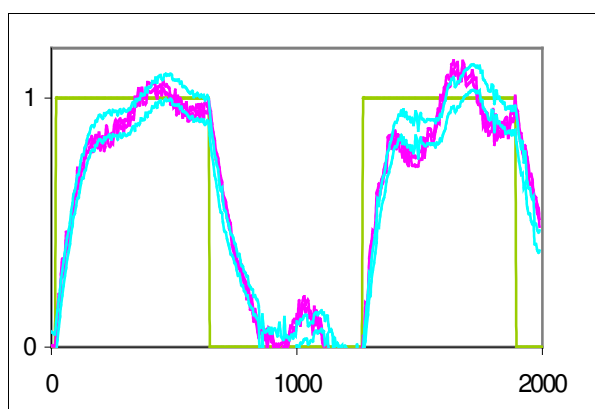


Figure 3-27 : Cas 4, méthode Int&Flou

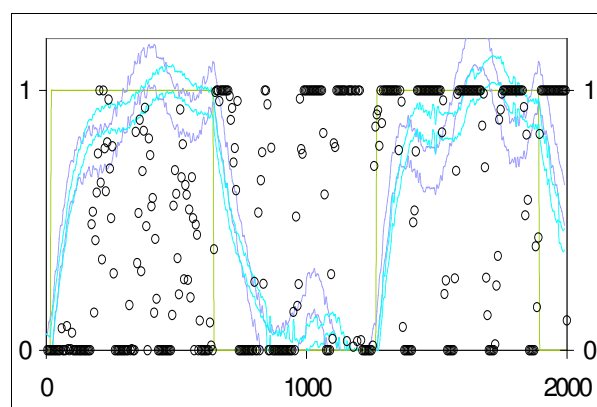


Figure 3-28 : Cas 4, méthode Int&Flou(2)

3.5.5 Conclusion

Dans cette section, nous avons présenté deux méthodes de détection issues de la littérature et qui permettent de prendre en compte les incertitudes du modèle et des mesures : une approche à base de logique floue et une approche ensembliste.

Nous les avons comparées en mettant en avant leurs avantages et leurs inconvénients. Nous avons finalement développé une méthode permettant de combiner ces deux techniques pour profiter des avantages de chacune et palier les inconvénients : elle présente l'avantage de garantir le résultat tout en proposant une évolution graduelle de la décision.

3.6 Localisation des défauts

3.6.1 Introduction

Les alarmes, décrites dans la section 3.3, permettent de déterminer quelles sont les variables qui ont un comportement anormal ou non vis-à-vis des variables exogènes ou vis-à-vis de leurs causes directes. Nous avons vu (cf. section 3.2.2) que l'alarme globale a un sens plus utile pour la conduite que pour le diagnostic. Contrairement aux alarmes globales, les supports qui sous-tendent les relations locales sont faibles et permettent donc d'obtenir des diagnostics concis.

Dans toute la démarche qui suit, seules les alarmes locales exprimant une incohérence entre une variable et ses causes directes sont utilisées pour la localisation. Les alarmes locales nous permettent de savoir quelle est la variable source qui explique toutes les déviations. Dans cette section nous cherchons le composant défaillant.

Les alarmes locales sont générées à l'aide de relations de redondance analytique : lorsqu'un composant qui sous tend une RRA se comporte anormalement, alors le résidu de cette relation peut être grand. En conséquence de toutes ces observations, nous proposons de raffiner le diagnostic en effectuant la localisation des défauts [Iserman et Ballé 1997].

Cette section présente une méthode permettant de localiser les composants défectueux : l'étape de localisation consiste à établir des diagnostics³⁻²⁵.

Nous nous sommes appuyés sur les travaux [Reiter 1987] et [Travé-Massuyès *et al.* 2001] pour générer les diagnostics à partir des conflits, selon un algorithme dit de «hitting-set». La théorie logique du diagnostic [Reiter 1987] repose sur le concept du conflit. Un conflit décrit un ensemble de composants dont un au moins doit être en panne pour que les observations soient consistantes avec le modèle. Les conflits sont calculés à partir des observations. A partir de cet ensemble de conflits, il est possible de calculer les diagnostics. Un diagnostic décrit un ensemble de composants dont le dysfonctionnement explique les observations. Dans ce cadre, la section 3.6.2 présente des stratégies de localisation.

3.6.2 Stratégies de diagnostic

3.6.2.1 Introduction

L'approche IA explicite le lien entre la relation et son support (cf. section 1.5.3.1) selon l'implication suivante : «si les composants du support ne sont pas anormaux alors les observations sont cohérentes avec la relation»³⁻²⁶. Un composant C se comportant anormalement est noté $AB(C)$ ³⁻²⁷ : il se différencie d'un composant ne se comportant pas anormalement qui est $\neg AB(C)$. Le système observé est représenté par son modèle causal approché (MCA), l'ensemble des composants (COMPS) et l'ensemble des variables connues (OBS). Il est détecté anormal si et seulement si l'ensemble $MCA \cup OBS \cup \{\neg AB(C) \mid C \in \text{COMPS}\}$ ne peut être satisfait.

Un diagnostic minimal est le plus petit ensemble de composants permettant d'expliquer les observations. Il est obtenu par interprétation des conflits.

[Travé-Massuyès *et al.* 2001] distinguent les «hard conflicts» et les «soft conflicts». Les «hard conflicts» sont générés à partir des variables dont le

³⁻²⁵ Rappelons qu'un diagnostic est un ensemble de composants dont toutes les entités sont suspectées d'avoir un comportement anormal.

³⁻²⁶ Ce prédicat suppose que le modèle représente parfaitement le composant en bon fonctionnement.

³⁻²⁷ $Ab(\text{Anormal})$, en Anglais.

comportement est anormal et les «soft conflicts» sont générés à partir des variables dont le comportement n'est pas anormal. Le comportement des variables est jugé selon la valeur d'une alarme.

Exemple : Cette section sera illustrée sur l'exemple de la Figure 3-29.

COMPS = {A1, A2, M1, M2, M3}, où A_i représente un additionneur et M_i un multiplicateur. Les variables sont a, b, c, d, e, f, g, t, u, v et w. Les valeurs mesurées de ces variables sont indiquées (**OBS**={a, b, c, d, e, f, g}). Les variables t, u et w ne sont pas mesurées. Notons qu'en bon fonctionnement, nous avons $f = 12$ et $g = 12$.

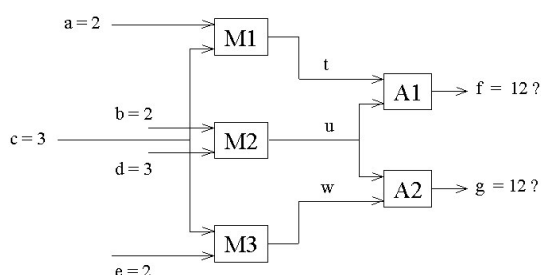


Figure 3-29 : Exemple de système

De cet exemple, se déduit directement le MCR.

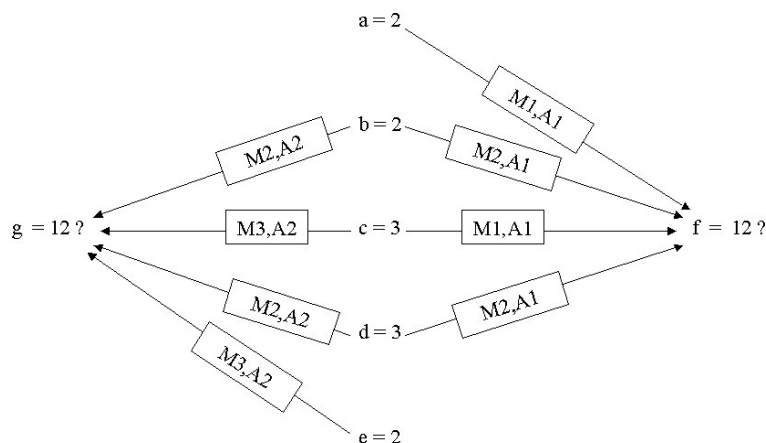


Figure 3-30 : Exemple après réduction

Remarque 3-6 : Comme toutes les relations sont statiques, alors d'autres représentations causales que celle de la Figure 3-30 sont envisageables. Les supports des relations de ces représentations sont différents et apportent une information

supplémentaire pour la localisation des défauts. Nous n'avons pas envisagé cette démarche mais elle pourrait constituer une perspective de ce manuscrit.

Génération des “hard conflicts”

Aucune exonération³⁻²⁸ n'est effectuée pour générer les «hard conflict». Considérons une variable V dont le comportement est anormal. Ce comportement est jugé anormal car une incohérence est observée entre la référence locale et la mesure (noté $LO(V)=1$). Un «hard conflict» est établi comme l'ensemble des composants physiques qui sous tendent ce modèle local :

$$LO(V)=1$$

$$\{\text{Hard Conflict}\} = \cup \text{Support}_{\text{Influence}} \left(\text{Var}_{V \rightarrow V, \text{causes.directes}}(V), V \right) \quad (3-31)$$

Exemple :

- Cas n°1 : Si f est localement anormal et g n'est pas localement anormal, alors le « hard conflict » généré à partir de f est $\{A1, M1, M2\}$.
- Cas n°2 : Si f est localement anormal et g est localement anormal, alors le « hard conflict » généré à partir de f est $\{A1, M1, M2\}$ et le « hard conflict » généré à partir de g est $\{A2, M2, M3\}$.

Génération des “soft conflicts”

Si l'hypothèse d'exonération des fautes simples est envisagée, alors il est possible de générer des «soft conflicts». Sous cette hypothèse, les variables n'ayant pas un comportement anormal apportent une information supplémentaire pour le diagnostic. Ces conflits sont obtenus dans certains cas particuliers à partir des variables dont le comportement n'a pas été jugé anormal. La suite de cette section illustre l'élaboration des «soft conflicts».

Soit Δ , l'intersection des supports des influences de deux variables V et Y :

$$\Delta = \text{Support}_{\text{Influence}} \left(\text{Var}_{V \rightarrow V, \text{causes.directes}}(V), V \right) \cap \text{Support}_{\text{Influence}} \left(\text{Var}_{V \rightarrow V, \text{causes.directes}}(Y), Y \right) \quad (3-32)$$

³⁻²⁸ L'exonération consiste à considérer qu'une faute se manifeste toujours.

Si l'alarme locale de V vaut 0 et l'alarme locale de Y vaut 1 et Δ contient un unique composant, alors la différence symétrique entre l'ensemble des antécédents de Y et l'ensemble des antécédents de V constitue un « soft conflict » :

$$\begin{aligned} \text{Si } LO(V) = 0 \text{ et } LO(Y) = 1 \text{ et } Card(\Delta) = 1 \\ \text{Alors } \{Soft\ Conflict\} = \bigcup_{X=V,Y} Support_{Influence} (Var_{V \rightarrow V, causes.directes}(X), X) \setminus \Delta \end{aligned} \quad (3-33)$$

Si l'exonération des fautes simples est adoptée, la génération des «soft conflicts» peut aussi être étendue au cas où une variable jugée non anormale possède un unique composant en commun avec plusieurs variables jugées anormales.

Dans le module que nous avons développé, nous ne générons pas de «soft conflict» qui correspondent à un cas très particulier qu'en pratique nous n'avons pas encore rencontré.

Exemple :

Dans le cas n°1, le «soft conflict» est l'ensemble {A1,A2,M1,M3}. En effet si M3 et A2 ont un comportement normal, alors M2 a nécessairement un comportement normal (sinon g serait anormal). L'hypothèse d'exonération d'une faute simple sur M2 permet de faire ce raisonnement. De plus si M1 et A1 ont un comportement normal, comme M2 est nécessairement normal alors f ne peut être anormal. L'ensemble {A1,A2,M1,M3} définit donc bien un conflit.

3.6.2.2 Génération des diagnostics

Les diagnostics sont obtenus par intersection des conflits. Deux stratégies sont ensuite présentées pour la génération des diagnostics.

Exemple :

- Cas n°1 : les conflits sont {A1,M1,M2} et {A1,A2,M1,M3}. Les diagnostics sont donc {A1},{M1},{A2,M2} et {M2,M3}. Dans ce cas, l'exonération des fautes simples est effectuée car le «soft conflict» {A1,A2,M1,M3} est considéré.
- Cas n°2 : les conflits sont {A1,M1,M2} et {A2,M2,M3}. Les diagnostics sont par conséquent {M2}, {A1,A2}, {A1,M3}, {A2,M1} et {M1,M3}.

Exonération et élimination de diagnostics

L'hypothèse d'exonération permet d'éliminer certains diagnostics. En effet, si une variable est jugée non anormale et si l'hypothèse d'exonération de fautes simples et de fautes multiples est considérée, alors les composants qui sous tendent le comportement de cette variable sont considérés non anormaux.

Exemple :

Dans le cas n°1, les diagnostics obtenus (avec exonération des fautes simples) sont $\{A1\}, \{M1\}, \{A2, M2\}$ et $\{M2, M3\}$. Comme g n'est pas anormal alors l'hypothèse d'exonération des fautes simples et multiples, entraîne à considérer que $A2$, $M2$ et $M3$ sont normaux. Par conséquent, les diagnostics retenus sont $\{A1\}$ et $\{M1\}$. En effet les diagnostics $\{A2, M2\}$ et $\{M2, M3\}$ consistent respectivement à considérer que $A2$ et $M2$ ou $M2$ et $M3$ sont anormaux. L'hypothèse d'exonération des fautes multiples³⁻²⁹ a permis d'éliminer les diagnostics $\{A2, M2\}$ et $\{M2, M3\}$.

Dans le cas n°2, l'hypothèse d'exonération n'apporte pas plus d'information car les deux variables f et g sont anormales.

3.6.2.3 Conclusion

Il est possible de développer des stratégies de localisation qui sont basées sur l'exonération ou la non exonération des fautes. Dans cette section, nous avons présenté deux stratégies et nous les avons illustrées dans le cas n°1. Aucune stratégie n'est présentée dans le cas n°2 car aucune exonération n'est possible.

- Cette hypothèse consiste à considérer que les fautes simples et multiples ne se manifestent pas toujours. Dans le cas n°1, seul³⁰ le «hard conflict» $\{A1, M1, M2\}$ est considéré. Par conséquent les diagnostics sont $\{A1\}, \{M1\}, \{M2\}$. Nous avons appliqué cette stratégie pour tous les défauts sauf les défauts capteurs qui, en pratique, sont traités selon la stratégie suivante.

³⁻²⁹ Une faute multiple correspond à plusieurs composants défailants en même temps.

³⁰ Le «soft conflit» $\{A1, A2, M1, M3\}$ ne peut être pris en compte car il fait intervenir l'hypothèse d'exonération des fautes simples.

- La seconde stratégie consiste à supposer que les fautes simples et multiples se manifestent toujours. Dans le cas n°1, les diagnostics obtenus sont donc {A1} et {M1}. En pratique, nous avons exonéré les fautes sur les capteurs : nous avons constaté sur le pilote de FCC que les défauts capteurs sont particulièrement brutaux et tous les résidus sont sensibles à ces défauts.

3.6.3 Conclusion

Dans cette section, nous avons présenté une méthode de localisation. Elle permet de localiser quels sont les composants susceptibles d'être défectueux. Nous proposons les perspectives suivantes d'amélioration :

- Nous avons décidé, dans un premier temps, de ne considérer que les alarmes locales pour établir les diagnostics car les supports des relations qui les génèrent sont de faible cardinal. Les alarmes globales apporteraient une information supplémentaire, en particulier concernant les défauts capteurs et une amélioration de la méthode consisterait à les prendre en compte pour la génération des conflits.
- Nous aurions pu utiliser des modèles de fautes : ils consistent à associer les relations non pas à des composants, mais à des défaillances particulières sur des composants. Cette démarche apporterait une amélioration à la méthodologie [Travé-Massuyès, Dague et Guerrin 1998].

En pratique, les diagnostics possibles peuvent comporter plusieurs composants. On utilise alors, dans la suite de ce chapitre, la connaissance de mauvais fonctionnement de ces composants pour que l'opérateur se focalise sur des composants particulièrement suspects.

3.7 Identification des défauts

3.7.1 Introduction

Nous avons choisi de nous appuyer sur un système à base de règles (cf. 1.5.4.2) pour effectuer l'identification de défauts et nous avons élaboré une méthode permettant de prendre en compte la connaissance préalablement obtenue lors de la détection et de la localisation de défauts. Nous avons associé à chaque composant une base de règles contenant la connaissance des experts sur le comportement anormal de ce composant. Lorsque le composant fait partie d'un diagnostic alors sa base de règles est activée.

La base de règles associe des symptômes à une défaillance particulière sur ce composant. Lorsque les symptômes sont observés selon la règle adéquate alors cette défaillance particulière est identifiée. Le résultat de l'identification se traduit par un message en langage naturel transmis à l'opérateur. Ce message décrit la défaillance sur un composant, les actions à entreprendre (en ligne ou pour la maintenance) et les conséquences de cette défaillance sur le fonctionnement du procédé.

Le principe de cette méthode, détaillée dans [Heim *et al.* 1999], est présenté dans la section 3.7.2.

3.7.2 Modèles qualitatifs de fonctionnement anormal

Dans un premier temps, nous associons chaque composant à une base de règles qui est complétée par l'expert selon le modèle du Tableau 3-5 :

- chaque ligne correspond à une défaillance sur un composant,
- la première colonne indique l'identité du composant (vanne, conduite, réacteur, etc.),
- la seconde colonne contient les symptômes sur les variables connues, liés entre eux par des liens logiques (ET et OU),
- la troisième colonne décrit la défaillance correspondante (fuite, bouchage, changement anormal des propriétés de fluides, etc.),

- la quatrième colonne décrit des observations qualitatives permettant de valider cette hypothèse (odeur, couleur, traces, etc.), qui n'ont pas pu être prises en compte par le système de détection local.
- la cinquième colonne indique les répercussions de cette défaillance sur le système (instabilité d'une variable, ou d'une partie du système, etc.),
- la sixième colonne présente les actions à entreprendre en ligne (changement des paramètres de régulation, purge, by pass, intervention sur l'unité, etc.).

Composant	Symptômes	Défaillance	Observations pour la validation	Répercussions	Actions à entreprendre
Vanne V108	Chute brutale de P108 et chute brutale de F108	Fuite externe	Présence de catalyseur sur V108	Mauvaise fluidisation	Remplacer V108 au prochain arrêt

Tableau 3-5 : Organisation de la connaissance experte

Classement des défauts

Les défauts sont classés selon la liste suivante [Busson *et al.* 1998] :

- **Défauts de capteurs.** Ils se caractérisent par un écart entre la valeur réelle de la grandeur et sa mesure. On classe ces défauts en fonction de leur type en : biais, dérive, modification du gain de mesure, valeurs aberrantes, blocage du capteur à une valeur constante ou coupure électrique du capteur.
- **Défauts d'actionneurs.** Ils se traduisent par une incohérence entre les commandes et la sortie (la pompe délivre un débit incohérent avec sa caractéristique hydraulique, vieillissement d'une vanne).
- **Défauts du processus physique.** Ces défaillances sont dues à des modifications de la structure (fuite interne, fuite externe, rupture d'un organe, etc.) ou des paramètres du modèle (encrassement d'un tube d'un four, bouchage partiel d'une conduite, bulles de gaz, etc.).
- **Défauts du système (ou de l'algorithme) de commande.** Ils se caractérisent par un écart entre la valeur réelle de la sortie du contrôleur (selon l'algorithme implémenté) et sa valeur théorique.

Un critère de classification consiste à distinguer les défauts internes et les défauts externes. La liste précédente a présenté des fautes internes. A cette liste nous ajoutons les fautes externes suivantes :

- Les **problèmes de services** de l'installation (manque de pression du réseau d'air, d'alimentation électrique, etc.).
- Les **erreurs de conduite** (bacs vides, consignes aberrantes, etc.).

Classement des symptômes

Nous avons établi la liste suivante de symptômes [Heim et al. 1999] pour les variables mesurées :

- **Tendances** (augmentation ↗, diminution ↘). Une simple moyenne glissante sur une fenêtre de grande taille est utilisée pour reconnaître ces symptômes.
- **Dépassement de seuils** (seuil haut >, seuil bas <). Ces symptômes sont obtenus par simple dépassement de seuils.
- **Formes** (à-coup ↗↘, oscillations ↗↘, palier ↗↘). L'obtention de ces symptômes est plus complexe que celle des précédents : elle nécessite de définir deux niveaux de formes intermédiaires et la décision sur le symptôme finalement retenu est retardée dans le temps³⁻³¹. La suite de cette section présente comment ces cinq derniers symptômes sont reconnus.

Un traitement simple des mesures permet aussi de reconnaître les formes de saut haut ↑, et de saut bas ↓, que nous qualifions de «formes primaires». Puis une interprétation temporelle et quantitative des formes primaires va nous permettre de définir les formes «secondaires» qui sont encore des formes intermédiaires. Finalement les formes finales, qui correspondent aux symptômes retenus (↗↘, ↗↘, ↗↘, ↗↘, ↗↘), sont obtenues, avec du retard, par interprétation des formes secondaires.

³⁻³¹ Il est difficile de distinguer un début de pompage et un à coup. Si la décision est prise avec du retard, alors cette distinction peut être effectuée car un pompage est persistant dans le temps alors qu'un à coup est éphémère.

- Formes primaires : saut haut et saut bas.
Nous définissons un saut haut \uparrow par une augmentation brutale de la mesure dans un temps «court» et un saut bas \downarrow par une diminution brutale de la mesure dans un temps court. Une moyenne glissante calculée sur une fenêtre de taille courte est comparée à un seuil pour détecter ces sauts.
- Formes secondaires : à-coup et pompage rapide.
Nous définissons un à-coup comme l'enchaînement dans un intervalle de temps suffisamment court des trois formes primaires $\uparrow, \downarrow, \uparrow$ ou des trois formes primaires $\downarrow, \uparrow, \downarrow$.³²
Nous définissons un pompage rapide comme une succession de plusieurs sauts opposés (au moins six) dans un intervalle de temps «suffisamment court».

Nous définissons les formes finales à partir des formes secondaires. Les formes finales sont obtenues chronologiquement plus tard que les formes secondaires car elles nécessitent de savoir comment la forme secondaire évolue dans le temps.

- Forme finale : à-coup, palier.
Nous définissons un palier haut $\uparrow\rightarrow$ comme un saut haut tel que la différence entre la moyenne après ce saut et la moyenne avant ce saut est supérieure à un seuil. Le palier bas $\downarrow\rightarrow$ est défini de manière symétrique.

Le symptôme à-coup $\uparrow\rightarrow$ ou $\downarrow\rightarrow$ est reconnu si une forme secondaire à-coup est observée et si aucun saut n'est observé suite à l'à-coup.

Si la forme secondaire est un pompage rapide alors la forme finale est aussi un pompage rapide $\uparrow\rightarrow$.

Il est bien entendu que cette technique de détection nécessite de fixer un nombre important de paramètres. Cependant, pour l'application au pilote de FCC, nous avons classé les variables par groupes (les pressions, les débits, les niveaux, etc.) et

³² Nous distinguons un à-coup commençant par un saut bas et un à-coup commençant par un saut haut. Ils sont, en général, liés à des défaillances de natures différentes.

nous avons fixé les mêmes paramètres pour les variables d'un même groupe. Nous avons préalablement remarqué que les variables d'un groupe ont sensiblement les mêmes symptômes. Tous les seuils qui permettent de détecter les symptômes sont fixés par les opérateurs. Une amélioration de la technique consisterait à fixer ces seuils de manière systématique.

Actions à entreprendre

Les actions à entreprendre dépendent en général des défaillances. Nous proposons dans la liste suivante des actions correctrices permettant de conduire le procédé en mode dégradé :

- changement des consignes du procédé,
- changement des paramètres des régulateurs,
- passage de la régulation en mode manuel,
- by pass d'un composant,
- en cas de bouchage certaines procédures permettent de dé-colmater en ligne, etc.

Implémentation

Nous avons implémenté les bases de règles sous la forme de graphes experts : les symptômes, les liens logiques et les dysfonctionnements sont représentés par des objets liés entre eux par des arcs selon la Figure 3-31. Le graphe expert d'un composant constitue un modèle qualitatif de mauvais fonctionnement qui est activé si et seulement si le composant est membre d'un diagnostic. Ce traitement préalable permet d'améliorer considérablement l'utilisation des graphes experts : des défauts différents ayant souvent des symptômes identiques, ce traitement préalable permet de focaliser l'étude d'un composant en particulier.

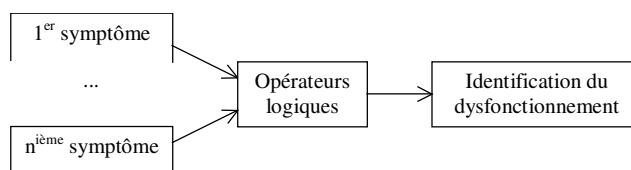


Figure 3-31 : Illustration d'un graphe expert

Un traitement du signal permet de reconnaître les symptômes à partir des observations. Lorsque les symptômes sont observés avec le respect des liens logiques, un message est envoyé sur l'interface de l'opérateur. Ce message décrit :

- la défaillance sur un composant,
- les confirmations (observations qualitatives pour la validation),
- les répercussions,
- les actions à entreprendre pour maintenir l'opération ou faire la maintenance de l'installation.

3.7.3 Conclusion

Dans cette section, nous avons décrit la méthode d'identification de défauts que nous avons développée. Nous utilisons des bases de règles que nous associons à chaque composant et qui sont activées si le composant est membre d'un diagnostic. Cette méthode présente l'avantage de prendre en compte les informations issues de la localisation lors de l'identification. Cela autorise des bases de règles très petites qui sont activées très rapidement en temps réel. Par contre, cela nécessite un travail d'extraction de l'expertise des opérateurs sur le mauvais fonctionnement des divers composants, ce qui peut se révéler assez long sur une installation complexe. Elle peut cependant être améliorée en extrayant des généralités des bases de règles qui ont été établies de manière purement experte.

3.8 Conclusion

Ce chapitre a présenté les trois méthodes que nous avons décidé d'utiliser successivement pour le diagnostic. Chaque méthode apporte individuellement une information utilisable pour la conduite (interface, graphe causal, liste des composants défectueux, messages) :

- La première méthode permet de générer des alarmes. Elle s'appuie sur un modèle quantitatif de bon fonctionnement. Ce modèle prend la forme explicative d'un modèle causal. Une comparaison entre la sortie du modèle et les mesures permet de générer les alarmes. Pour chaque variable, une alarme globale permet de savoir si la variable est cohérente avec les consignes du procédé et une alarme locale permet de savoir si la variable est cohérente vis-à-vis de ses causes directes. L'apport du raisonnement par logique floue et du

calcul ensembliste est évalué pour la prise en compte des incertitudes. Une comparaison avec une approche statistique simple permet cette étude. Ces deux techniques proposent des paramètres permettant de s'adapter à la dynamique des défaillances et à la persistance des défauts. Le raisonnement par logique floue étant graduel, il permet d'anticiper les défaillances en utilisant des codes de couleurs caractérisant la présence ou l'absence de défauts. Le calcul ensembliste est pessimiste, la détection est alors peu sensible, mais le résultat est garanti.

- La seconde méthode permet la localisation de la défaillance sur des composants. Les alarmes locales précédemment générées et la connaissance qualitative sous-tendue par le modèle causal (composants associés aux influences) sont interprétées par un algorithme de «hitting-set» pour la génération de diagnostics. Nous avons considéré que les fautes simples et multiples de tous les composants (sauf les capteurs) ne se manifestent pas toujours. Une étude de sensibilité des résidus aux défauts permettrait d'exonérer certaines fautes sur certains composants, mais cette étude nécessite une connaissance très approfondie des dysfonctionnements du procédé. Nous avons exonéré les fautes sur les capteurs car elles sont brutales dans le cadre du pilote de FCC.
- La troisième méthode permet d'identifier la défaillance. Des bases de règles, associées à chaque composant suspecté anormal, sont activées. Elles permettent de générer des messages identifiant la panne sur un composant. L'opérateur, lisant ces messages, est prévenu et le message lui indique les actions à entreprendre pour valider les hypothèses ou maintenir l'opération ; il peut alors prendre ses décisions en fonction de ces informations. Le module de diagnostic n'est donc pas bouclé car il n'agit pas sur le système, ce qui constitue une spécification très raisonnable.

Le chapitre suivant présente l'application de cette méthode à un pilote de FCC.

Chapitre 4

Mise en œuvre sur un pilote de FCC

Ce chapitre présente le procédé de craquage catalytique puis décrit le fonctionnement de l'unité pilote du CEDI⁴¹. Comme nous envisageons d'appliquer la méthode dans un cadre plus général, nous décrivons aussi le fonctionnement des unités industrielles de craquage catalytique. Afin d'évaluer les différences entre ces unités aux échelles différentes, nous avons établi une liste des dysfonctionnements classiquement rencontrés dans les deux cas. Nous verrons qu'ils sont sensiblement différents. Plus précisément, nous détaillons, dans les sections suivantes :

- le modèle causal du pilote de FCC,
- les résultats de diagnostics obtenus sur le pilote de FCC,
- les problèmes rencontrés et les apports de nos travaux dans le domaine du diagnostic de ce type de procédé,
- le contexte du développement du module informatique dans le projet Européen CHEM,
- les résultats escomptés de la méthode sur un FCC industriel,
- l'application de la méthode à d'autres procédés.

⁴¹ Institut français du pétrole (IFP) : Etablissement de Solaize – Centre d'Etude et de Développement Industriel - René Navarre

4.1 Le procédé de craquage catalytique

4.1.1 Introduction

Un des problèmes majeurs de l'industrie pétrolière est de convertir les produits lourds en produits légers et facilement valorisables comme les essences, les kérosènes, les gazoles utilisables pour la fabrication des carburants et les oléfines utilisées en pétrochimie. L'objectif des procédés de conversion est de craquer les molécules d'hydrocarbures lourds, c'est-à-dire de rompre les liaisons Carbone-Carbone dans les chaînes hydrocarbonées, pour diminuer la taille des molécules et donc transformer une charge lourde en produits légers.

Initialement (de 1915 à 1936), cette conversion des produits pétroliers lourds/légers était réalisée par un craquage thermique. Puis, ce mode de conversion a été remplacé par celui du craquage catalytique - FCC (Fluid Catalytic Cracking) - qui garantit à la fois la qualité des produits obtenus après traitement et un très bon rendement.

Les procédés en lit fluidisé représentent actuellement l'essentiel de la capacité de craquage catalytique dans le monde. Dans ces procédés, le catalyseur est une fine poudre fluidisée par un courant de vapeur, d'air ou d'azote qui permet au mélange gaz/solide de circuler comme un fluide entre le régénérateur et le réacteur (riser). La production de coke⁴⁻² inhérente au procédé exige une régénération en continu du catalyseur, d'où la nécessité de développer une technologie mettant en œuvre la circulation du catalyseur entre la zone réactionnelle et les capacités (réservoirs). D'un point de vue fonctionnel, le principe de fonctionnement d'un FCC est représenté par le diagramme de la Figure 4-1.

⁴⁻² Coke : produits lourds présents dans la charge ou issus de réactions secondaires de polymérisation lors du craquage.

Le catalyseur circule dans :

- Le riser dans lequel les réactions chimiques de craquage de la charge au contact du catalyseur se produisent à des températures de l'ordre de 500°C et à basse pression (2 à 3 Bars). Le coke se dépose alors sur le catalyseur ce qui arrête son activité catalytique.
- Deux régénérateurs, en série, dans lesquels le coke déposé sur le catalyseur brûle grâce à l'injection d'air de combustion. La chaleur ainsi générée permet le réchauffage du catalyseur jusqu'à 700-750°C le plus souvent. Le catalyseur régénéré, débarrassé du coke, peut donc être réutilisé dans le réacteur.
- Des liaisons réacteur régénérateur qui permettent la circulation du catalyseur d'une capacité à l'autre. La chaleur récupérée par le catalyseur dans les régénérateurs est en quelques sortes restituée au réacteur.

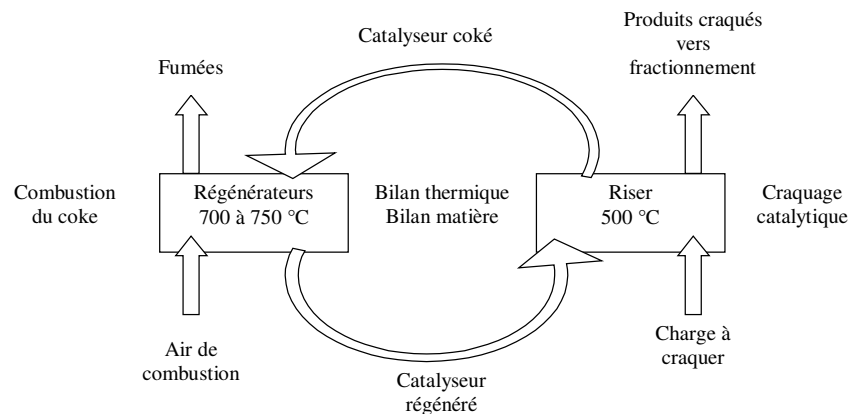


Figure 4-1 : Principe du craquage avec circulation de catalyseur

4.1.2 Unités industrielles de FCC

En moyenne, un procédé industriel de FCC permet de traiter 6500 m³/j (40000 barils/j). Le niveau de catalyseur dans les régénérateurs est de l'ordre de 15 m et le riser est de l'ordre de 35 mètres. Le débit de catalyseur est de l'ordre de 25 tonnes/min. Les procédés industriels sont parfaitement adiabatiques et peuvent fonctionner trois ans sans arrêt.

4.1.3 Le pilote de FCC du CEDI

L'unité pilote de FCC de l'IFP Solaize est un modèle réduit d'une unité industrielle de FCC. Elle a été construite pour développer le procédé de craquage catalytique en se rapprochant du comportement des unités industrielles. Elle permet d'envisager divers modes de fonctionnement, des séries de tests de charges, tests de catalyseurs qui permettent d'envisager des améliorations transposables aux unités industrielles. La durée d'un cycle de catalyseur est proche d'une heure. Le pilote de FCC peut traiter deux barils par jour. Les lignes de circulation du catalyseur ont un diamètre de l'ordre de 2 centimètres. Les capacités ont un diamètre de l'ordre de 20 centimètres. Le niveau de catalyseur dans les capacités est proche d'un mètre.

Nous proposons dans le Tableau 4-1 des ordres de grandeur des dimensions d'une unité industrielle et de l'unité pilote de FCC.

Caractéristiques	Unité industrielle	Unité pilote
Capacité	40000 barils/jour	2 barils/jour
Hauteur de chaque régénérateur	15 m	1 m
Diamètre des régénérateurs	8 m	20 cm
Hauteur du riser	35 m	7 m
Diamètre du riser	1 m	2 cm
Débit de charge	245 t/h	6 kg/h
Débit de catalyseur	1500 t/h	40 kg/h
Inventaire de catalyseur	300 tonnes	40 kg
Temps de séjour du catalyseur riser	2 à 4 sec	1 sec
Temps de séjour du catalyseur dans le stripper	1 min	15 min
Temps de séjour du catalyseur dans chaque régénérateur	5 min	20 min

Tableau 4-1 : Comparaison quantitative d'une unité pilote et industrielle de FCC

La Figure 4-2 présente un schéma simplifié du pilote de FCC. Le catalyseur circule dans une boucle physique : il circule depuis le stripper (R_3), vers le premier régénérateur (R_1) puis le second régénérateur (R_2) et finalement retourne vers le stripper (R_3). Le transport du catalyseur d'une capacité à l'autre est assuré par des tuyaux verticaux ou inclinés : le lift (T_2), le stand pipe (T_3), la jambe du riser (T_4) et le riser (T_1). Le riser constitue la zone réactionnelle. La charge est mise en contact avec le catalyseur chaud en bas du riser au niveau de la vanne (V_6), réagit et traverse le riser en quelques secondes. Les produits de la réaction et le catalyseur sont entraînés dans le stripper (R_3). Le stripper sépare les produits pouvant être valorisés, du catalyseur et du coke. Les produits pouvant être valorisés sont entraînés puis séparés dans la colonne de séparation (C_1). Le coker (F_3), situé entre le stripper et la colonne, est utilisé pour prévenir les bouchages par de grosses molécules entraînées depuis le stripper. Le catalyseur et le coke sont entraînés, par gravité, via la vanne de circulation de catalyseur V_4 vers le premier régénérateur R_1 . Le coke est partiellement brûlé dans le premier régénérateur. Le second régénérateur (R_2) achève la combustion (la présence de deux régénérateurs permet d'effectuer la combustion à une température plus faible que s'il n'y avait qu'un seul régénérateur). Le catalyseur chaud, sortant du second régénérateur est entraîné jusqu'au riser (T_1) et réagit immédiatement avec la charge. Les deux filtres F_1 et F_2 empêchent le catalyseur d'aller en aval des lignes à la sortie des régénérateurs. Les trois vannes V_1 , V_2 et V_3 permettent de réguler respectivement la pression du premier régénérateur, du second régénérateur et du stripper. Des débits d'azote dans les conduites et les capacités sont assurés par les vannes V_7 à V_{12} permettent d'entraîner et de fluidiser le catalyseur. Les deux vannes V_{13} et V_{14} permettent d'assurer des débits d'air pour la combustion du coke dans les régénérateurs. La vanne V_4 permet de réguler le niveau de catalyseur dans le stripper et la vanne V_5 , le niveau des essences et des produits lourds dans la colonne. Contrairement au pilote de FCC du CEDI, le FCC industriel contient une vanne située sous chaque régénérateur.

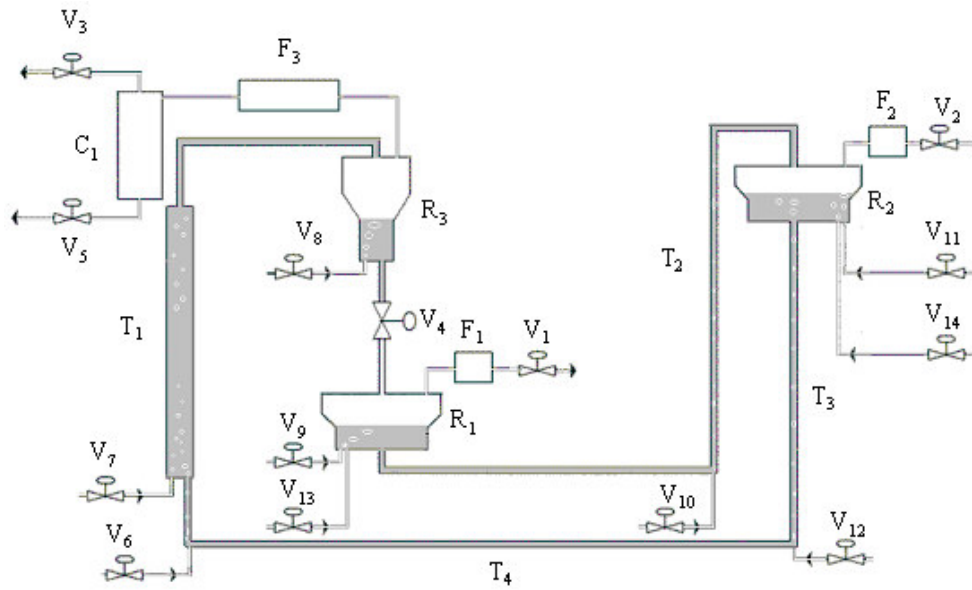


Figure 4-2 : Schéma simplifié d'une unité de FCC

La Figure 4-3 illustre une photographie du pilote.



Figure 4-3 : Photographie du pilote

4.1.4 Conclusion

Le pilote de FCC est un modèle réduit d'une unité industrielle : il fonctionne sur le même principe mais son objectif n'est pas de produire de l'essence mais d'étudier différents modes de fonctionnement. Ces modes de fonctionnement sont transposables aux unités industrielles et permettent d'en optimiser le fonctionnement.

4.2 Les dysfonctionnements du procédé de FCC

4.2.1 Introduction

L'automatisation des processus, les systèmes de supervision et la mise en place des calculs en ligne de bilans (matière et thermique) facilitent l'observation des unités de FCC. Néanmoins leur conduite reste difficile :

- Les informations sont mises à disposition des opérateurs en grande quantité et ils doivent les interpréter rapidement afin de comprendre ce qui se passe réellement. Cette interprétation et les décisions qui en découlent, demandent une connaissance approfondie du procédé.
- Les défauts se propagent rapidement dans la boucle physique du FCC, ce qui rend la localisation du défaut source extrêmement difficile. A fortiori, les premières alarmes (à seuils classiques) déclenchées ne situent souvent pas la défaillance.
- Des régulations complexes, à plusieurs cascades de régulateurs, sont appliquées sur le pilote de FCC. Ces régulations entraînent des relations de causalité supplémentaires qui complexifient l'interprétation du comportement.

Les problèmes rencontrés sur les unités industrielles de FCC, présentés dans la section 4.2.2 sont différents de ceux rencontrés sur les pilotes de FCC présentés dans la section 4.2.3. La petite taille du pilote de FCC est source de problèmes particuliers.

4.2.2 Problèmes rencontrés sur les unités industrielles de FCC

Cette section énumère les défauts fréquemment rencontrés sur les unités industrielles de FCC qui nous ont été communiqués par un expert de FCC.

Défauts du processus physique

Les défauts fréquents du processus physique d'un FCC industriel sont :

- inversion de la circulation du catalyseur et entraînement de charge dans le second régénérateur,
- entraînement d'air depuis le second régénérateur vers le riser,
- post combustion (after-burning) dans le régénérateur. Quand la production de coke diminue, le débit d'air est trop important par rapport à la quantité de coke à brûler. L'oxygène en surplus développe alors des réactions exothermiques dans la partie supérieure de la capacité se traduisant par une élévation rapide de la température des fumées,
- départ en coke. Lorsque la quantité de coke produite au riser devient plus grande que celle brûlée aux régénérateurs, il y a accumulation de coke dans l'unité,
- niveau haut de catalyseur dans le stripper induisant une perte de catalyseur vers la colonne de fractionnement,
- blocage des vannes de circulation de catalyseur,
- bouchage d'un cyclone (les cyclones sont des séparateurs rapides dans les capacités),
- perte des aérations du puits de soutirage,
- perte de la circulation de catalyseur par dé-fluidisation ou formation de bulles d'air rendant la circulation erratique,
- formation de coke dans le bas du riser par perte de la vapeur d'atomisation dans un injecteur,
- chute de coke accumulé en haut du riser bouchant le stripper,
- blocage des vannes de régulation de pression,
- bouchage par du catalyseur des générateurs de vapeur sur le circuit en fond de tour,
- dépôts de catalyseur sur les plateaux au-dessus de l'alimentation,
- blocage par du coke sur la ligne entre le réacteur et le fractionnement.

Problèmes de services

Les problèmes de services fréquents d'un FCC industriel sont :

- perte du débit d'air des régénérateurs induisant l'arrêt de la circulation et le remplissage de la ligne d'alimentation,
- perte de niveau d'eau dans les ballons de chaudière.

4.2.3 Problèmes rencontrés sur le pilote de FCC du CEDI

Nous avons consulté le cahier de l'unité pilote de FCC et en avons extrait les problèmes les plus fréquemment rencontrés sur le pilote de FCC. Ces problèmes sont majoritairement liés à la petite taille des tuyaux.

Défauts du processus physique

Les défauts fréquents du processus physique du pilote de FCC sont :

- mauvaise fluidisation du catalyseur dans les lignes (vidage du stand pipe ou de la jambe du lift),
- formation d'agglomérats de catalyseur dans les lignes,
- bouchage des vannes de régulation de pression,
- bouchage par du catalyseur des lignes de gaz en aval des capacités (filtres, tuyaux, vannes),
- bouchage des piquages des prises de pression,
- vieillissement du matériel (vanne de circulation de catalyseur).

Problèmes de services

Les problèmes de services fréquents du pilote de FCC sont :

- manque d'énergie électrique,
- manque d'air réseau.

Erreurs de conduite

Nous avons recensé une erreur de conduite sur le pilote de FCC :

- cavitation de la pompe de charge car la température des lignes d'alimentation est trop élevée.

4.2.4 Conclusion

Le FCC est un procédé particulièrement difficile à conduire car il contient une boucle physique et de nombreuses boucles de régulation imbriquées.

Les tests de catalyseur et de charge effectués sur le pilote de FCC sont transposables aux unités industrielles, par contre, les défaillances se produisant sur le pilote de FCC sont de nature différente de celles qui se produisent sur les unités industrielles, en particulier, les bouchages sont très fréquents à cause de la dimension des tuyaux.

4.3 Application au pilote de FCC du CEDI

4.3.1 Introduction

La méthodologie que nous avons développée et présentée dans ce document a été appliquée au pilote de FCC du CEDI. Cette section présente :

- le modèle causal approché du pilote de FCC,
- le positionnement de ce modèle par rapport aux autres modèles existants,
- les résultats obtenus de l'application du module de diagnostic sur le pilote de FCC.

4.3.2 Le modèle causal du pilote de FCC

4.3.2.1 Introduction

Après une approche pragmatique développée dans [Heim, Cauvin et Gentil 2000a], le modèle causal approché du pilote de FCC du CEDI a été construit selon la méthode présentée dans le Chapitre 2. Cette section présente les différents graphes causaux que nous avons construits. [Heim *et al.* 2002a] détaillent cette démarche.

4.3.2.2 Modélisation causale du pilote de FCC

Nous avons effectué la modélisation causale de l'ensemble du pilote de FCC. Il a tout d'abord été divisé en 25 composants qui sont présentés sur la Figure 4-2.

Plusieurs configurations sont envisageables sur le pilote de FCC selon l'agencement des cascades de régulation [Heim, Cauvin et Gentil 2000a] : le mode de fonctionnement de chaque régulateur doit être considéré.

Nous avons choisi de ne représenter que les boucles internes des cascades de régulation. La consigne de cette boucle étant fixée par l'opérateur ou calculée par un autre régulateur. Ainsi une seule configuration du pilote est envisagée.

Nous n'avons pas représenté les boucles externes des cascades de régulation. Cela nous aurait entraînés à envisager un très grand nombre de configurations et aurait seulement permis de détecter les défauts du système (ou de l'algorithme) de commande, défauts que nous ne rencontrons pas sur le pilote de FCC.

Nous avons sélectionné 323 variables, dont 83 mesurées, pour décrire le pilote de FCC. 40 sont exogènes et 283 sont endogènes. Nous avons écrit 274 relations entre ces variables. Il nous a donc été nécessaire de définir 9 variables pseudo exogènes [Heim *et al.* 2003b]. Nous avons ensuite effectué les deux opérations de réduction, puis d'approximation sur le MCS :

Nous avons ensuite effectué les deux opérations de réduction, puis d'approximation sur le MCS :

- le modèle causal réduit contient 86 variables dont 49 exogènes et 37 endogènes,
- le modèle causal approché contient 42 variables (13 exogènes, 2 pseudo exogènes et 27 endogènes).

Le modèle causal approché du pilote de FCC est présenté sur la Figure 4-5.

- Les variables pressions sont notées $P_{\text{composant}}$:

P_{R_1} , P_{R_2} , P_{R_3} et P_{C_1} sont respectivement la pression des gaz dans le premier et le second régénérateurs et R_3 et en tête de colonne.

- Les variables débits sont notées $F_{\text{composant}}$:

F_{V_1} et F_{V_2} sont respectivement les débits de fumées à la sortie de R_1 , R_2 et F_{V_3} est le débit de légers qui sortent en tête de colonne,

F_{V_6} est le débit de charge,

$F_{V_{7..12}}$ sont les débits d'azote, entrées du pilote de FCC,

$F_{V_{13}}$ et $F_{V_{14}}$ sont les débits d'air, entrées du pilote de FCC.

- Les variables pertes de charge sont notées $\Delta P_{\text{composant}}$:
 ΔP_{F_1} et ΔP_{F_2} sont les pertes dans les filtres à la sortie du premier et du second régénérateurs,
 ΔP_{F_3} est la perte de charge dans le coker,
 ΔP_{T_1} , ΔP_{T_2} , ΔP_{T_3} et ΔP_{T_4} sont respectivement les pertes de charge dans le riser, le lift, le stand pipe et la jambe du riser,
 ΔP_{V_4} est la perte de charge de la vanne de circulation de catalyseur.
- Les variables niveaux sont notés $L_{\text{composant}}$:
 L_1 , L_{R_2} et L_{R_3} sont les niveaux de catalyseur dans le premier et le second régénérateur et L_{R_3} est le niveau de catalyseur dans le stripper.
- Les variables ouvertures de vannes sont notés $OP_{\text{composant}}$.
- Les variables consignes sont notées $SP_{\text{variable_régulée}}$.

La Figure 4-4 positionne ces variables sur le schéma simplifié du pilote de FCC. La Figure 4-11 présente le MCA implémenté du pilote de FCC.

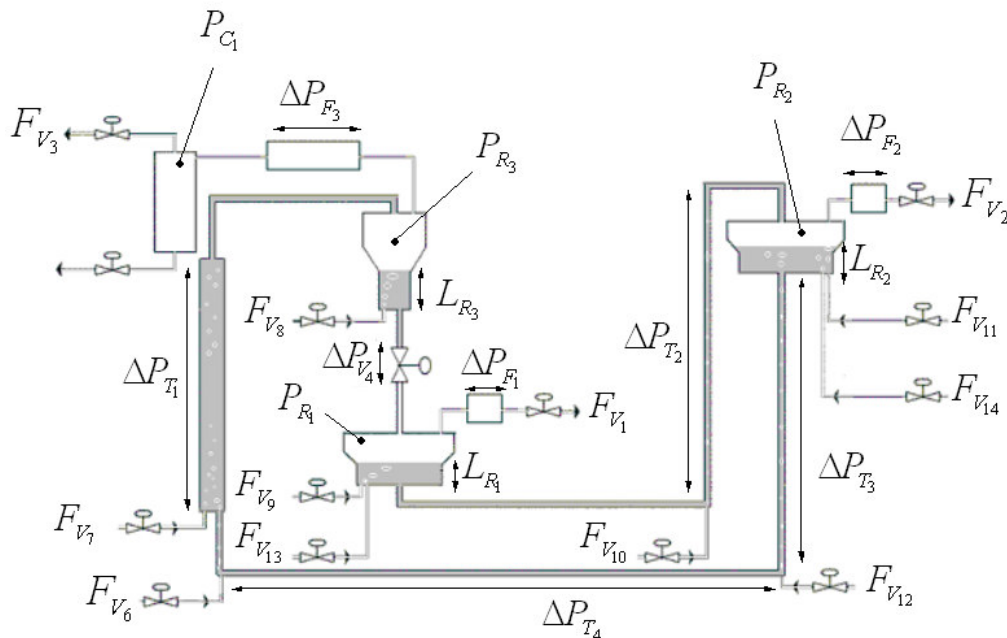


Figure 4-4 : Position des variables du pilote de FCC

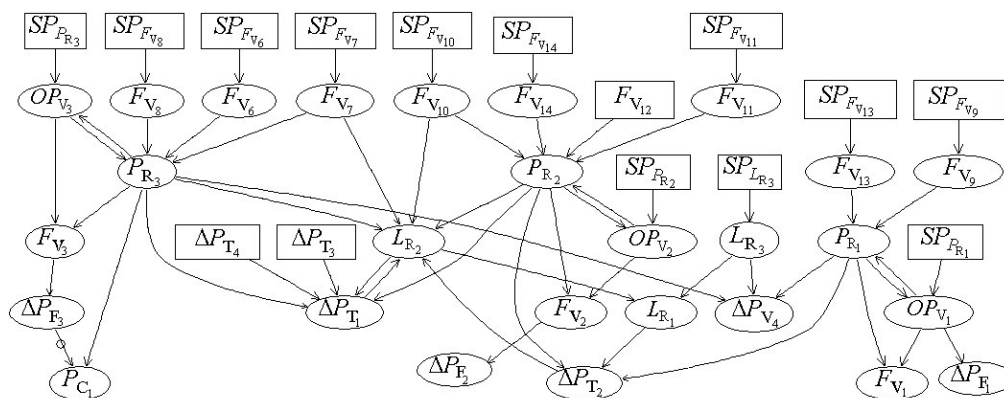


Figure 4-5 : Modèle causal approché du pilote de FCC

4.3.2.3 Positionnement du modèle développé par rapport aux modèles existants

De nombreux modèles de FCC sont disponibles dans la littérature. Parmi les derniers, citons les modèles de [Farlane *et al.* 1993], [Han et Chung 2001a] et [Han et Chung 2001b]. Tous les modèles existants s’attachent essentiellement à représenter des cinétiques plus ou moins complexes des phénomènes se déroulant dans le riser et dans les régénérateurs. Par contre, ils ne décrivent pas l’environnement complet du FCC (vannes d’alimentations, de régulation de pressions et de niveaux, colonne de séparation, filtres, etc.).

Le modèle cinétique s’appuie sur une description classique et simple des réactions chimiques [Heim *et al.* 2002a]. Le modèle décrit l’environnement complet du pilote de FCC, incluant les 25 composants précédemment cités.

Nous considérons que l’objectif d’un modèle dédié au diagnostic n’est pas de proposer une description fine des phénomènes, qui en pratique est difficile à paramétrer, mais plutôt de proposer une description complète permettant d’interpréter une grande quantité de données et de prendre en compte tout l’environnement du pilote.

4.3.2.4 Conclusion

Nous avons construit le modèle causal approché du pilote de FCC. Il est constitué de 42 variables. Pour obtenir automatiquement l'ordonnancement causal, nous avons utilisé le logiciel « Causalito »^{4.3} qui a été élaboré au LAAS^{4.4} [Travé-Massuès et Pons 1997]. Ce modèle a ensuite été implémenté pour le diagnostic.

4.3.3 Résultats obtenus sur le pilote de FCC

4.3.3.1 Introduction

Nous avons testé, sur le pilote de FCC, le module de diagnostic présenté dans ce manuscrit. Nous avons effectué ces tests sur des données réelles de mauvais fonctionnement archivées puis en ligne.

A l'aide de la base de données du pilote de FCC et de son cahier de vie, nous avons identifié 13 scénarios de mauvais fonctionnement. Nous n'avons pas pu créer ces scénarios car il ne nous a pas été permis d'intervenir sur le pilote de FCC.

Nous avons décrit tous ces scénarios dans [Promonet 2002] et cette section présente un scénario particulier pour illustrer la méthode. Dans [Heim *et al.* 2003c], nous décrivons précisément trois autres scénarios.

Le Tableau 4-2 récapitule les treize scénarios de défauts que nous avons recensés sur le pilote de FCC.

- La première colonne donne le numéro du scénario.
- La seconde colonne décrit la défaillance.
- La troisième colonne donne la durée totale du scénario.
- La quatrième colonne donne le temps de fonctionnement normal avant que la défaillance ne s'installe.
- La cinquième colonne indique le temps s'écoulant entre l'instant où la défaillance est installée et l'instant où elle est isolée par le système ASCO.

^{4.3} Les grands principes de fonctionnement de « Causalito » sont présentés dans la section 2.3.6

^{4.4} Laboratoire d'Analyse et d'Architecture des Systèmes, Toulouse, France.

N°	Scénarios	Durée totale	Durée de fonctionnement normal	Temps pour la localisation
1	Bouchage de la ligne entre le coker et la colonne de séparation	70 minutes	50 minutes mais perturbé	50 minutes
2	Bouchage de la vanne de régulation de la pression du premier régénérateur	25 minutes	15 minutes	5 minutes
3	Défaut capteur du niveau du 1 ^{er} régénérateur	50 minutes	40 minutes mais perturbé	5 minutes
4	Purges successives de la colonne par l'opérateur entraînant l'instabilité du procédé	35 minutes	25 minutes	instantané
5	Filtre du 1 ^{er} régénérateur colmaté	50 minutes	25 minutes	instantané
6	Filtres du 2 nd régénérateur colmatés	30 minutes	25 minutes	instantané
7	Manque d'air pour actionner les vannes d'air à l'entrée des deux régénérateurs	10 minutes	-	moins de 5 minutes
8	Bulles de gaz dans le circuit de pompage	70 minutes	50 minutes	instantané
9	Vidage du stand pipe	60 minutes	10 minutes	moins de 5 minutes
10	Entraînement de catalyseur entre le stripper et la colonne	120 minutes	20 minutes	60 minutes
11	Passage de gaz du premier régénérateur vers le second régénérateur	30 minutes	20 minutes	moins de 5 min
12	Comportement anormal de la vanne de régulation de niveau du stripper	50 minutes	-	instantané
13	Bouchage de la buse d'injection de charge	2 heures	-	1 heure

Tableau 4-2 : Tableau récapitulatif des scénarios de défauts du pilote de FCC

4.3.3.2 Description du premier scénario

Le scénario que nous avons choisi de décrire dans cette section correspond à un bouchage par du catalyseur de la ligne entre le coker F_3 et la colonne de séparation C_1 (cf. Figure 4-2). Le catalyseur est entraîné dans cette ligne par le flux de gaz sortant du stripper R_3 et allant vers la colonne. La Figure 4-6 illustre le support des influences qui entrent en jeu dans ce scénario (capteurs exclus). Par exemple, le composant RV3 (régulateur de la vanne V3) est dans le support des influences SPR3->OPV3 et PR3->OPV3.

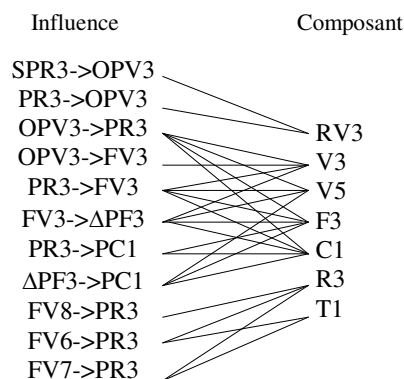


Figure 4-6 : Support des relations du scénario

4.3.3.3 Génération des alarmes

Apport du raisonnement par logique floue à la génération des alarmes

Nous avons implémenté les algorithmes de [Evsukoff 1998], puis nous les avons testés hors ligne et en ligne, sur les données du pilote de FCC, selon la méthode présentée dans la section 3.4. Un des principaux avantages du raisonnement par logique floue est qu'il permet d'obtenir une interface faisant apparaître de manière graduelle l'évolution des alarmes :

- Nous avons converti les appartenances à AL en code de couleurs sur les nœuds du MCA. Si l'appartenance à AL est 0 alors le nœud est vert et si l'appartenance à AL est 1, alors le nœud est rouge. Des couleurs de transition allant du vert au jaune puis à l'orange et au rouge sont utilisées pour traduire les appartenances à AL comprises entre 0 et 1.
- Nous avons aussi élaboré une interface qui présente l'historique des couleurs de chaque nœud selon la Figure 1-2. Cette interface permet de visualiser très rapidement l'historique du comportement de la variable.

Ces deux dernières interfaces ont été antérieurement proposées dans [Evsukoff 1998]. Elles permettent d'obtenir une vision massive des alarmes dans l'espace (via le graphe causal) et dans le temps (via l'interface des historiques). Elles permettent aussi d'anticiper les défauts.

L'algorithme de raisonnement par logique floue est plus robuste qu'une décision binaire car il permet d'interpréter simultanément les deux derniers résidus et leurs trois dernières variations. Nous avons observé quelques limites à l'application de la méthode.

- L'interprétation de la variation des résidus par le raisonnement par logique floue n'est pas nécessaire sur les scénarios du pilote de FCC.
- Lorsque les entrées du modèle sont très excitées, nous avons observé des fausses détections. Les paramètres (cf. Figure 3-6) qui définissent les classes floues devraient donc dépendre de l'excitation des entrées. Ils devraient augmenter lorsque le système est excité car l'incertitude sur la sortie du modèle augmente.

Apport des intervalles modaux à la génération d'alarmes

- Nous avons testé hors ligne les algorithmes de [Armengol 1999] mais nous ne les avons pas implémentés dans le module ASCO. Nous avons effectué ces tests dans le cadre du projet CHEM (cf. 4.3.5.2) et nous nous sommes appuyés sur une boîte à outils (TB3.2) réalisée par ailleurs.

Cette technique est donc difficile à mettre en œuvre. Malgré cela nous avons extrait des observations encourageantes.

- La méthode engendre des retards à la détection, par contre elle est plus robuste que la méthode par logique floue. En particulier, nous n'observons pas de fausse détection lorsque les entrées du système sont très excitées.
- Les défauts se déroulant sur le pilote de FCC sont de grande amplitude et la méthode a permis de les détecter.

Finalement, cette technique n'est pas encore assez mûre et les résultats obtenus reflètent l'état de l'art.

Nous avons testé la méthode combinant le flou et les intervalles sur les données simulées mais ne l'avons pas testée sur des données du pilote de FCC pour les raisons précédentes.

4.3.3.4 Interfaces

Les résidus locaux des variables PR3 et FV3 sont sensibles aux bouchages sur les composants C1 et F3. Les alarmes locales des variables se déclenchent donc au cours du scénario. La Figure 4-6 montre que les composants C1 et F3 sont dans le support d'influences sur les variables PR3 et sur FV3.

Le bouchage de lignes entre C1 et F3 déstabilise le pilote de FCC et un grand nombre de variables s'écartent de leur comportement attendu étant données les consignes et perturbations extérieures (référence globale). Les alarmes globales des variables L_{R2} , L_{R1} , P_{R3} , F_{V3} , ΔP_{F3} , P_{C1} , OP_{V3} , ΔP_{T1} , ΔP_{T2} , ΔP_{V4} se déclenchent.

La Figure 4-7 et la Figure 4-8 représentent respectivement l'évolution du résidu global et du résidu local de PR3 et les seuils sur ces résidus (cf. Remarque 3-6). La Figure 4-9 et la Figure 4-10 représentent respectivement l'évolution du résidu global et du résidu local de PR3. Nous remarquons que le défaut est brutal.

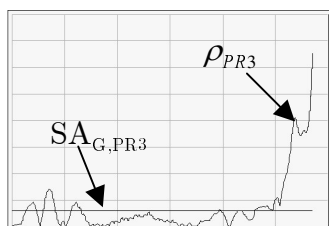


Figure 4-7 : Résidu global de PR3

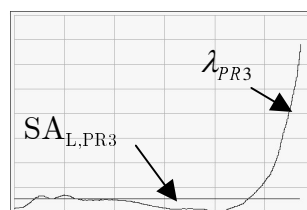


Figure 4-8 : Résidu local de ΔP_{V4}

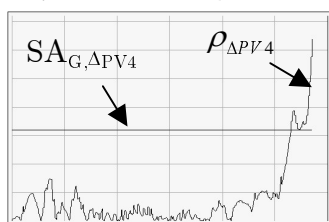


Figure 4-9 : Résidu global de ΔP_{V4}

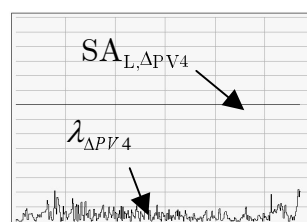


Figure 4-10 : Résidu local de ΔP_{V4}

Dans un tel scénario, les alarmes classiques à seuil de toutes ces variables sont susceptibles de se déclencher (en fonction du seuil choisi par l'opérateur), ce qui illustre la complexité du suivi en ligne d'un tel procédé.

La Figure 4-11 illustre le graphe causal du pilote de FCC au cours du scénario. Une variable est en gris si son alarme globale associée est déclenchée et les arcs qui l'influencent sont en gras si son alarme locale est déclenchée. On remarque que la méthode proposée permet bien de se focaliser sur le sous système en défaut.

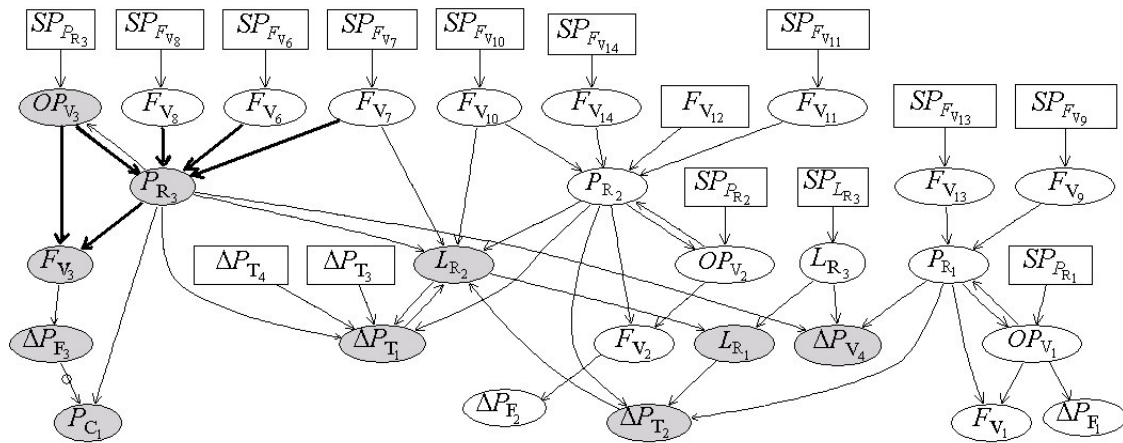


Figure 4-11 : Graphe causal du pilote de FCC : alarmes déclenchées

4.3.3.5 Localisation des défauts

L’alarme locale de PR3 est déclenchée, donc l’union des composants sur les supports des influences de PR3 constitue un premier conflit : $\{R3, V1, V3, V5, F3, C1, FV6^s, FV7^s, FV8^s, PR3^s\}$. L’alarme locale de FV3 est déclenchée, donc le second conflit $\{V3, V5, F3, C1, PR3^s\}$ est aussi généré. Les diagnostics possibles sont donc $\{V3\}, \{V5\}, \{F3\}, \{C1\}, \{PR3^s\}$ (fautes simples et non exonération).

L’expérience du pilote de FCC conduit à exonérer les fautes sur les capteurs. Le capteur PR3^s est sur le support de relations (par exemple influences sur LR3) dont l’alarme locale n’est pas déclenchée, par conséquent, et en exonérant, ce composant est supprimé des diagnostics.

4.3.3.6 Identification des défauts

Les graphes experts des composants $\{V3\}, \{V5\}, \{F3\}$ et $\{C1\}$ sont activés. Toutes les règles de ces graphes experts sont automatiquement analysées et seule la règle de $\{V3\}$ illustrée par la Figure 4-13 est vérifiée.

Au début du scénario, l'analyse du signal permet de détecter que P_{R3} et OP_{V3} augmentent. Puis à la fin du scénario, P_{R3} et OP_{V3} sont anormalement élevés. La Figure 4-12 illustre l'évolution de P_{R3} au cours du scénario.

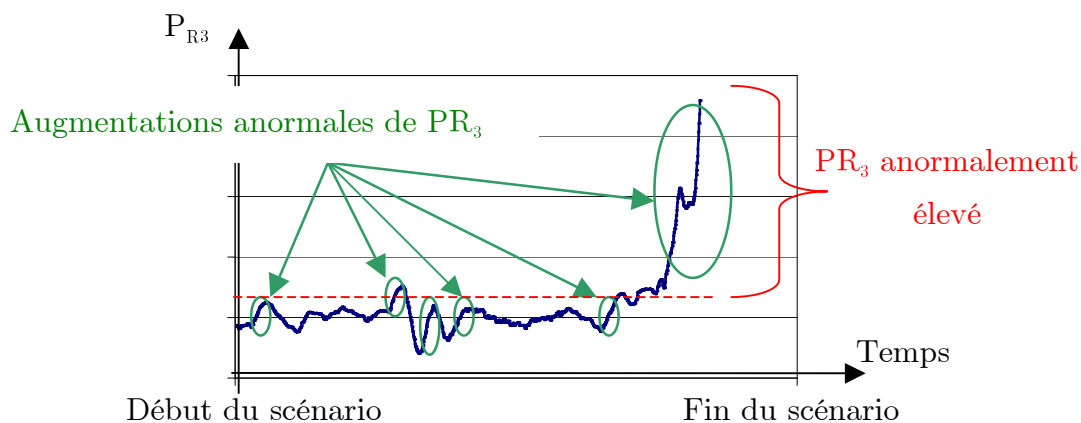


Figure 4-12 : Evolution de P_{R3} au cours du scénario

Les règles étant observées tout au long du scénario avec la logique souhaitée, le message qui est décrit dans la suite de cette section apparaît sur l'interface des opérateurs.

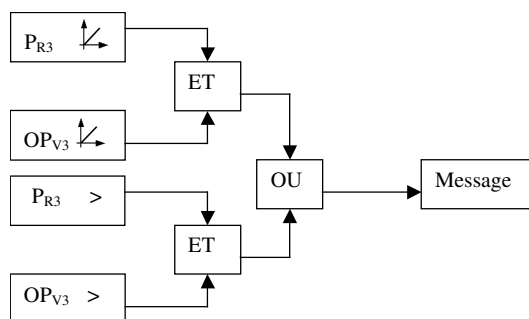


Figure 4-13 : Exemple de règle vérifiée lors du scénario

Le message envoyé à l'opérateur contient les informations suivantes:

- **Description de la défaillance** : bouchage de la vanne V3 ou du lavage à la soude (en aval de V3).
- **Actions à entreprendre pour vérifier cette hypothèse et affiner le diagnostic** : Court-circuiter le lavage à la soude pour déterminer l'origine du bouchage. Si l'ouverture OP_{V3} ne diminue pas alors le bouchage se trouve au niveau de V3.

- **Conséquences à long terme de cette défaillance** : déstabilisation du bilan pression de l'unité pilote. Risque d'arrêt brutal de l'unité.

Les autres règles contenues dans les autres composants ne sont pas vérifiées, parmi ces règles nous citons par exemple :

- Bouchage des filtres en tête de colonne C1 : cette hypothèse est rejetée car la température des filtres n'est pas suffisamment élevée,
- Bouchage du coker F3: cette hypothèse est rejetée car la perte de charge du coker n'est pas suffisamment élevée.

4.3.3.7 Conclusion

Cette section a présenté l'application du module de diagnostic que nous avons effectuée au pilote de FCC. Les trois étapes de génération d'alarmes, de localisation et d'identification apportent une information nécessaire et de plus en plus précise au fur et à mesure que l'amplitude du défaut augmente.

- Les prémices du défaut sont d'abord observés sur le graphe causal grâce aux codes de couleurs. Le graphe causal est donc une interface qui permet d'anticiper et de visualiser les défauts et de comprendre les phénomènes de propagation de défauts.
- Quand l'amplitude du défaut est suffisamment grande alors le module de localisation permet de suspecter une liste de composants. Cette information est plus précise que les alarmes mais est antérieure.
- Finalement l'étape d'identification permet d'identifier la panne sur un composant. Cela se traduit par un message qui est l'information la plus précise mais qui est transmise le plus tardivement à l'opérateur.

Nous avons observé, dans le cadre du pilote de FCC, qu'il est souvent possible d'effectuer des confirmations pour valider les hypothèses (odeur, couleur, etc.). Ces confirmations permettent généralement d'identifier de manière univoque le défaut. Par contre, il n'y a pas de reconfiguration possible sur le pilote de FCC et, par conséquent, il est fréquemment nécessaire d'arrêter l'unité suite au diagnostic. Dans une telle situation, le système ASCO aide la compréhension à posteriori des problèmes.

4.3.4 Problèmes rencontrés et apports des travaux au domaine

4.3.4.1 Introduction

Cette section présente les principaux problèmes que nous avons rencontrés au cours de l'élaboration et de l'application du module de diagnostic. Ces problèmes reflètent les avantages et les inconvénients des différentes méthodes envisagées. Ils sont aussi représentatifs de problèmes liés à l'application de telles techniques sur un procédé industriel. Cette étude *a posteriori* permet de mettre en avant les originalités de la méthodologie de diagnostic et les choix envisagés.

4.3.4.2 Organisation du module de diagnostic

Initialement [Heim *et al.* 1999], le module de diagnostic contenait uniquement le module d'identification et les règles étaient associées à des variables et non à des composants. Toutes les règles étaient étudiées simultanément sans aucun traitement préalable. L'inconvénient de cette approche que nous avons observé est que des défaillances avaient souvent les mêmes symptômes et qu'il était très difficile de les différencier. Nous avons alors envisagé deux voies d'amélioration :

- La première consistait à construire des bases de règles plus détaillées et permettant de bien différencier les défaillances. L'inconvénient de cette démarche est qu'il est nécessaire de fixer des seuils dans les bases de règles et qu'il est nécessaire d'avoir un très bon historique des défaillances pour pouvoir l'envisager.
- La seconde solution, qui a été choisie, a consisté à envisager un traitement préalable permettant de localiser la variable source des défauts. Le graphe expert d'une variable en défaut est activé si et seulement si l'alarme sur la variable est déclenchée.

Nous nous sommes ensuite appuyés sur les travaux de [Evsukoff 1998] pour effectuer ce traitement préalable. Ces travaux ont apporté plusieurs améliorations [Heim, Cauvin et Gentil 2000b]:

- les graphes causaux permettent de construire de manière systématique des résidus,
- la notion d'alarme locale a permis de localiser la variable source des défauts observés et d'activer les graphes experts associés à la variable,

- la notion d’alarme globale a apporté une information supplémentaire pour la conduite de l’unité . La plupart des méthodes de détection n’utilisent que la notion d’alarme locale. Pour illustrer le rôle important de l’alarme globale nous prenons l’exemple de l’alarme locale de la mesure de la température d’une serre qui vaudrait 0 alors que cette température serait de 80 °C.
- la logique floue, proposée dans [Evsukoff 1998] permet de prendre en compte les incertitudes du modèle et d’anticiper les défaillances.

A ce stade de développement, le module de diagnostic était plus performant mais les bases de règles n’étaient pas encore organisées de manière satisfaisante : il n’était pas judicieux de les associer aux variables⁴⁻⁵ et nous avons cherché une nouvelle méthode d’organisation. Parallèlement, nous avons étudié les travaux de [Travé-Massuyès *et al.* 2001] pour voir dans quelle mesure ils pouvaient apporter des améliorations à la méthode. Cette étude nous a permis d’élaborer la structure finale du module de diagnostic : nous avons associé les composants aux influences et nous avons organisé les graphes experts non plus en les associant aux variables mais aux composants. Nous avons apporté, grâce à cette démarche, une amélioration supplémentaire à la méthode [Heim *et al.* 2002b], [Heim *et al.* 2003a].

4.3.4.3 Modélisation causale

Nous avons élaboré une première méthode de génération de modèles causaux et nous l’avons appliquée sur le pilote de FCC [Heim, Cauvin et Gentil 2000b]. Cette méthode était basée sur une approche experte de recherche des phénomènes de causalité.

L’inconvénient de cette démarche est que les experts effectuent souvent des fermetures transitives (ils oublient des phénomènes intermédiaires pour décrire la causalité) et transgressent le principe de localité (les relations ne sont pas élémentaires et le support des relations est grand).

L’avantage de l’approche experte est qu’elle s’applique rapidement. L’inconvénient majeur est que la connaissance experte des opérateurs est souvent imprécise ou

⁴⁻⁵ Ce choix d’organisation avait été effectué par le technicien responsable de l’unité pilote de FCC qui attachait de l’importance aux variables clés qu’il suivait généralement en ligne.

confuse. Par exemple, lorsqu'ils décrivent un phénomène, les opérateurs experts omettent souvent de décrire des phénomènes intermédiaires. Contrairement au raisonnement des opérateurs experts, le raisonnement des ingénieurs qui se base sur des équations, est beaucoup plus systématique et il permet de structurer les raisonnements. Nous avons donc jugé préférable de nous orienter vers la méthode systématique de modélisation causale profonde, décrite dans le Chapitre 2.

L'approche profonde garantit une modélisation correcte mais nous avons soulevé un certain nombre de problèmes à son application :

- il est nécessaire en pratique de considérer des variables pseudo-exogènes,
- nous avons amélioré⁴⁻⁶ l'algorithme de réduction initialement proposé dans [Travé-Massuyès *et al.* 2001],
- nous avons créé un algorithme d'approximation (non implémenté) car en pratique il est nécessaire de faire des approximations. L'algorithme d'approximation permet d'effectuer des hypothèses de simplification tout en ne pénalisant pas le diagnostic.

4.3.4.4 Quantification du modèle

Nous avons rencontré plusieurs problèmes lors de la quantification du modèle.

- Dans le cadre de l'approche experte, la structure des équations du modèle n'est pas connue et nous avons dû faire des hypothèses avant d'effectuer l'identification [Heim, Cauvin et Gentil 2001]. Dans le cadre de l'approche profonde, la structure des équations du modèle à identifier est connue.
- Il est difficile en pratique de disposer de données pertinentes pour effectuer l'identification des paramètres du modèle. Il est préférable d'effectuer des manipulations en boucle ouverte sur le procédé pour l'exciter suffisamment. Des contraintes de temps ou de sécurité empêchent en général ces manipulations.
- Nous avons observé, sur le pilote de FCC des jeux systématiques sur les vannes. Ces jeux sont source d'incertitudes importantes sur le modèle. Nous avons utilisé une prédiction à un pas (cf. Formule 3-3) qui est moins sensible aux incertitudes de modèle.

⁴⁻⁶ Dans le cas où la variable inconnue à éliminer agit sur elle même via une autre variable.

4.3.4.5 Application à d'autres procédés

La méthode de diagnostic a été appliquée à un générateur de vapeur. Initialement, il a été envisagé d'en établir un modèle linéarisé autour du point de fonctionnement. Cependant, les régulateurs fonctionnant en tout ou rien, cette démarche a été abandonnée car elle n'était pas adaptée. L'introduction de relations non linéaires a été nécessaire.

4.3.4.6 Conclusion

La méthode de diagnostic présentée dans ce document est originale car elle fait successivement appel à différentes techniques pour le diagnostic d'un processus. Cet enchaînement de techniques, qui ont été chronologiquement introduites, a été en pratique nécessaire pour préciser l'information générée par le module de diagnostic. De nombreux problèmes pratiques ont été rencontrés pour appliquer la méthode (obtention de données pertinentes pour l'identification, jeux des vannes, etc.).

4.3.5 Contexte de l'étude et du développement du module informatique

4.3.5.1 Introduction

La première boîte à outil, intitulée TB3.4 (Toolbox), comprend la partie génération du graphe causal, génération des alarmes et interprétation des résidus par la logique floue.

La seconde entité, TB5.4, comprend l'association des composants aux influences du graphe causal et la génération des diagnostics. La partie graphes experts n'est pas incluse dans le projet CHEM.

Le support informatique choisi dans le cadre du projet CHEM est G2 : ce logiciel qui permet de développer rapidement des applications en temps réel est présenté dans la section 4.3.5.3. Les sections suivantes présentent respectivement les documents de spécification et le manuel d'utilisation du module développé.

4.3.5.2 Le projet CHEM

L'IFP coordonne le projet européen CHEM "Advanced Decision Support System for Chemical and Petrochemical Processes" dont l'objectif est de développer et de faire collaborer des outils de diagnostic et d'aide à la décision pour la conduite d'unités industrielles de chimie, de raffinage et de pétrochimie. CHEM intègre des "boîtes à outils informatiques" qui mettent en œuvre des technologies variées - automatique, statistique, graphes causaux, logique floue, à base d'intervalle, etc. - afin de donner une information ciblée aux opérateurs et d'améliorer le fonctionnement des unités. Le système sera testé en conditions réelles sur plusieurs sites, notamment sur le pilote FCC de l'IFP. CHEM, se déroule sur 3 ans (2001-2004) et réunit quinze partenaires industriels et universitaires européens de huit pays différents : Espagne, Finlande, France, Grande-Bretagne, Norvège, Pays-Bas, Pologne et Suède. Le projet CHEM est divisé en huit Workpackages (WP) :

- WP1 : Concepts généraux et méthode d'intégration (l'objectif est de définir une méthodologie de communication entre les boîtes à outil développées).
- WP3 : Analyse du signal et évaluation de la situation (l'objectif est d'extraire des informations qualitatives ou semi-qualitatives du signal, de générer des alarmes et des modes de fonctionnement, cette information est indispensable pour les autres WPs).
- WP5 : Diagnostic et gestion des alarmes (l'objectif est d'obtenir des modèles dédiés au diagnostic et de les utiliser pour localiser les défaillances).
- WP7 : Aide à la décision (l'objectif est de disposer de plusieurs structures permettant de fournir à l'opérateur des informations sur des actions à entreprendre étant donné les diagnostics réalisés précédemment).
- WP9 : Planification en ligne (l'objectif est de générer des plans de production en fonction des informations disponibles sur le procédé).
- WP10 : Intégration des boîtes à outil.
- WP11 : Tests et validations industrielles. Les boîtes à outil et des groupes de boîtes à outils sont testés sur différents procédés industriels allant du pilote de FCC à un procédé de fabrication de pâte à papier.
- WP12 : Gestion de projet.

Les deux boîtes à outil (Toolboxes) développées par l'IFP (TB3.4 et TB5.4) s'insèrent respectivement dans WP3 et WP5. D'autres boîtes à outils ayant des objectifs similaires aux TB3.4 et TB5.4 sont incluses dans les WP5 et WP4.

4.3.5.3 Environnement informatique de développement

Le module de diagnostic a été élaboré dans l'environnement G2 de la société Gensym⁷. G2 est un langage orienté objet qui permet de développer rapidement et facilement des applications en temps réel. La création et le développement du graphe causal s'effectue de manière graphique : les nœuds du graphe causal sont des objets graphiques créés et déplacés facilement, les influences entre les nœuds sont aussi créées graphiquement. Les objets G2 (nœuds, arcs, etc.) ont des attributs (mesure, référence locale, paramètres des fonctions de transferts, composants industriels des supports, etc.) qui peuvent être modifiés très facilement en sélectionnant l'objet adéquat. G2 permet d'acquérir et de traiter en temps réel des données industrielles ou de traiter des données *a posteriori* à partir d'une base de données.

4.3.5.4 Document de spécification

Le document de spécification du module informatique a été rédigé selon le formalisme de représentation IDEFØ. IDEFØ (Integrated Definition Language), (Structured Analysis and Design Technic) est une méthode d'analyse et de spécification fonctionnelle. Cette méthode, conforme aux normes FIPS (Federal Information Processing Standard 183 et 184), a été mise au point en 1976 aux États-Unis par Douglas T. Ross (Société SOFTECH avec ITT) pour des applications aérospatiales. De renommée mondiale, la méthode IDEFØ est largement utilisée en France dans les domaines de l'informatique, de l'aéronautique, du spatial et des télécommunications. Elle peut tout aussi bien s'appliquer à d'autres activités. Le formalisme IDEFØ est utilisé pour analyser et décrire un système sous son aspect fonctionnel. La description peut porter sur :

- les fonctions remplies par un élément matériel, (offrir les fonctionnalités d'un réveil matin, réguler la température, ...),
- la description d'un processus (organiser un voyage, former des stagiaires à IDEFØ, etc.); les deux à la fois (transporter des voyageurs, etc.),

⁷ <http://www.gensym.com>

Il est possible d'utiliser la méthode à plusieurs stades d'avancement dans la réalisation d'un projet. Ainsi peuvent être réalisés :

- le modèle d'un système existant,
- le modèle du système souhaité,
- le modèle du système à réaliser,
- un modèle théorique général, etc.

L'élément de base d'un diagramme IDEFØ est le **module**. Dans un diagramme, celui-ci représente une activité fonctionnelle du système que l'on souhaite décrire. Cette fonction est identifiée par un verbe à l'infinitif précisé éventuellement par un complément (réguler la température, organiser une visite, offrir les fonctionnalités d'un réveil, etc.). Les modules communiquent entre eux et avec leur environnement par l'intermédiaire de leurs **interfaces**. Ils échangent ainsi des flux matériels (nourriture, etc.) ou d'information (liste des amis; etc.). Les interfaces ont une fonction différente selon la position de leurs points de connexion (cf. Figure 4-14).

Les entrées, connectées à gauche, sont consommées par l'activité pour produire les sorties. Elles ne conditionnent pas le comportement du module. Elle ne peuvent pas non plus déclencher son activité. Un module peut ne pas avoir d'entrée quand ce qu'il produit ne nécessite l'apport d'aucun flux fonctionnel extérieur (ex : production de l'heure par une montre).

Les sorties, connectées à droite, sont le résultat de l'activité fonctionnelle du module. Celui-ci doit, bien entendu, avoir au moins une sortie.

Les contrôles, connectés en haut, ne sont jamais consommés par l'activité. Il agissent sur son déroulement en la déclenchant (occasion à fêter) ou en influençant fortement son comportement (recette de cuisine). Une activité doit avoir au moins un contrôle qui la déclenche.

Les mécanismes, connectés en bas, ne sont pas à considérer comme des éléments du modèle fonctionnel. Ils offrent la possibilité de décrire des éléments physiques mis en œuvre pour réaliser la fonction. Un diagramme peut ne pas comporter de mécanismes.

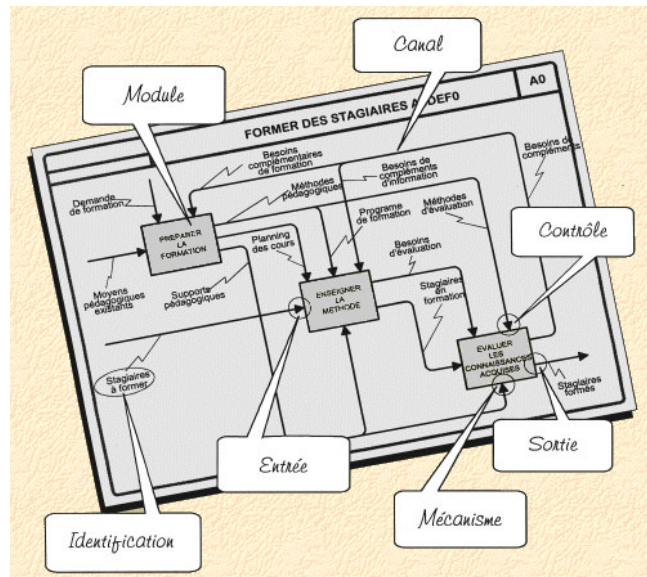


Figure 4-14 : Illustration du formalisme IDEF0

4.3.5.5 Conclusion

La méthode présentée dans ce document a donné lieu à deux boîtes à outils développées dans le cadre du projet Européen CHEM. Cette démarche a nécessité l'établissement de documents de spécification répondant à un formalisme générique et pouvant être facilement interprétés par les autres membres du projet. La rédaction de manuels utilisateurs a permis de communiquer la méthode et d'avoir un échange avec d'autres chercheurs dans le domaine du diagnostic.

4.3.6 Résultats escomptés sur un FCC industriel

4.3.6.1 Introduction

Cette section présente les limites du modèle causal du pilote de FCC, puis les efforts à fournir pour appliquer la méthode à un premier procédé industriel de FCC. Enfin, les efforts à fournir pour l'appliquer ensuite à d'autres procédés de FCC sont présentés.

4.3.6.2 Limites du modèle du pilote de FCC

La priorité sur le pilote de FCC est d'atteindre le plus rapidement possible la stabilité du bilan pression. Le bilan thermique de l'unité pilote de FCC est très long à s'établir et ses relations, bien que décrites dans le MSR, n'ont pas été introduites dans le module de diagnostic pour les raisons suivantes :

- elles n'apportent pas d'information pertinente aux utilisateurs pour le diagnostic du pilote de FCC,
- la charge du pilote de FCC changeant souvent, il est nécessaire de la caractériser à chaque changement, ce qui serait à ce niveau de développement fastidieux,
- le pilote de FCC n'étant pas adiabatique, il est nécessaire de modéliser les pertes thermiques.

4.3.6.3 Efforts nécessaires pour la déclinaison de la méthode sur un FCC industriel

- Les efforts identifiés pour appliquer la méthode à un procédé industriel sont les suivants : les procédés industriels de FCC contiennent une vanne supplémentaire de circulation de catalyseur à l'entrée du riser, il est donc nécessaire de légèrement modifier le modèle du pilote de FCC.
- Les variables mesurées sur un FCC industriel sont différentes de celles mesurées sur le pilote de FCC, il est donc nécessaire d'effectuer à nouveau l'opération de réduction.
- Il est nécessaire d'identifier à nouveau les paramètres du modèle.
- Il est nécessaire d'établir les règles des graphes experts de tous les composants.

L'application à une unité industrielle sera pénalisée pour les raisons suivantes :

- La qualité des mesures est très médiocre sur les unités industrielles et les étalonnages sont très peu fréquents.
- Etant donné qu'il y a moins de mesures disponibles sur une unité industrielle de FCC, l'information fournie par le module de diagnostic sera moins précise.
- Il est très rare que les exploitants d'unités industrielles acceptent d'effectuer des manipulations en boucle ouverte sur leurs procédés en vue d'une identification de paramètres. Il sera alors judicieux soit de rechercher dans l'historique⁴⁻⁸ de fonctionnement des données excitées soit d'effectuer une identification en boucle fermée qui seraient mieux acceptée par les utilisateurs.

Des raisons humaines peuvent entraver l'application. En effet, il faut que le produit soit accepté par les utilisateurs finaux et que ceux-ci soient convaincus de son utilité. Il sera donc nécessaire d'effectuer un travail profond et persuasif de communication et certainement de formation auprès des utilisateurs finaux. Le module ne doit être envisagé que comme une aide à la conduite pour anticiper les dysfonctionnements, mais la prise de décision reste sous la responsabilité de l'opérateur.

On doit former les opérateurs à la sémantique très différente dans les deux cas.

- Dans le cadre de l'approche ensembliste, il faut mettre l'effort lors de l'élaboration du modèle.
- Dans le cadre de l'approche par logique floue, l'effort doit être centré sur l'interprétation correcte des alarmes par les opérateurs.

4.3.6.4 Efforts nécessaires pour la déclinaison de la méthode sur un FCC industriel

Tous les procédés industriels de FCC sont différents les uns des autres et en particulier ils font appel, dans certaines zones physiques, à des technologies différentes. D'autres part, les mesures ne sont pas effectuées au même endroit d'un procédé de FCC à un autre.

⁴⁻⁸ Dans les bases de données, certaines quantités correspondent à des moyennes journalières et sont inutilisables pour paramétrer un modèle dynamique.

L'application de la méthode à un procédé industriel nécessite donc systématiquement une révision du modèle. Il est donc nécessaire de modifier légèrement le modèle, d'effectuer à nouveau les opérations de réduction, d'approximation et d'identification du Chapitre 2.

4.3.6.5 Conclusion

La méthode décrite dans ce document a été appliquée à un pilote de FCC et a donné des résultats très encourageants.

Ce développement, qui a été effectué dans le cadre du projet Européen CHEM, est à la source de documents techniques conformes à la méthode d'analyse et de spécification fonctionnelle IDEFØ.

Ces documents ont été partagés par des chercheurs européens, partenaires du projet, dans le domaine du diagnostic. Il est envisagé d'appliquer la méthode à un procédé industriel de FCC, mais un certain nombre de points doivent être abordés concernant l'adaptation du modèle développé sur le pilote de FCC, le suivi, la communication et la formation des utilisateurs finaux.

4.3.7 Conclusion

Nous avons appliqué à un pilote de FCC la méthode décrite dans le Chapitre 2 et dans le Chapitre 3. Nous avons obtenu des résultats très encourageants sur des cas réels de fautes.

Ce travail a été effectué dans le cadre du projet Européen CHEM et nous disposons du système informatique de diagnostic ASCO qui peut être connecté à un autre procédé.

Cette première application nous a permis d'élaborer une méthodologie de diagnostic et de lister les problèmes que nous avons rencontrés et ceux que nous pensons rencontrer lors de l'application à un FCC industriel.

Malgré les connaissances acquises et le développement du système ASCO, des efforts importants restent à fournir car le pilote de FCC n'est pas complètement

représentatif d'un FCC industriel. Le chapitre suivant présente les perspectives de travaux.

Chapitre 5

Conclusion et perspectives

Dans ce mémoire, nous avons développé une méthodologie de diagnostic de procédés et nous l'avons appliquée à un pilote de FCC. Elle fait successivement appel à deux techniques s'appuyant sur un modèle causal pour la génération d'alarmes et la localisation de défauts. Ensuite elle s'appuie sur un système à base de règles pour l'identification des défauts.

Le modèle causal, qui est un modèle de bon fonctionnement et qui met en valeur les influences entre les variables du procédé, est alimenté en parallèle avec le procédé et calcule des références. Des alarmes sont générées en comparant ces références aux mesures. Ces alarmes déterminent si chaque variable a un comportement anormal ou non vis-à-vis des variables mesurées qui en sont la cause directe ou vis-à-vis des consignes et perturbations mesurables.

Si une variable est incohérente avec ses causes directes alors l'ensemble des composants qui sous tendent ces influences constitue un conflit. L'étape de localisation, qui consiste à étudier les conflits, génère des diagnostics possibles.

Finalement, l'étape d'identification de défauts consiste à activer les graphes experts des composants membres des diagnostics. Ces graphes experts contiennent des bases de règles qui constituent des modèles de mauvais fonctionnement de ces

composants. Elles associent les défaillances aux symptômes. Les symptômes sont créés par un traitement du signal approprié sur les variables mesurées du procédé et sont analysés au moyen d'expressions logiques booléennes simples. Lorsque l'on a caractérisé un certain type de symptômes selon la logique convenable, alors un message qui identifie la défaillance sur un composant est transmis à l'opérateur. Ce message contient des confirmations qualitatives à effectuer sur le procédé pour valider l'identification. Il propose des actions à entreprendre en ligne ou pour la maintenance du procédé. Enfin, il décrit les répercussions à long terme de la défaillance sur le fonctionnement du procédé.

La méthodologie de diagnostic proposée est très liée à la notion de causalité : le modèle causal reflète la connaissance des ingénieurs et les processus de réflexion qui sont effectués pour comprendre la propagation des phénomènes normaux et anormaux dans une installation industrielle.

L'obtention du modèle causal n'est évidemment pas facile et nous avons décrit soigneusement les différentes étapes permettant de l'obtenir à partir d'une connaissance profonde sur l'installation. Nous avons élaboré sept étapes préliminaires permettant de générer le modèle structurel des relations (cf. Chapitre 2). Le MSR est un système d'équations élémentaires, répondant au principe de localité, entre des variables connues et non connues. Le support de chaque relation du MSR est connu.

Nous appliquons ensuite un algorithme d'ordonnement causal au MSR. Il permet de générer le modèle causal structurel. Si le système d'équations du MSR est complet alors un MCS est assurément obtenu (cf. Chapitre 2), même en présence de boucles de rétroaction. Le MCS contient les mêmes informations que les MSR, cependant les influences entre les variables sont orientées. Plusieurs solutions sont possibles.

Nous effectuons ensuite deux opérations dites de réduction puis d'approximation sur le MCS (cf. Chapitre 2). L'opération de réduction appliquée au MCS permet d'extraire les influences entre les variables connues : le modèle causal réduit est obtenu. L'opération d'approximation est finalement appliquée au MCR. Elle permet de générer le modèle causal approché : cette opération permet de négliger des

influences tout en garantissant un diagnostic acceptable (le MCA est utilisé comme support du module de diagnostic).

Les influences du MCA ont été associées à des fonctions de transfert classiques. Leurs paramètres sont estimés ou peuvent être obtenus via une connaissance physique (données de la littérature, paramètres connus ou mesurés, etc.). Le MCA est donc un modèle dynamique de bon comportement du système. Comme il permet de prédire le comportement normal des variables autour du point de fonctionnement et d'expliquer les phénomènes de propagation des défauts, il constitue un outil pertinent pour le diagnostic. Si toutefois le modèle est très fortement non linéaire, alors des relations non linéaires peuvent être implémentées dans le MCA, mais nous n'avons pas eu besoin de suivre cette démarche dans le cadre du pilote de FCC.

Les résidus spécifiques calculés grâce au MCA sont les résidus locaux et globaux qui permettent respectivement de générer des alarmes locales et globales. La comparaison de la mesure et des références ne peut s'effectuer de manière simple à cause des incertitudes du modèle et des mesures. Il a été nécessaire d'utiliser des techniques permettant de prendre en compte ces imprécisions : une technique à base de logique floue et une technique faisant appel aux intervalles modaux ont été testées et comparées.

L'approche par logique floue propose une interprétation graduelle de l'évolution de la variable vers un état anormal. Elle permet donc d'anticiper les défaillances. Son inconvénient principal est qu'elle s'appuie sur des seuils fixes de détection. Nous avons implémenté les algorithmes de logique floue dans un module informatique que nous avons nommé ASCO.

L'approche par intervalles génère des enveloppes dont la taille dépend de l'excitation des entrées, par contre la décision est binaire. L'approche par intervalles garantit le résultat mais elle retarde la décision. L'approche par logique floue est différente : des fausses détections sont possibles, par contre elle permet d'anticiper les problèmes. Nous avons testé les algorithmes par intervalles mais cette technique n'est pas encore assez mûre et il est encore nécessaire de lui apporter des améliorations.

Nous avons aussi testé une méthode d'identification ensembliste pour obtenir les paramètres nécessaires pour l'approche par intervalles mais cette méthode a donné des résultats trop pessimistes qui n'ont pas pu être utilisés par l'approche de détection de défauts par intervalles. Le couplage algorithme d'estimation des paramètres et diagnostic est un problème qui mériterait beaucoup d'attention.

Nous avons élaboré une méthode permettant de combiner les avantages de l'approche par logique floue et les avantages de l'approche ensembliste. Deux enveloppes sont générées via l'approche ensembliste et leurs positions par rapport à la mesure sont interprétées via un raisonnement par logique floue. La taille de ces enveloppes dépend de l'excitation de l'entrée (avantage des intervalles modaux), et l'interprétation de leurs positions par rapport à la mesure est graduelle (avantage du raisonnement flou).

L'étape de localisation permet, en interprétant les supports des influences et les alarmes locales, de déterminer un diagnostic, i.e. de cibler un groupe de composants suspects. Le support des relations est constitué d'un nombre de composants d'autant plus grand que les mesures sont rares. Par conséquent, nous proposons une étape supplémentaire qui identifie la défaillance et focalise l'attention sur un composant particulier. L'étape d'identification est alors nécessaire pour caractériser plus précisément la défaillance d'un composant.

Des bases de règles sont utilisées pour l'identification des défauts. L'inconvénient des systèmes à base de règles est que des défauts différents peuvent avoir des symptômes identiques. Ce phénomène est lié à la propagation des défauts et aux phénomènes de causalité. Par conséquent, s'ils sont utilisés sans aucun traitement préalable, les systèmes à base de règles apportent une information très souvent erronée. Dans ce mémoire un système à base de règles est associé au composant dont il décrit les défaillances. La base de règles d'un composant est activée si et seulement si le composant fait partie d'un diagnostic. Ce traitement préalable permet d'utiliser de manière judicieuse les systèmes à base de règles et évite d'apporter des informations erronées.

Finalement, ce travail laisse un nombre de problèmes ouverts.

- Les résidus du modèle causal approché ne permettent pas de prévenir la situation suivante. Il serait possible qu'un jeu de consignes indépendamment cohérentes entraîne le processus dans un mode de fonctionnement dangereux mais tout à fait normal étant donné les consignes. Les alarmes globales ne permettent pas de détecter ces défaillances humaines. Des alarmes à seuils classiques permettent de pallier simplement à ce problème.
- Dans la méthode de localisation, il serait intéressant de considérer que le support des relations est composé de défaillances types (bouchage, fuite, etc.) sur des composants plutôt que simplement des composants⁵⁻¹. De tels modèles de défauts permettraient de mieux cibler les défaillances et de proposer un diagnostic encore plus précis.
- Dans la méthode d'identification de défaut, il serait judicieux d'extraire des bases de règles issues du pilote de FCC, ou d'autres bases de connaissances, des règles plus générales associant des symptômes à des défaillances pour des composants spécifiques. On pourrait imaginer qu'un industriel vende un composant avec sa base de règles.
- Afin d'en étudier les avantages et les limites, la méthodologie décrite dans ce document devrait être appliquée à un FCC industriel. Quelques problèmes sont déjà identifiés par une telle démarche : tous les procédés de FCC sont légèrement différents, d'autre part les mesures sont moins fréquentes que sur le pilote et situées à des endroits différents. Il sera donc nécessaire de reprendre l'ensemble de la méthode et d'adapter le modèle aux particularités d'un autre FCC.

⁵⁻¹ Une RRA peut être sensible à une fuite sur un composant et ne pas être sensible à un bouchage sur ce composant.

Bibliographie

- [Adrot 2000] Adrot O., Diagnostic à base de modèles incertains utilisant l'analyse par intervalles : l'approche bornante. Doctorat de l'Institut National Polytechnique de Lorraine, France, 2000.
- [Armengol 1999] Armengol J., Application of Modal Interval Analysis to the simulation of the behaviour dynamic system with uncertain parameters. Thèse de doctorat de l'Université de Girone, Espagne, 1999.
- [Ahriz 1998] Ahriz H., Modélisation automatique de systèmes physiques. Application au diagnostic. Thèse de doctorat de l'Université de Savoie, France, 1998.
- [Bakhtazad, Palazoglu et Romagnoli 1998] Bakhtazad A., Palazoglu A., Romagnoli J.A, Process Trend Analysis Using Wavelet Based De-noising. 3rd IFAC Workshop on On-Line Fault Detection and Supervision in the Chemical Process Industries, France, 1998.
- [Benquילו *et al.* 1999] Benquילו C., Ruiz D., María Nougues J., Puigjaner L., A Hybrid Neural Network-First Principles Approach for Process Modeling. 8th Mediterranean Congress of Chemical Engineering, Espagne, 1999.
- [Bonarini et Bontempi 1994] Bonarini A., Bontempi G., A Qualitative Simulation Approach for Fuzzy Dynamic Models. ACM Transactions on Modeling and Computer Simulation, Vol. 4, No. 4, pp. 258-313, 1994.
- [Busson *et al.* 1998] Busson F., Aïtouche A., Ould Bouamama B., Staroswiecki M., Sensors Failure Detection in Steam Condensers. 3rd IFAC Workshop on On-Line Fault Detection and Supervision in Chemical Process Industries, France, 1998.

- [Bullemer et Nimmo 1996] Bullemer P.T., Nimmo, I., A Training Perspective on Abnormal Situation Management: Establishing an Enhanced Learning Environment. AICHE conference on Process Plant Safety, Etats-Unis, 1996.
- [Cauvin 1995] Cauvin S., Un environnement générique à base de connaissances pour la supervision de procédés de raffinage et de pétrochimie. Thèse de Doctorat du Conservatoire National des Arts et Métiers, France, 1995.
- [Clément 1984] Clément T., Estimation d'incertitude paramétrique dans un contexte de bruit de sortie inconnu mais borné. Thèse de doctorat de l'Institut National Polytechnique de Grenoble, France, 1984.
- [Colomer, Mendelez et Gamero 2002] Colomer J., Melendez J., Gamero F.C., Qualitative Representation of Process Trends for Situation Assessment Based on Cases. 15th Triennial World Congress of the International Federation of Automatic Control, Espagne, 2002.
- [Cordier *et al.* 2000a] Cordier M.O., Dague P., Dumas M., Levy F., Montmain J., Staroswiecki M., Travé-Massuyès L., AI and automatic control approaches of model based diagnosis: links and underlying hypotheses. 4th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), pp. 274-279, Hongrie, 2000.
- [Cordier *et al.* 2000b] Cordier M.O., Dague P., Dumas M., Levy F., Montmain J., Staroswiecki M., Travé-Massuyès L., A comparative analysis of AI and control theory approaches to model-based diagnosis. European Conference on Artificial Intelligence, Allemagne, 2000.
- [Dauphin-Tanguy 2000]. Dauphin-Tanguy G., Les Bonds Graphes. Hermès Sciences , 2000.
- [De Kleer et Williams 1987] de Kleer J., Williams, B. Diagnosing Multiple Faults, Artificial Intelligence. Vol. 32, pp. 97-130, 1987.
- [De Kleer et Brown 1990] De Kleer J., Brown J.S., A Qualitative physics based on confluences, Artificial Intelligence, Vol. 24, pp. 7-83, 1984.
- [Demerle et Siguerdidjane 2002] Demerle M., Siguerdidjane H., Analyse des systèmes linéaires. Hermès Sciences, 2002.

-
- [Diego *et al.* 2001] Ruiz D., Cantón J., Nogués J.M., Espuña A., Puigjaner L., On-line fault diagnosis system support for reactive scheduling in multipurpose batch chemical plants. *Computers & Chemical Engineering*, Vol. 25, pp. 829-837, 2001.
- [Dubuisson 2000] Dubuisson B., *Diagnostic par intelligence artificielle et reconnaissance des formes*. Hermès Sciences, 2000.
- [Dvorak et Kuipers 1989] Dvorak D., B. J. Kuipers B.J., Model-based monitoring of dynamic systems. 11th International Joint Conference on Artificial Intelligence (IJCAI), Etats-Unis, 1989.
- [Dziopa 1996] Dziopa P., *Représentation Multi Modèles pour la Supervision de Procédés Industriels Continus*. Thèse de doctorat de l'Institut National Polytechnique de Grenoble, France, 1996.
- [Evsukoff, Montmain et Gentil 1997] Evsukoff A., Montmain J., Gentil S., Dynamic model based supervising and causal knowledge-based fault detection and isolation. 3rd IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), pp. 699-704, Angleterre, 1997.
- [Evsukoff 1998] Evsukoff A., *Le raisonnement approché pour la surveillance de procédés*. Thèse de doctorat de l'Institut National Polytechnique de Grenoble, 1998.
- [Fagarasan, Ploix et Gentil 2001] Fagarasan I., Ploix S., Gentil S., Bounding approach for fault detection and diagnosis, IFAC Symposium Large Scale Systems LSS 2001, Roumanie, 2001.
- [Farlane *et al.* 1993] McFarlane R.C., Reinemann R.C., Bartee J.F., Georkakis C., Dynamic simulation for a model IV fluid catalytic cracking unit. *Computer and Chemical Engineering*, Vol. 17, pp. 275-300, 1993.
- [Fenu et Parisini 1998] Fenu G, Parisini T., Kernel regression and neural networks for model-free fault diagnosis. 2nd IFAC workshop on line fault detection and supervision in the chemical industries, France, 1998.
- [Flaus et Gentil 2003] Flaus J.M., Gentil S., Hybrid System Automatic Modeling for Diagnostic Purposes. 5th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), Etats Unis, 2003.
- [Forbus 1984] Forbus K.D., *Qualitative Process Theory*. *Artificial Intelligence*, Vol. 24, No. 1, pp. 85-168, 1984.

- [Ford et Fulkerson 1956] Ford L.R., Fulkerson D.R., Maximal flow through a network. *Canadian Journal Mathematics*, Vol. 8, pp. 399-404, 1956.
- [Frank 1990] Frank P.M., Fault Diagnosis in Dynamic Systems Using Analytical and Knowledge-based Redundancy, A Survey and Some New Results. *Automatica*, Vol. 26, No. 3, pp. 459-474, 1990.
- [Frank 1996] Frank P.M., Analytical and Qualitative Model-Based Fault Diagnosis - A survey and some new results. *European Journal of Control*, Vol. 2, No. 1, pp. 6-23, 1996.
- [Frank et Ding 2000] Frank P., Ding S., Current development in the theory of FDI. 4th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (SafeProcess), pp. 16-27, 2000.
- [Gardeñes *et al.* 2001] Gardeñes E., Sainz M.A, Jorba I., Calm R., Estela R., Mielgo H., Trepas A., Modal Intervals. *Reliable Computing*. Vol. 7, No.2, pp. 77-111, 2001.
- [Gertler et Singer 1990] Gertler J., Singer D., A New Structural Framework for Parity Equation-based Failure Detection and Isolation. *Automatica*, Vol. 26, 1990.
- [Gondran et Minoux 1979] Gondran M., Minoux M., *Graphes et algorithmes*. Eyrolles, Paris, pp. 160-163, 1979.
- [Goldsztejn 2000] Goldsztejn A., Applicabilité des Intervalles Modaux. Rapport de DEA d'Informatique, Institut Supérieur d'Electronique du Nord, France, 2000.
- [Han et Chung 2001a] Han I.S., Chung C.B., Dynamic Modeling and Simulation of a Fluidized Catalytic Cracking Process: Part 1. Process Modeling. *Chemical Engineering Science*, Vol. 56, pp. 1951-1971, 2001.
- [Han et Chung 2001b] Han I.S., Chung C.B., Dynamic Modeling and Simulation of a Fluidized Catalytic Cracking Process: Part 2. Property Estimation and Simulation. *Chemical Engineering Science*, Vol. 56, pp. 1973-1990, 2001.
- [Heim *et al.* 1999] Heim B., Cauvin S., Duplan J.L., Girardon S., Système d'aide à la conduite de l'unité pilote de FCC. Rapport interne, Institut Français du Pétrole, Réf : 54032, France, 1999.

-
- [Heim, Cauvin et Gentil 2000a] Heim B., Cauvin S., Gentil S., Causal and fuzzy reasoning methodology for cascaded loops diagnosis. 2000 IAR ICD Workshop devoted to Intelligent Control and Diagnosis, France, 2000.
- [Heim, Cauvin et Gentil 2000b] Heim B., Cauvin S., Gentil S., Raisonement approximatif et causal pour le diagnostic et la supervision de procédés de raffinage. Application à un pilote de FCC. Rapport interne, Institut Français du Pétrole, Réf : 54278, France, 2000.
- [Heim, Cauvin et Gentil 2001] Heim B., Cauvin S., Gentil S., FCC process diagnosis based on a causal and heuristic approach. 4th Workshop on On-Line Fault Detection and Supervision in the Chemical Process Industries, Corée, 2001.
- [Heim *et al.* 2002a] Heim B., Cauvin S., Gentil S., Travé-Massuyès L., Modélisation causale pour le diagnostic de procédés. Application à un pilote de FCC. Rapport interne, Institut Français du Pétrole, Réf : 56669, France, 2002.
- [Heim *et al.* 2002b] Heim B., Gentil S., Cauvin S., Travé-Massuyès S., Braunschweig B., Fault diagnosis of a chemical process using causal uncertain model. 15th European Conference on Artificial Intelligence, France, 2002.
- [Heim *et al.* 2003a] Heim B., Gentil S., Celse B., Cauvin S., Travé-Massuyès L., FCC diagnosis using several causal and knowledge models. 5th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), Etats-Unis, 2003.
- [Heim *et al.* 2003b] Heim B., Travé-Massuyès L., Gentil S., Celse B., Cauvin S., Causal modeling Methodology. Application to a Fluid Catalytic Cracking Plant. Rapport interne, Institut Français du Pétrole, Réf : 57760, France, 2003.
- [Heim *et al.* 2003c] Heim B., Gentil S., Celse B., Cauvin S., Travé-Massuyès S., Application of a model based diagnostic module to an FCC pilot process. Rapport interne, Institut Français du Pétrole, Réf : 57761, France, 2003.
- [Iri *et al.* 1979] Iri M., Aoki K., O'Shima E., Matsuyama H., An algorithm for diagnosis of system failures in the chemical process. *Computers & Chemical Engineering*, Vol. 3, pp. 489-493, 1979.
- [Isermann 1993] Isermann R., Fault diagnosis of machines via parameter estimation and knowledge processing - tutorial paper. *Automatica* Vol. 29, No. 4, pp. 813-835, 1993.

- [Isermann et Ballé 1997] Isermann R., Ballé P., Trends in the application of model-based fault detection and diagnosis of technical processes. *Control Engineering Practice*, Vol. 5, pp. 5, 709-719, 1997.
- [Iwasaki et Simon 1986] Iwasaki S., Simon H., Causality in device behavior. *Artificial Intelligence*, Vol. 29, No. 1-3, pp. 3-32, 1986.
- [Janati, Adrot et Ragot 2001] Idrissi J., Adrot O., Ragot J., Multi-Fault Detection of Systems with Bounded Uncertainties 2001 40th IEEE Conference on Decision and Control (CDC), Floride, 2001.
- [Jia, Martin et Morris 1998] Jia F., Martin E.B., Morris A.J., Non-linear Multi-Way Principal Components Analysis for Processes Monitoring. *Computers and Chemical Engineering*, Vol. 22, pp. 851-854, 1998.
- [Klir et Yuan 1995] Klir G.J, Yuan B., *Fuzzy Sets and Fuzzy Logic. Theory and Applications*. Prentice Hall, 1995.
- [Kramer et Palowitch] Kramer M. A. et Palowich Jr.B.L., A rule-based approach to fault diagnosis using the signed directed graph. *AIChE Journal*, Vol. 33, No.7, pp. 1067-1078, 1987.
- [Krysander et Nyberg 2002] Krysander M., Nyberg M., Structural Analysis utilizing MSS Sets with Application to a Paper Plant. Thirteenth International Workshop on Principles of Diagnosis, Australie, 2002.
- [Kuipers 1994] Kuipers B. J., *Qualitative Reasoning: Modeling and Simulation with Incomplete Knowledge*. Cambridge, MA: MIT Press, 1994.
- [Leseq et Barraud 2000] Leseq S., Barraud A., Fault detection using on-line wavelet analysis : application to induction motor. 4th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), Hongrie, 2000.
- [Leseq, Barraud et Tran-Dinh 2003] Leseq S., Barraud A., Tran-Dinh K.Q., Numerical Accurate Computations for Ellipoidal State Bounding, Note interne du Laboratoire d'Automatique de Grenoble, AP03-009, 2003.
- [Leyval, Gentil et Feray-Beaumont 1994] Leyval L., Gentil S., Feray-Beaumont S., Model based causal reasoning for process supervision. *Automatica*, Vol. 30, No. 8, pp. 1295-1306, 1994.

-
- [Loiez et Taillibert 1997] Loiez E., Taillibert P., Polynomial Temporal Band Sequences for Analog Diagnosis. IJCAI 1997, pp. 474-479, Nagoya, Japon, 1997.
- [Lucas 1994] Lucas, B., Méthode d'Aide à la Modélisation par Graphes de Liaison et Utilisation pour le Diagnostic Qualitatif de Systèmes Physiques. Thèse de doctorat de l'Université Paris XI-Orsay, 1994.
- [Milanese et Belforte 1982] Milanese M., Belforte G., Estimation theory and uncertainty intervals evaluation in presence of unknown but bounded errors: Linear families of models and estimators. IEEE Transactions on Automatic Control, Vol. 27, No. 2, pp. 408-414, 1982.
- [Molina 2000] Molina J.A.J, The Design and Implementation of the Fault-Isolation Module for the Causal Simulator Ca-En of the L.A.A.S-C.N.R.S. Final project of the Computer Engineering diploma of the University of Saragosse, Espagne, 2000.
- [Montmain et Gentil 1993] Montmain J., Gentil S., Decision-making in Fault Detection: a fuzzy approach. International Conference on Fault Diagnosis (Tooldiag), Toulouse, 1993.
- [Montmain et Gentil 2000] Montmain J., Gentil S., Dynamic causal model diagnostic reasoning for on-line technical process supervision. Automatica, Vol. 36, pp. 1137-1152, 2000.
- [Mosterman, Biswas et Narasimham 1997] Mosterman P., Biswas G., Narasimham S., Measurement selection and diagnosability of complex physical systems. Eight International Workshop on Principles of Diagnosis (DX), pp. 79-86, France, 1997.
- [Mylaraswamy et Venkatasubramanian 1997] Mylaraswamy D., Venkatasubramanian V., A Hybrid Framework for Large-scale Process Fault Diagnosis. Computer and Chemical Engineering, Vol. 21, pp. 935-940, 1997
- [Ould Bouamama 2003] Ould Bouamama B., Modélisation et Supervision des Systèmes en Génie des Procédés -- Approche Bond Graphs, Mémoire à Diriger les Recherches N°H360, Université des Sciences et Technologies de Lille, HDR soutenue le 20 Décembre 2002.
- [Patton et Chen 1997] Patton R., Chen J., Observer based fault detection and isolation: robustness and applications. Control Engineering Practice, Vol. 5, No. 5, pp. 671-682, 1997.

- [Pessi et Luparia 1992] Pessi E, Luparia C., MORSAF : Un système expert temps réel d'aide aux opérateurs de salle de contrôle d'installation de distillation atmosphérique. Génie Logiciel Temps réel et systèmes à base de connaissances, Vol. 28, 1992.
- [Ploix et Follet 2001] Ploix S., Follet C., Fault diagnosis reasoning for set-membership approaches and application. Common Component Architecture, International Symposium on Intelligent Control (CCA, ISIC), Mexique, 2001.
- [Porté *et al.* 1988] Porté N., Boucheron S., Sallantin S., Arlabosse F., An algorithmic view at causal ordering. 2nd International Workshop on Qualitative Physics (QR), Paris, 1998.
- [Ploix, Gentil et Leseq 2003] Ploix S., Gentil S., Leseq S., Isolation decision for a multi-agent-based diagnostic system, 5th IFAC Symposium on Fault Detection, Supervision and Safety of Technical Processes (Safeprocess), Etats-Unis, 2003.
- [Promonet 2002] Promonet I., Amélioration du système d'aide à la conduite pour le pilote de FCC. Rapport de stage de fin d'étude, Analyse Industrielle et Informatique du Lycée La Martinière Terreaux, 2002.
- [Rasmussen 1993] Rasmussen J., Diagnostic Reasoning in Action. IEEE Transactions on Systems, Man, and Cybernetics, Vol. 23, No 4, pp. 981-992, 1993.
- [Reiter 1987] Reiter T., A theory of diagnosis from first principles. Artificial Intelligence Vol. 32, pp. 57-95.
- [Sanmarti, Puigjaner et Friedler 1998] Sanmarti E., Friedler F., Puigjaner L., Combinatorial technique for short term scheduling of multipurpose batch plants based on schedule-graph representation. Computers & Chemical Engineering, Vol. 22, pp. 847-850, 1998.
- [Sedda 1998] Sedda E., Estimation en ligne de l'état et des paramètres d'une machine asynchrone par filtrage à erreur bornée et par filtrage de Kalman. Thèse de doctorat de l'Ecole Nationale Supérieure de Cachan, 1998.
- [Schweppe 1968] Schweppe F.C., Recursive state estimation: unknown but bounded errors and system inputs. IEEE Transactions on Automatic Control, Vol. 13, pp. 22, 1968.
- [Thirion et Lesaffre 1999] Thirion C., Lesaffre , F.M., Mise en place de SACHEM, système d'aide à la conduite des hauts fourneaux : retour d'expériences et

-
- perspectives. 5th International Symposium on Artificial Intelligence, Hollande, 1999.
- [Taillibert 1998] Taillibert P., Various Improvements to Diagnosing with Temporal Bands. 9th Int. Workshop on Principles of Diagnosis (DX), Cape Cod, USA, 1998.
- [Travé-Massuès et Pons 1997] Travé-Massuyès L., Pons R., Causal ordering for multiple mode systems. 11th Int. Workshop on Qualitative Reasoning about Physical Systems, Italie, 1997.
- [Travé-Massuyès, Dague et Guerrin 1998] L. Travé-Massuyès, P. Dague, F. Guerrin. Le raisonnement qualitatif. Hermès, 1998.
- [Travé-Massuyès et Gentil 1999] Travé-Massuyès L., Gentil S., Artificial intelligence approaches for supervision and alarm interpretation in industrial environment. (ECC), Allemagne, 1999.
- [Travé-Massuyès *et al.* 2001] Travé-Massuyès L., Escobet T., Pons R., Tornil S., The Ca-En diagnosis system and its automatic modelling method. *Computation and Systems Journal*, Vol. 5, No. 2, 2001.
- [Villemeur 1988] Villemeur A., Sûreté de fonctionnement des systèmes industriels : fiabilité, facteur humain, informatisation. Eyrolles, 1988.
- [Walter et Pronzato 1994] Walter E., Pronzato L., Identification de modèles paramétriques à partir de données expérimentales. Masson, 1994.

Annexe A

Algorithmes de réduction et d'approximation

Représentation systématique du modèle causal structurel

Ce paragraphe propose une nomenclature qui simplifiera la présentation des algorithmes dans les paragraphes suivants. Cette nomenclature permet aussi d'identifier exactement la connaissance nécessaire à l'application de la méthode de modélisation causale.

Les instances V_i de l'ensemble V de variables sont exactement identifiées par :

$$V = \{V_i, i \in \{1..n_V\}\} \quad (A-1)$$

Informations nécessaires à la connaissance des modèles causaux

Neuf matrices de dimensions $n_V * n_V$ sont définies. Elles caractérisent exactement le MCS, le MCR et le MCA.

Soit la variable générique *Model* désignant le MCS ou le MCR ou le MCA. Soit la variable générique *Application* désignant l'application *Exist*, l'application *Support* ou l'application *Contribution*.

Les modèles causaux sont exactement définis par les neuf matrices *Tab_Application_Model* définies par la relation générique suivante :

$$\forall (i, j) \in \{1..n_V\}^2, Tab_Application_Model[i, j] = Application(I_{Model}(V_i, V_j)) \quad (A-2)$$

Si *Application* = *Exist*, alors *Tab_Application_Model* contient des booléens sinon (*Application* = *Support* ou *Contribution*) *Tab_Application_Model* contient le support des relations ou les contributions.

Finalement, une représentation plus concise est proposée. Le MCS, le MCR et le MCA sont exactement identifié par les 3 tableaux :

$$\begin{aligned} Model &= MCS, MCR, MCA \\ Tabs_Model &= \cup_{Application=\{Exist, Support, Contribution\}} Tab_Application_Model \end{aligned} \quad (A-3)$$

Information concernant les variables connues

Les deux applications *Capteur* et *Connues* précédemment définies (cf. formule 2-4 et section 2.3.3.5) sont respectivement représentées par les deux vecteurs *Tab_Capteurs* et *Tab_Connues* de dimension n_V^*1 , contenant respectivement les noms des capteurs et des booléens :

$$\forall i \in \{1 \dots n_V\}, \text{Tab_Capteurs}[i] = \text{Capteur}(V_i) \quad (A-4)$$

$$\forall i \in \{1 \dots n_V\}, \text{Tab_Connues}[i] = \text{Connues}(V_i) \quad (A-5)$$

Information concernant les phénomènes négligés

Deux tableaux qui contiennent la connaissance des phénomènes négligés sont finalement définis. Négliger l'influence d'une variable V_i sur une variable V_j revient à faire l'hypothèse que la variable V_i ne varie pas au cours du temps (au moins pas de manière significative). Par conséquent si $I(V_i, V_j)$ est négligée, alors l'ensemble des composants qui sous tendent une valeur constante pour V_i doit être précisément identifié par l'expert. V_i est la variable source de l'influence négligée.

Le tableau *Tab_Négligé* de dimension $n_V^*n_V$ contenant des valeurs booléennes indique quelles sont les influences du MCR à négliger :

$$\begin{aligned} \forall i, j \in \{1 \dots n_V\}^2, I(V_i, V_j) \text{ à négliger} &\Rightarrow \text{Tab_Négligé}[i, j] = \text{vrai} \\ \forall i, j \in \{1 \dots n_V\}^2, I(V_i, V_j) \text{ à conserver} &\Rightarrow \text{Tab_Négligé}[i, j] = \text{faux} \end{aligned} \quad (A-6)$$

Le dernier tableau, *Tab_Amont* $n_V^*n_V$ contient l'identité des composants qui sous tendent une valeur constante pour la variable V_i , source des influences négligées :

$$\begin{aligned} \forall i \in \{1 \dots n_V\}, \\ \text{SI } (V_i \text{ source d'une influence négligée}) \\ \text{ALORS } \text{Tab_Amont}(i) = \{\text{Composants qui sous tendent une valeur constante pour } V_i\} \\ \text{SINON } \text{Tab_Amont}(i) = \emptyset \end{aligned} \quad (A-7)$$

Algorithme principal

Cette section présente l'algorithme principal nommé *Algorithme_Principal* qui effectue successivement les opérations de réduction (via l'algorithme *Reduction_Principal*) et d'approximation (via l'algorithme *Approximation*).

La connaissance nécessaire à l'algorithme principal réside dans les variables globales (variables communes à tous les algorithmes) :

Tabs_MCS, *Tabs_MCR*, *Tabs_MCA*, *Tab_Connues*, *Tab_Capteurs*, *Tab_Négligé*, *Tab_Amont* et *nb_var=n_v*.

Initialement l'expert du procédé affecte les valeurs de *Tabs_MCS*, *Tab_Connues*, *Tab_Capteurs*, *Tab_Négligé*, *Tab_Amont* et *nb_var=n_v*.

```

Algorithme_Principal( )
BOUCLE PRINCIPALE
  Tabs_MCR=APPEL Reduction_Principal(Tabs_MCS,Tab_Connues)
  Tabs_MCA=APPEL Approximation(Tab_Négligé,Tab_Amont,Tabs_MCR,Tab_capteurs)
FIN

```

Algorithme A-1 : Algorithme_Principal

Algorithme de réduction

Ce paragraphe décrit les différents algorithmes nécessaires à la réduction. L'algorithme principal de réduction est *Reduction_Principal*.

➤ Algorithme *Reduction_Principal* (cf. Algorithme 2)

Entrées : le MCS (*Tabs_MCS*) et l'identité des variables connues (*Tab_Connues*).

Principe : les influences entre les variables endogènes connues et les variables exogènes (connues ou inconnues) sont extraites du MCS.

Sortie : le MCR (*Tabs_MCR*).

Algorithmes appelés : *Rech_Variable_Inconnue* et *Reduction*.

➤ Algorithme *Rech_Variable_Inconnue* (cf. Algorithme A-3).

Entrées : j , Tab_Exist et $Tab_Connues$.

Principe : cet algorithme détermine si une variable V_j est endogène inconnue.

Sortie : si V_j est endogène inconnue, la sortie est « vrai », sinon « faux ».

Algorithmes appelés : Aucun.

➤ Algorithme *Reduction* (cf. l'Algorithme A-4 et Figure A-1).

Entrées : j , Tab_Exist , $Tab_Support$, et $Tab_Relation$.

Principe : V_j est une variable inconnue et endogène. Cet algorithme transforme Tab_Exist , $Tab_Support$, et $Tab_Relation$. Pour chaque cause directe V_i de V_j et pour chaque conséquence directe V_z de V_j , une nouvelle influence est créée de V_i sur V_z (la relation est obtenue en effectuant des manipulations algébriques). Les influences de V_i sur V_j et de V_j sur V_z sont éliminées.

Sorties : Tab_Exist , $Tab_Support$, et $Tab_Relation$.

Algorithmes appelés : *Casser_Boucle* et *Eliminer*.

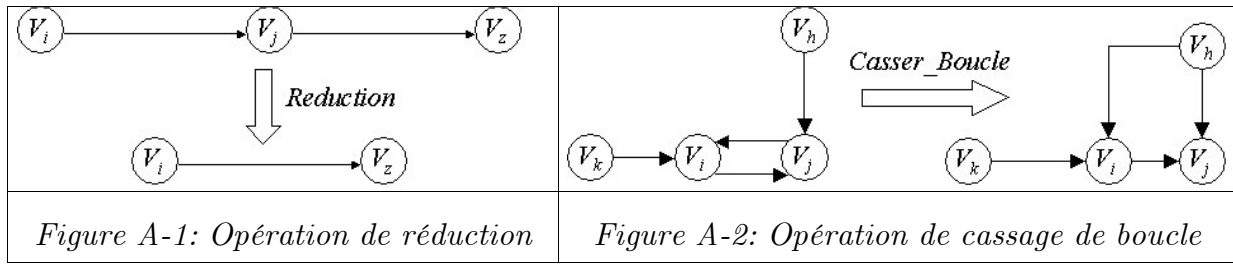
➤ Algorithme *Casser_Boucle* (cf. l'Algorithme A-5 et Figure A-2)

Entrées : j , Tab_Exist , $Tab_Support$ et Tab_Rel .

Principe : cet algorithme transforme les influences dans le cas particulier où il existe une variable endogène inconnue V_j qui agit sur elle même via une autre variable V_i . Ce cas particulier n'est pas considéré dans [Travé *et al.* 2001]. Cet algorithme transforme donc Tab_Exist , $Tab_Support$ et Tab_Rel . Des compositions de relations sont effectuées de telle manière que V_i est influencée par les causes directes de V_j et que l'influence de V_j sur V_i est éliminée. Cet algorithme ne modifie pas la causalité : les causes (directes et indirectes) et les conséquences de V_i et V_j ne changent pas. Si cet algorithme n'est pas inséré dans l'algorithme réduction, alors l'algorithme réduction ne converge jamais dans le cas particulier de la Figure A-2.

Sorties : Tab_Exist , $Tab_Support$, et Tab_Rel .

Algorithmes appelés : *Casser_Boucle* et *Eliminer*.



➤ Algorithme *Eliminer* (cf. Algorithme A-6)

Entrées : $i, j, Tab_Exist, Tab_Support$ et $Tab_Relation$.

Principe : l'influence de V_i sur V_j est éliminée et $Tab_Exist, Tab_Support$ et Tab_Rel sont donc transformés.

Sorties : $Tab_Exist, Tab_Support$ et Tab_Rel .

Algorithmes appelés : Aucun

```

Reduction_Principale(Tab_Exist, Tab_Support, Tab_Rel, Tab_Connues)
BOUCLE PRINCIPALE
POUR j:=1 A nb_var
    Var_inconnue_endo:=Rech_Variable_Inconnue(j, Tab_Exist, Tab_Connues)
    SI Var_inconnue_endo=vrai ALORS
        (Tab_Exist, Tab_Support, Tab_Rel) := APPEL Reduction(j, Tab_Exist, Tab_Support, Tab_Rel)
    FIN SI
FIN POUR j
RETOURNER (Tab_Exist, Tab_Support, Tab_Rel)

```

Algorithme 2 : Reduction_Principale

```

Rech_Variable_Inconnue(j, Tab_Exist, Tab_Connues)
BOUCLE PRINCIPALE
variable_endogene = faux
variable_inconnue_endogene := faux
POUR i := 1 A nb_var
    SI Tab_Exist[i, j] = vrai ALORS variable_endogene := vrai
    FIN SI
FIN POUR i
SI (variable_endogene = vrai) ET (Tab_Connues[j] = faux) ALORS
    variable_inconnue_endogene := vrai
FIN SI
RETOURNER(variable_inconnue_endogene)

```

Algorithme A-3 : Rech_Variable_Inconnue

```

Reduction(j,Tab_Exist,Tab_Support,Tab_Rel)
BOUCLE PRINCIPALE
POUR i :=1 A nb_var
  SI Exist[i,j]=vrai
    {Si  $V_i$  est une cause directe de  $V_j$  }
    (Tab_Exist,Tab_Support,Tab_Rel) :=APPEL Casser_Boucle(j,Tab_Exist,Tab_Support,Tab_Rel)
    POUR z :=1 A nb_var
      SI Exist[j,z]=vrai ALORS
        {Si  $V_j$  est une cause directe de  $V_z$  }
        Tab_Exist[i,z]:=vrai
        Tab_Rel[i,z]:=Tab_Rel[j,z] ◦ Tab_Rel[i,j]
        Tab_Support[i,z]:=Tab_Support[j,z] ∪ Tab_Support[i,j]
        {Création d'une influence de  $V_i$  sur  $V_z$  }
      FIN SI
    FIN POUR z
  FIN SI
FIN POUR i
POUR i :=1 A nb_var
  (Tab_Exist,Tab_Support,Tab_Rel) :=APPEL Eliminer(i,j,Tab_Exist,Tab_Support,Tab_Rel)
  (Tab_Exist,Tab_Support,Tab_Rel) :=APPEL Eliminer(j,i,Tab_Exist,Tab_Support,Tab_Rel)
FIN POUR i
RETOURNER (Tab_Exist,Tab_Support,Tab_Rel)

```

Algorithme A-4 : Reduction

```

Casser_boucle ( $j, Tab\_Exist, Tab\_Support, Tab\_Rel$ )
BOUCLE PRINCIPALE
POUR  $i := 1$  A  $nb\_var$ 
{Pour toute variable  $V_i \in E$  }
SI ( $Tab\_Exist[i, j] = vrai$ ) ET ( $Tab\_Exist[j, i] = vrai$ ) ALORS
  {Si  $V_i$  est une cause directe de  $V_j$  et si  $V_j$  est une cause directe de  $V_i$  }
  POUR  $k := 1$  A  $nb\_var$ 
    SI ( $Tab\_Exist[k, i] = vrai$ ) ET ( $k \neq j$ ) ALORS
      {Pour toutes les causes directes  $V_k$  de  $V$  }
       $Tab\_Rel[k, i] := (Id_{V_i} - Tab\_Rel[j, i] \circ Tab\_Rel[i, j])^{-1} \circ Tab\_Rel[k, i]$ 
       $Tab\_Support[k, i] := Tab\_Support[k, i] \cup Tab\_Support[j, i] \cup Tab\_Support[i, j]$ 
      {Modification de l'influence de  $V_k$  sur  $V$  par composition de relations}
    FIN SI
  FIN POUR  $k$ 
  POUR  $h := 1$  A  $nb\_var$ 
    SI ( $Tab\_Exist[h, j] = vrai$ ) ET ( $h \neq j$ ) ALORS
      {Pour  $V_h \in E \setminus \{V_j\}$ , cause directe de  $V_j$  }
       $Tab\_Rel[h, i] := (Id_{V_i} - Tab\_Rel[j, i] \circ Tab\_Rel[i, j])^{-1} \circ Tab\_Rel[j, i] \circ Tab\_Rel[h, j]$ 
       $Tab\_Support[h, i] := Tab\_Support[j, i] \cup Tab\_Support[h, j] \cup Tab\_Support[i, j]$ 
      {Creation d'une influence de  $V_h$  sur  $V_i$  par composition de relations}
    FIN SI
  FIN POUR  $h$ 
  FIN SI
  ( $Tab\_Exist, Tab\_Support, Tab\_Rel$ ) := APPEL Eliminate( $j, i, Tab\_Exist, Tab\_Support, Tab\_Rel$ )
  {Elimination de l'influence  $I(V_j, V_i)$  }
  FIN POUR  $i$ 
RETOURNER  $Tab\_Exist, Tab\_Support, Tab\_Rel$ 

```

Algorithme A-5: *Casser_Boucle*

```

Eliminer( $i, j, Tab\_Exist, Tab\_Support, Tab\_Rel$ )
BOUCLE PRINCIPALE
   $Tab\_Exist[i, j] := faux$ 
   $Tab\_Support[i, j] := ""$ 
   $Tab\_Rel[i, j] := ""$ 
RETOURNER( $Tab\_Exist, Tab\_Support, Tab\_Rel$ )

```

Algorithme A-6 : *Eliminer*

Algorithme d'approximation

Ce paragraphe présente et illustre l'algorithme d'approximation. Il est important de remarquer que l'algorithme d'approximation doit notamment et obligatoirement être appliqué aux variables exogènes non connues.

➤ Algorithme : *Approximation* (cf. l'Algorithme A-7 et Figure A-3)

Entrées : *Tab_Négligé*, *Tab_Amont*, *Tab_MCR* et *Tab_Capteurs*.

Principe : Si l'influence de V_k sur V_j est à négliger alors :

- $Support(V_k, V_j)$ est ajouté à $Comps(V_i, V_j)$ et
- $Tab_Amont[k]$ est ajouté à $Comps(V_i, V_j)$ et
- $Capteurs(V_i)$ est retiré de $Comps(V_i, V_j)$ et
- $Relation(V_k, V_j)$ est éliminée.

Sortie : le MCA (*Tab_MCA*).

Algorithmes appelés : *Eliminer*.

```

Approximation(Tab_Négligé, Tab_Amont, Tab_Exist, Tab_Support, Tab_Rel, Tab_Capteurs)
BOUCLE PRINCIPALE
POUR  $i := 1$  A  $nb\_var$ 
  POUR  $j := 1$  A  $nb\_var$ 
    SI  $Tab\_Negligé[k, j] = vrai$  ALORS
      {Si l'influence de  $V_k$  sur  $V_j$  est négligeable}
       $Tab\_Support[i, j] = Tab\_Support[k, j] \cup Tab\_Support[i, j]$ 
      {Ajouter le composant  $Comps(V_k, V_j)$  à  $Comps(V_i, V_j)$ }
       $Tab\_Support[i, j] = Tab\_Support[i, j] \cup Tab\_Amont[k]$ 
      {Add components that determine the constant value of  $V_k$ }
       $Tab\_Support[i, j] = Tab\_Support[i, j] \setminus Tab\_Capteurs[k]$ 
      {Si  $V_k$  est mesuré alors éliminer  $V_k^S$  de  $Comps(V_i, V_j)$ }
      ( $Tab\_Exist, Tab\_Support, Tab\_Rel$ ) := APPEL Eliminer( $k, j, Tab\_Exist, Tab\_Support, Tab\_Rel$ )
      {Eliminer l'influence de  $V_k$  sur  $V$ }
    FIN SI
  FIN POUR  $j$ 
FIN POUR  $i$ 
RETOURNER( $Tab\_Exist, Tab\_Support, Tab\_rel$ )

```

Algorithme A-7 : *Approximation*

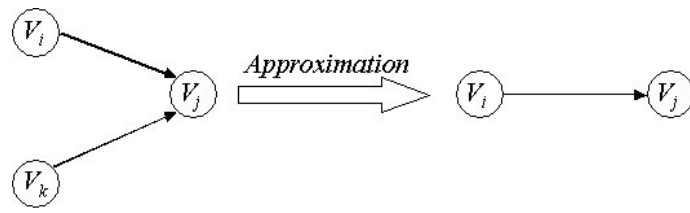


Figure A-3: Illustration de l'algorithme d'Approximation

Annexe B

Identification des paramètres

Introduction

Deux types d'erreurs sont inévitables lors de l'estimation des paramètres à partir de données expérimentales. La première correspond aux erreurs de mesure et autres perturbations qui ont pour conséquences de rendre les données incertaines. La seconde correspond aux erreurs structurelles, qui traduisent que le modèle utilisé est toujours une approximation de la réalité.

Ces erreurs doivent être prises en compte lors de l'élaboration du modèle. Dans le cadre d'une approche stochastique, la confiance à apporter aux paramètres est envisagée de manière statistique. Dans le cadre de l'approche ensembliste, des ensembles englobants qui garantissent de contenir les paramètres sont proposés.

Les méthodes d'identification stochastiques étant usuelles, l'objectif de cette section est surtout d'évaluer l'apport des techniques d'identification ensembliste et en particulier celles basées sur des ellipsoïdes. En effet, la méthode de détection proposée par [Armengol 1999] suppose que l'on connaît chaque paramètre du modèle sous la forme d'un intervalle. Comme le modèle que nous utilisons est discret, il est assez difficile de relier les paramètres à des paramètres physiques dont on pourrait, pour des raisons matérielles, connaître l'ensemble de définition. Il nous a donc paru nécessaire d'utiliser une méthode d'estimation ensembliste qui nous apporte ce que nous cherchons : un modèle dont les paramètres sont des intervalles.

Cette annexe présente :

- les grandes familles de méthodes d'identification paramétrique,
- des algorithmes d'identification ensembliste ellipsoïdale,
- une étude quantitative de l'apport de cette méthode.

Méthodes d'identification

Introduction

Certaines méthodes d'identification s'attachent à identifier à la fois un modèle du procédé et un modèle des bruits ou des perturbations environnantes. Ce document se limite aux méthodes dont l'unique objectif est la détermination du modèle du procédé.

Dans ce cadre sont distinguées les méthodes basées sur l'erreur d'équation et sur l'erreur de sortie. Pour illustrer ces méthodes dans la suite de cette annexe, la relation entre une entrée u_k et une sortie y_k est donnée par :

$$y_k = -(a_1 y_{k-1} + \dots + a_n y_{k-n}) + (b_0 u_{k-d} + \dots + b_m u_{k-d-m}) \quad (B-1)$$

Les a_i et b_i constituent les paramètres du modèle. Les entiers d , m et n représentent respectivement le retard, le nombre de zéros et le nombre de pôles. Les deux sections suivantes rappellent deux grandes méthodes d'identification : la méthode d'erreur d'équation et la méthode d'erreur de sortie.

Méthodes basées sur l'erreur d'équation

Le principe de la méthode est illustré par la Figure B-1. La sortie mesurée est comparée à la sortie prédite au même instant selon :

$$y_k = -(a_1 y_{k-1} + \dots + a_n y_{k-n}) + (b_0 u_{k-d} + \dots + b_m u_{k-d-m}) + e_t \quad (B-2)$$

(où y et u sont des mesures)

Le problème ainsi posé est linéaire par rapport aux paramètres car il s'écrit sous la forme :

$$y_k = \Phi_k \cdot \theta + e_k \quad (B-3)$$

où le vecteur Φ_k est parfaitement connu, $\Phi_k^T = (y_{k-1}, \dots, y_{k-n}, u_{k-d}, \dots, u_{k-d-m})$, et $\theta = (-a_1, \dots, -a_n, b_0, \dots, b_m)$ est le vecteur des paramètres. Un traitement global peut être envisagé hors ligne : cette démarche consiste à étudier toutes les données en même temps. Un traitement séquentiel, par rapport au temps, peut aussi être envisagé :

cette démarche «récursive» peut être appliquée en ligne. Ce traitement est donc particulièrement avantageux.

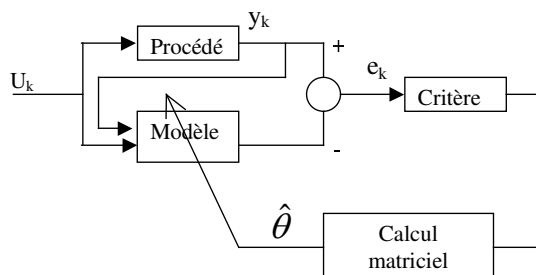


Figure B-1 : Méthode de l'erreur d'équation

Dans le cadre stochastique, la procédure d'identification repose généralement sur la minimisation du critère quadratique $J(\theta) = \sum_k e_k^2$.

La minimisation de ce critère conduit aux méthodes des moindres carrés et de ses dérivées (Moindres Carrés Etendus, Moindre Carrés Généralisés, ...). Toutes ces méthodes présentent l'avantage de pouvoir être utilisées en ligne et permettent d'obtenir une forme explicite de la solution. Avec une méthode de moindres carrés simples, l'estimation des paramètres est biaisée mais ne l'est pas avec les moindres carrés étendus.

Dans le cadre ensembliste, [Milanese et Belforte 1982], le problème d'identification avec une méthode basée sur l'erreur d'équation peut s'exprimer sous la forme de la recherche de l'espace paramétrique Θ vérifiant :

$$\Theta = \left\{ \theta \in \mathcal{R}^{n+m+1} / |y_k - \Phi_k \cdot \theta| \leq \delta_k \right\} \quad (B-4)$$

où δ_k reflète le niveau de bruit la mesure.

Méthodes basées sur l'erreur de sortie

Le principe de la méthode est illustré par la Figure B-2. La sortie du modèle η_t est comparée à la sortie du procédé au même instant selon :

$$\begin{aligned} \eta_k &= -(a_1 \eta_{k-1} + \dots + a_n \eta_{k-n}) + (b_0 u_{t-k} + \dots + b_m u_{t-k-m}) \\ \varepsilon_k &= \eta_k - y_k \end{aligned} \quad (B-5)$$

Le problème ainsi posé est plus représentatif de la réalité mais il est non linéaire par rapport aux paramètres. Il est généralement résolu par des méthodes de programmation non linéaires PNL qui impliquent un traitement hors ligne.

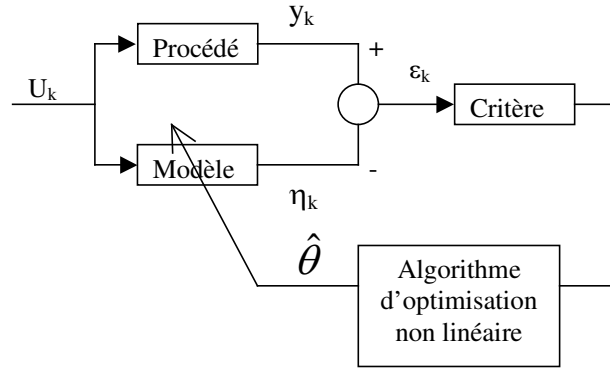


Figure B-2 : Méthode de l'erreur de sortie

Dans le cadre stochastique, la méthode dite du modèle consiste à minimiser le critère quadratique $J(\theta) = \sum_k \varepsilon_k^2$. D'autres méthodes comme la méthode de maximum de vraisemblance s'attachent à identifier un modèle du procédé et du bruit. Ces approches sont fondamentalement sans biais. La minimisation du critère peut s'effectuer avec des méthodes bien connues (méthode du gradient, de Newton, de Levenberg et Marquart ...). [Walter et Pronzato 1994] propose une revue des techniques d'identifications.

Dans le cadre ensembliste, le problème d'identification par la méthode du modèle peut s'exprimer sous la forme de la recherche de l'espace paramétrique Θ vérifiant :

$$\Theta = \left\{ \theta \in \mathfrak{R}^{n+m+1} \mid |\eta_k - y_k| \leq \varepsilon_k^M \right\} \quad (B-6)$$

où ε_k^M reflète le bruit de la mesure à l'instant k et l'imprécision du modèle.

Sous l'hypothèse restrictive de la connaissance à-priori des signes des paramètres autorégressifs a_1, \dots, a_n , la recherche des intervalles paramétriques, correspondants à une borne ε_k^M peut être menée avec une technique similaire à celles proposées pour l'erreur d'équation.

Identification ensembliste ellipsoïdale

Introduction

Les méthodes d'identification paramétriques permettent une estimation à partir de signaux d'entrée quelconques, dans la mesure où les entrées sont suffisamment riches pour permettre une convergence de l'estimation. Les erreurs peuvent être envisagées de deux manières différentes.

- Dans le cadre statistique, les données sont supposées entachées par des erreurs pouvant être caractérisées par un vecteur de densité de probabilité connue ou paramétrée. La qualité de l'estimée obtenue est le plus souvent exploitée en utilisant les propriétés de la matrice d'information de Fisher [Demerle et Siguerdidjane 2002]. Ces méthodes sont cependant mal adaptées dans le cas où la seule information disponible pour les erreurs se présente sous la forme de bornes.
- Dans le cadre ensembliste, l'erreur est inconnue mais supposée bornée : l'approche est dite à erreur bornée. L'objectif n'est plus de s'intéresser à une valeur privilégiée du paramètre mais consiste à chercher l'ensemble des valeurs des vecteurs paramètre qui sont acceptables pour que les erreurs associées restent comprises entre des bornes données à-priori. Cette idée a été introduite par [Schweppe 1968].

Ce chapitre présente

- les grands principes de l'identification ensembliste et en particulier ceux des méthodes basées sur des ellipsoïdes.
- L'algorithme OBE d'identification ellipsoïdale.

Grands principes de l'identification ensembliste ellipsoïdale

L'objectif des méthodes d'identification ensemblistes est d'obtenir une partie englobante de Θ compatible^{B-1} avec (B-4) ou (B-6). Cet ensemble englobant garantit de contenir les paramètres, il explique l'ensemble des mesures, compte tenu le bruit et l'incertitude du modèle.

Les relations (B-4) et (B-6) peuvent s'écrire simplement sous la forme $\Theta \in H_k$ où H_k est le volume compris entre deux hyperplans.

Dans un espace à deux dimensions, chaque mesure permet de construire deux droites (hyperplans).

L'espace des paramètres solution de (B-4) est l'intersection des surfaces H_k , comprises entre deux droites parallèles, selon la Figure B-3 :

$$\Theta = \bigcap H_k \quad (B-7)$$

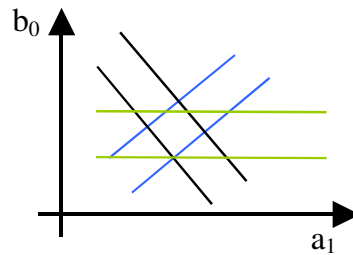


Figure B-3 : Espace des paramètres Θ compatibles avec (B-6)

Dans un espace à deux dimensions, l'espace des paramètres solution de (B-6) s'exprime comme l'intersection de droites non parallèles S_k .

^{B-1} dans le cadre de l'erreur d'équation

L'ensemble englobant peut se présenter sous la forme d'ellipsoïdes, d'orthotropes, de parallélotopes ou encore de polytopes.

- Les ellipsoïdes sont rapides en terme de temps de calcul et donnent des résultats satisfaisants.
- Les polytopes coûteux en calculs mais donnent de meilleurs résultats.

Ces techniques sont généralement appliquées à l'erreur d'équation, mais de travaux originaux ont été effectués dans le cadre de l'erreur de sortie [Clément 1984]. [Leseq, Barraud et Tran-Dinh 2003] ont ensuite apporté d'importantes améliorations à ces travaux. En particulier ils proposent une forme factorisée qui améliore nettement la convergence de l'algorithme. Nous allons brièvement décrire et illustrer ces travaux dans la suite de cette annexe.

Une ellipse est définie par sa matrice directrice P_k nécessairement définie positive (elle caractérise sa taille et son orientation), et son centre c_k selon la formule générale :

$$E_k = \left\{ x \in \mathfrak{R}^{n+m+1} : (x - c_k)^T P_k^{-1} (x - c_k) \leq \sigma_k^2 \right\} \quad (B-8)$$

L'algorithme d'identification peut être est récursif : un ellipsoïde E_k est exprimé en fonction de l'ellipsoïde précédent E_{k-1} et du volume H_k de telle manière que E_k est le plus petit ellipsoïde possible contenant $E_{k-1} \cap H_k$. L'ellipse initiale, E_0 , est choisi pour contenir Θ à coup sûr, donc très grande.

La Figure B-4 illustre les ensembles E_k et H_k dans le cadre de l'erreur d'équation et dans le cadre de l'erreur de sortie.

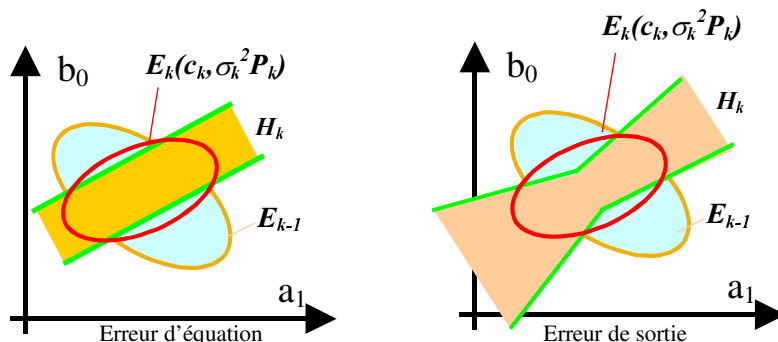
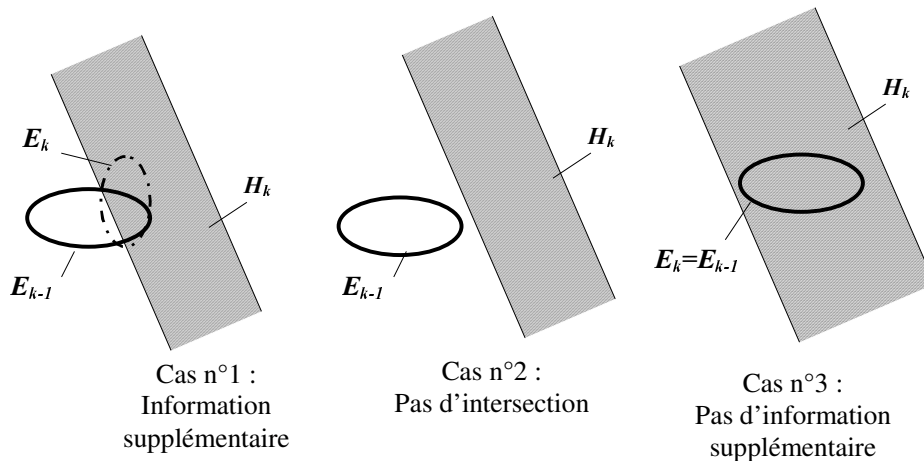


Figure B-4 : Principe récursif de génération des ellipsoïdes

Pour chaque mesure (qui définit deux hyperplans) trois cas peuvent se présenter (cf. FigureB-5) :

- Si le volume H_k intersecte l'ellipsoïde E_{k-1} alors il apporte de l'information. Les paramètres situés dans le sous espace intersection sont donc conservés. Les paramètres situés à l'extérieur sont rejetés. Ce cas est représenté par le cas n°1 de la FigureB-5. Un nouvel ellipsoïde E_k , le plus petit possible, passant par les deux points à l'intersection de H_k et E_{k-1} est construit.
- Le volume H_k n'intersecte pas l'ellipsoïde initial. Ce cas se présente si la mesure utilisée pour générer H_k est aberrante. Autrement dit la mesure et son niveau de bruit ne sont pas cohérent avec les paramètres précédemment estimés. Cet exemple est illustré par le cas n°2 de la FigureB-5.
- Le volume H_k contient l'ellipsoïde E_{k-1} . Aucune information supplémentaire n'est apportée par la mesure. Cet exemple est illustré par le cas n°3 de la FigureB-5.



FigureB-5 : Méthode itérative intersection du volume et de l'ellipsoïde

Différents critères existent pour caractériser la taille de la matrice. Les critères géométriques du déterminant ou de la trace sont souvent utilisés. Le critère du déterminant est proportionnel au volume au carré de l'ellipsoïde et la trace au carré des axes principaux. L'utilisation du déterminant peut conduire dans certains cas à la génération d'ellipsoïdes très allongés. L'utilisation de la trace peut conduire à des ellipsoïdes sphériques. Par exemple, en dimension 2, tout point $x=(x_1;x_2)^T$ appartenant à un ellipsoïde de centre $(c;c)^T$ et définie dans son repère propre par :

$$\frac{(x_1 - c)}{a_2} + \frac{(x_2 - c)}{b^2} \leq 1 \quad (B-9)$$

La Figure B-6 illustre un ellipsoïde dans un espace à 2 dimensions.

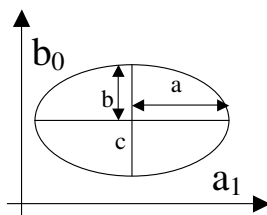


Figure B-6 : Illustration d'un ellipsoïde en dimension 2

Les algorithmes OBE au cœur de l'identification ensembliste ellipsoïdale

Dans la suite de cette annexe, nous posons les notations suivantes :

$$\begin{cases} y_k = d_k^T x + e_k \\ |e_k| \leq \delta_k \end{cases} \quad d_k, x \in \mathbb{R}^n \quad (B-10)$$

Les algorithmes OBE (Optimal Bounding Ellipsoid) proposent de déduire l'ellipsoïde E_k directement et simplement par combinaison linéaire de l'ellipsoïde E_{k-1} et du volume H_k :

$$E_k \supset E_{k-1} \cap H_k \Rightarrow E_k = \alpha_k \cdot E_{k-1} + \beta_k \cdot H_k \quad (B-11)$$

Les coefficients α_k et β_k permettent respectivement de pondérer l'information cumulée depuis l'instant initial, contenue dans E_{k-1} , et l'information apportée à l'instant k contenu dans le volume H_k . En général $\alpha_k + \beta_k = 1$. La relation (B-11) se traduit par les égalités suivantes :

$$\begin{aligned} P_k &= \frac{1}{\alpha_k} \left[P_{k-1} - \frac{\beta_k P_{k-1} d_k d_k^T P_{k-1}}{\alpha_k + \beta_k d_k^T P_{k-1} d_k} \right] \\ c_k &= c_{k-1} + \beta_k P_k d_k (y_k - d_k^T c_{k-1}) \\ \sigma_k^2 &= \alpha_k \sigma_{k-1}^2 + \beta_k \delta_k^2 - \frac{\alpha_k \beta_k (y_k - d_k^T c_{k-1})^2}{\alpha_k + \beta_k d_k^T P_{k-1} d_k} \end{aligned} \quad (B-12)$$

Les algorithmes OBE consistent à chercher un paramètre λ_k qui minimise un critère et qui permet de déterminer les deux paramètres $\alpha_k(\lambda_k)$ et $\beta_k(\lambda_k)$. Deux groupes d'algorithme OBE sont distingués.

- Le premier groupe s'attache à minimiser la taille géométrique de E_k . Cette taille est caractérisée soit par le critère du déterminant de P_k soit par le critère de la trace de $\sigma_k^2 \cdot P_k$. Dans ce cadre, l'algorithme de Fogel et Huang, [Leseq, Barraud et Tran-Dinh 2003] constitue le cœur de l'algorithme propose de chercher $\lambda_k > 0$ avec $\alpha_k = (\sigma_{k-1}^2)^{-1}$, $\beta_k = \lambda_k \cdot (\delta_k^2)^{-1}$. [Sedda 1998] propose de chercher $\lambda_k \in]0;1[$ avec $\alpha_k = \lambda_k \cdot (\sigma_{k-1}^2)^{-1}$, $\beta_k = (1 - \lambda_k) / \delta_k^2$.
- Le second consiste à faire converger l'argument σ_k^2 . Par exemple, l'algorithme de Arruda et Favier propose $\alpha_k = \beta_k = \lambda_k$, sans contrainte sur λ_k qui minimise σ_k^2 .

Les premiers algorithmes implémentés dans [Clément 1984] n'étaient pas numériquement stables : ils ne garantissaient pas de conserver P_k définie positive. Des formulations factorisées ^{B-2} de l'algorithme OBE ont été employées par [Leseq, Barraud et Tran-Dinh 2003] pour palier à ce problème : elle garantissent que P_k reste définie positive.

Conclusion

Les algorithmes d'identification ensembliste que nous avons utilisés sont basés sur l'algorithme OBE qui propose une forme explicite et permet de calculer récursivement la solution : un traitement séquentiel (par rapport au temps) est effectué contrairement aux méthodes proposant un traitement global de toutes les données à la fois.

Il est nécessaire d'utiliser une forme factorisée de l'algorithme OBE pour garantir sa stabilité en ligne. Le problème majeur de cet algorithme est le choix de la borne δ sur le bruit.

^{B-2} Factorisation d'une matrice (Cholesky) : une matrice M qui est symétrique et positive peut s'écrire sous la forme $M=LL^T$ où L est une matrice triangulaire inférieure.

Pour toutes ces méthodes il est nécessaire de spécifier plusieurs caractéristiques.

- L'ellipsoïde initial E_0 , qui peut être obtenu par une méthode d'identification classique (méthode du modèle et un critère quadratique). L'écart type sur les paramètres peut être utilisé pour caractériser la taille de E_0 ,
- La borne initiale δ sur la sortie qui quantifie le bruit de mesure. En simulation, la borne δ est fixée comme la borne supérieur de l'écart entre la sortie réelle bruitée $y_{\text{réelle}}$ et de la sortie du modèle non bruitée $y_{\text{modèle}}$:

$$\delta = |y_{\text{réelle}} - y_{\text{modèle}}| \quad (\text{B-13})$$

Si aucun modèle n'est disponible, la borne δ est obtenue lorsque aucune entrée n'est excitée par comparaison de la mesure y_{mes} et du point de fonctionnement y :

$$\delta = |y_{\text{mes}} - y| \quad (\text{B-14})$$

L'algorithme peut être appliqué successivement plusieurs fois sur le même jeu de données. Le traitement séquentiel (par rapport au temps) d'un ensemble de données est une **circulation**. Il est possible de traiter plusieurs fois un ensemble de données ; on parle alors de **re-circulation**.

L'ellipsoïde initial d'une circulation étant celui finalement obtenu à la circulation précédente. A chaque re-circulation tous les points sont interprétés et certains points permettent d'affiner l'ellipsoïde. Une mesure peut apporter de l'information lors de plusieurs itérations successives car le nouvel ellipsoïde est généralement sur évalué (cf. cas n°1 de la FigureB-5).

Nous testons les algorithmes présentés dans cette section dans la section suivante.

Application de la méthode d'estimation

Introduction

Cette section présente des tests que nous avons effectués pour évaluer les performances de la méthode d'identification ellipsoïdale.

Nous avons choisi la fonction de transfert $H(z)$ ^{B-3} de la Figure B-7. pour tester la méthode. La période d'échantillonnage est 5 secondes. Nous avons utilisé l'algorithme proposé par [Leseq, Barraud et Tran-Dinh 2003] avec la forme factorisée et avons envisagé plusieurs cas :

- Nous avons ajouté un bruit additif e à la sortie y . Ce bruit peut être blanc gaussien ou uniforme (borné).
- Nous avons testé différentes excitations de l'entrée u . Les signaux d'entrées étudiés sont un échelon puis une SBPA^{B-4}.
- Nous avons testé les deux formes de modèle (B-2) (erreur d'équation) ou (B-5) (erreur de sortie).

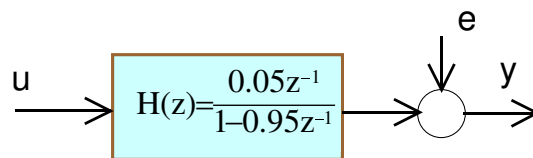


Figure B-7 : Fonction de transfert utilisée pour tester la méthode d'identification

Influence de la forme du bruit et de la forme du modèle

Dans un premier temps nous avons étudié l'influence du bruit et de la forme du modèle.

Un échelon unitaire est appliqué dans l'intervalle $t = 5$ sec. à $t = 200$ sec. La durée de la simulation est de 400 secondes. L'énergie du bruit e , calculée selon la formule (B-15) vaut 92 dB. 10 re-circulations sont effectuées.

^{B-3} La dynamique de $H(z)$ est représentative du comportement du pilote de FCC.

^{B-4} Séquence Binaire Pseudo Aléatoire

$$\frac{S}{B} = 20 * \log_{10} \left(\frac{\sum_{k=0}^{k=n} (Y_{\text{modèle}}(t_0 + k.Td) - \bar{Y}_{\text{modèle}})^2}{\sum_{k=0}^{k=n} (e(t_0 + k.Td) - \bar{e})^2} \right) \quad (B-15)$$

Dans le cas du bruit blanc gaussien $\delta = 0.14$ (δ est obtenu selon (B-13)) et dans le cas du bruit uniforme $\delta = 0.069$.

La Figure B-8 illustre les résultats obtenus.

- Les ellipsoïdes rouge et magenta sont obtenus avec un bruit blanc gaussien et avec respectivement la méthode d'erreur d'équation et de sortie,
- Les ellipsoïdes noir et vert sont respectivement obtenus avec un bruit uniforme et avec respectivement la méthode d'erreur d'équation et de sortie.

Ces tests illustrent qu'à niveau d'énergie identique, le bruit blanc gaussien est plus pénalisant pour l'identification que le bruit uniforme. Le bruit uniforme est borné, donc il est plus adapté à l'estimation ensembliste que le bruit blanc gaussien qui peut prendre ponctuellement des valeurs importantes.

La méthode par erreur de sortie est toujours plus pessimiste que la méthode d'erreur d'équation à niveau de bruit fixé. Pour la méthode par erreur de sortie, il est nécessaire de construire un volume H_k plus grand que celui qui est construit à partir de l'erreur d'équation [Clément 1984]. La méthode d'erreur de sortie est cependant plus représentative de la réalité que la méthode d'erreur d'équation.

La Figure B-8 illustre que tous les ellipsoïdes sont allongés selon l'axe $b_0=1-a_1$ qui définit les paramètres des fonctions de transfert dont le gain vaut 1.

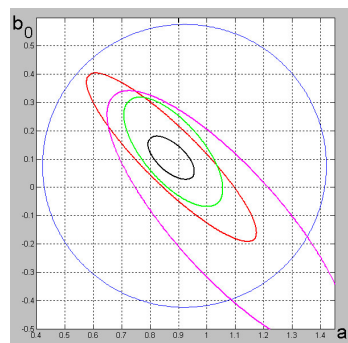


Figure B-8 : Influence de la forme du bruit et du modèle

Nous avons ensuite étudié l'influence de l'excitation de l'entrée. La durée de la simulation est de 300 secondes. Le bruit est uniforme et $\delta = 0.005$. Le modèle est donné par la forme (B-2) (erreur d'équation). Les signaux d'entrées étudiés sont représentés sur la Figure B-9. La première entrée est un échelon effectué à l'instant initial et la seconde est une SBPA. Les ellipsoïdes obtenus pour ces deux excitations sont représentés sur la Figure B-9 :

- en bleu pour l'échelon,
- en vert pour la SBPA.

L'ellipsoïde obtenu avec l'échelon unique est très allongé, mais il est nettement plus petit (selon la projection sur les axes) que l'ellipsoïde noir de la Figure B-8 car δ est plus petit.

L'autre ellipsoïde est plus petit et n'est plus incliné selon l'axe $b_0=1-a_1$: il permet de mieux caractériser les paramètres.

Le pourcentage de données utilisées pour estimer les paramètres est 14,3% pour l'échelon et 17,3%. Le pourcentage de données utilisées augmente avec la précision du résultat mais est faible.

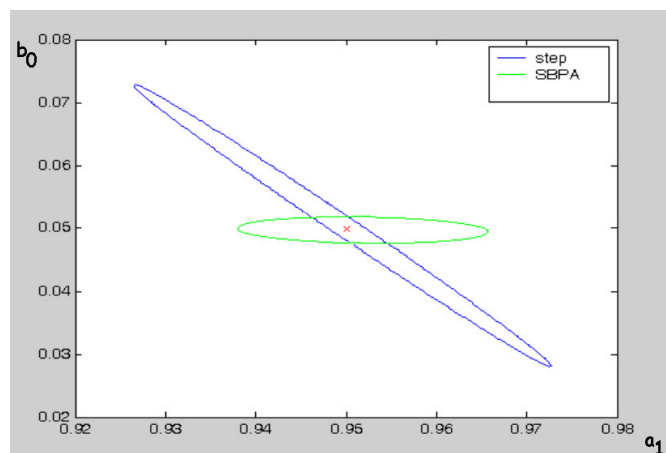


Figure B-9 : Ellipsoïdes obtenus pour différentes excitations du système

L'entrée appliquée au système a une influence importante sur le résultat obtenu. En pratique, sur notre procédé, il est cependant impossible d'effectuer une SBPA.

Nous avons testé cette méthode d'identification sur des données du pilote de FCC (cf. Figure B-10), mais les données n'étaient pas suffisamment excitées pour obtenir de bons résultats. L'entrée (en bleu) est une ouverture de vanne (en %) et la sortie (en rose) est un débit. Dans la plage de fonctionnement qui est représentative des variations des deux variables, la relation est linéaire et nous cherchons à l'identifier par une fonction de transfert du premier ordre.

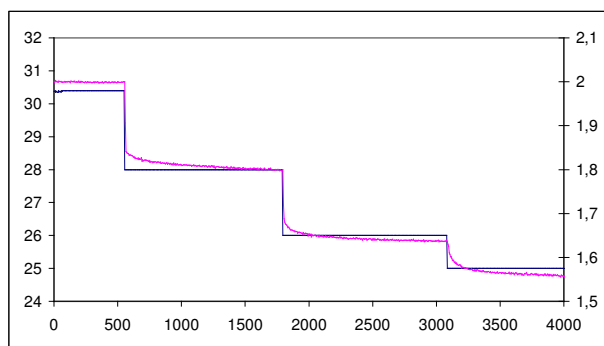


Figure B-10 : Des données du pilote de FCC

Nous avons appliqué la méthode de Foguel et Huang sur les données de la Figure B-10, avec la forme factorisée et la méthode d'erreur d'équation. Nous avons effectué 10 re-circulations et avons fixé $\delta = 0.0032$. Nous observons que l'ellipse est très allongée.

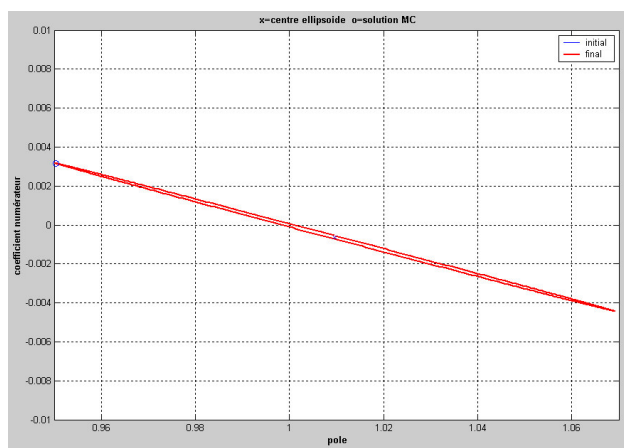


Figure B-11 : Résultats obtenus sur les données de la Figure B-11

La méthode d'identification présentée dans ce document fait actuellement l'objet d'importantes améliorations [Leseq, Barraud et Tran-Dinh 2003]. Les résultats présentés dans cette annexe ont été effectués sans ces améliorations. Ces travaux qui sont effectués en dehors du cadre de ce manuscrit ont donné de bien meilleurs résultats sur les données du pilote de FCC mais ne sont pas présentés.

Conclusion

Cette annexe a présenté une méthode d'identification basée sur des ellipsoïdes : le résultat est un ellipsoïde englobant les paramètres identifiés du modèle. Les algorithmes OBE sont classiques mais pour les utiliser, les données doivent avoir des propriétés :

- La technique est mieux adaptée à un bruit borné qu'à un bruit Gaussien,
- La technique est adaptée aux signaux peu bruités,
- Il est nécessaire de disposer de données dont les entrées sont bien excitées pour obtenir un résultat satisfaisant
- L'estimation par les ellipsoïdes n'est certainement pas la méthode d'identification ensembliste la plus adaptée à la méthode de diagnostic choisie, car la projection sur chaque axe de l'ellipsoïde entraîne un pessimisme supplémentaire sur les paramètres.

Annexe C

Ordre de calcul des nœuds

Les contraintes d'ordre de calcul sont engendrées par les arcs de retard nul. En effet, si une variable est influencée par un arc de retard non nul, alors la contribution de cet arc peut être évaluée aisément car elle ne dépend que de l'historique. Si le retard de l'arc est nul, alors la contribution de cet arc dépend de la valeur instantanée de l'arc amont qu'il est nécessaire de calculer au préalable.

Algorithmme

- **Première étape : Pré-traitement**

Le pré-traitement est utile dans le cas où les arcs influençant une même variable ont des retards différents dont un, au moins, est nul. Le pré-traitement consiste à considérer artificiellement que le retard de tous ces arcs est nul. Le retard nul étant le plus contraignant vis à vis de l'ordre de calcul. Sur la Figure C-1, le retard entre X_2 et X_5 est nul et celui entre X_3 et X_5 vaut 1. Le retard entre X_3 et X_5 est donc considéré artificiellement nul.

- **Seconde étape : Traitement des arcs de retard non nul**

Cette étape consiste à éliminer les arcs dont le retard est différent de 0 et les nœuds influencés par ces arcs. Ces nœuds sont dits de «contrainte» nulle. La Figure C-1 illustre cette étape.

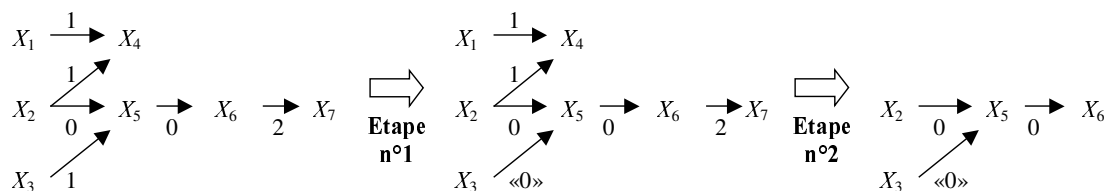


Figure C-1 : Illustration de l'ordre de calcul des nœuds

- **Troisième étape : Traitement des arcs de retard nul**

Cette étape consiste à calculer les contraintes des nœuds restants. La valeur de la contrainte est le nombre de nœuds entre la variable et sa variable cause de plus grande profondeur. Ainsi la contrainte de X_6 vaut 2, celle de X_5 vaut 1 et celles de X_2 et X_3 valent 0.

- **Troisième étape : Etablissement de l'ordre de calcul**

Ce rapport est basé sur les travaux de J.A.J Molina [Molina 2000]. Les valeurs de références globales sont calculées selon leurs contraintes croissantes.

Contrainte de valeur 0 : $X_1, X_2, X_3, X_4, X_7,$

Contrainte de valeur 1 : $X_5,$

Contrainte de valeur 2 : $X_6.$

L'ordre de calcul des nœuds de niveau de contrainte identique est quelconque.

Dans le cas où il existe une boucle de rétroaction dont tous les retards sont nuls, cet algorithme ne permet pas de déterminer l'ordre de calcul des nœuds : les contraintes des nœuds sont infinies. Ce cas particulier peut cependant être traité numériquement en considérant que toutes les variables sont mutuellement dépendantes et en calculant simultanément leurs valeurs. Par contre, la présence d'un retard non nul dans la boucle est suffisante pour déterminer les valeurs contraintes de tous ses nœuds.

RESUME

La thèse traite du diagnostic de procédés et se divise en trois parties.

La première partie présente une méthode de modélisation causale pour le diagnostic.

La seconde partie présente les différentes méthodes mises en œuvre pour la détection, la localisation et l'identification de défauts. Pour effectuer la détection, nous comparons les variables du modèle causal aux mesures. Nous évaluons les apports d'une approche ensembliste et d'une approche par logique floue pour la prise en compte des incertitudes. La localisation est réalisée grâce à un algorithme de hitting-sets afin de déterminer l'ensemble des composants suspects. Ceci permet de focaliser l'identification de défauts sur un nombre réduit de composants. Cette identification utilise l'expertise des exploitants sur les modes de défaillance, exprimée sous forme de bases de règles, chacune relative à un composant, et permettant d'aiguiller l'opérateur sur des actions et des vérifications à effectuer.

La troisième partie présente l'application. Le système de diagnostic informatique ASCO (Aide à la Supervision et à la Conduite pour les Opérateurs) a été testé sur le pilote de FCC de l'IFP.

MODAL INTERVAL AND FUZZY APPROACH FOR CAUSAL DIAGNOSIS OF REFINERY PROCESSES. APPLICATION TO AN FCC PILOT PROCESS

ABSTRACT

This thesis deals with process diagnosis and is divided into three parts.

The first part presents a causal modeling methodology dedicated to process diagnosis.

The second part presents the fault detection, isolation and identification methodologies. Faults are detected comparing process variables with measurements. We tested a fuzzy and an interval approach to take into account model uncertainties. Then a hitting-sets algorithm is used to generate a list of suspected components. This isolation enables to focus the identification on a reduced set of components. Expert knowledge is finally used to identify faults. It expresses, for each component, action and verification to undertake.

The third part presents the application. The developed diagnosis software ASCO has been tested on an FCC pilot process.

SPECIALITE : AUTOMATIQUE – PRODUCTIQUE

MOTS-CLES : Détection, Localisation et Identification de Défauts, Modélisation Causale, FCC, Raisonnement flou, Intervalles modaux.

LABORATOIRE D'AUTOMATIQUE DE GRENOBLE ENSIEG

Rue de la houille blanche – BP46
38402 SAINT MARTIN D'HERES CEDEX