

Sélection de modèle
pour la classification non supervisée.
Choix du nombre de classes.

Jean-Patrick Baudry

Directeur de thèse : Gilles Celeux

Université Paris-Sud 11
Projet SELECT (INRIA)

3 Décembre 2009

Model Selection for Clustering. How Many Classes?

Jean-Patrick Baudry

Advisor: Gilles Celeux

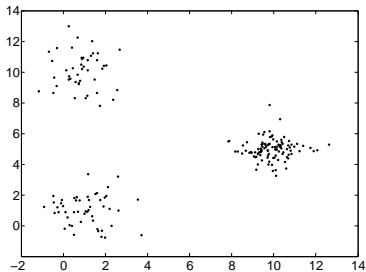
Université Paris-Sud 11
Project `SELECT` (INRIA)

December 3, 2009

Table of contents

- 1 Introduction
 - Clustering
 - Model-Based Clustering
 - Choosing the Number of Classes
- 2 Contrast Minimization for Clustering
 - Conditional Classification Likelihood
 - Estimation: MLccE
 - Model Selection
 - A New Light on ICL
 - Slope Heuristics
- 3 Simulations
- 4 Mixtures of Mixtures
- 5 Conclusion and Perspectives

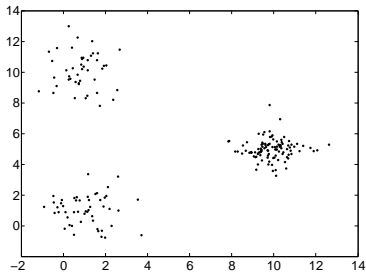
Clustering



Data: $x_1, \dots, x_n \in \mathbb{R}^d$.

Aim: designing K classes.

Clustering

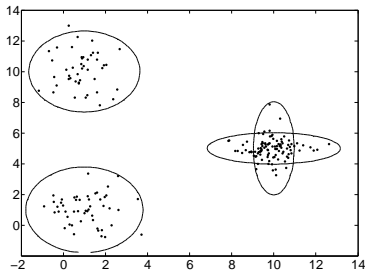


Data: $x_1, \dots, x_n \in \mathbb{R}^d$.

Aim: designing K classes

Two red arrows point from the circled K to two red question marks, one on the left and one on the right.

Clustering



200 observations from a four-component
Gaussian mixture

Data: $x_1, \dots, x_n \in \mathbb{R}^d$.

Aim: designing K classes.

Model-Based Clustering

Statistical Approach: (x_1, \dots, x_n) realization of (X_1, \dots, X_n) i.i.d. $\sim f^\theta$.

- Fit a mixture model to the data.

Model-Based Clustering

- Fit a mixture model to the data.

$$\mathcal{M}_K = \left\{ \sum_{k=1}^K \pi_k \phi(\cdot; \omega_k) \mid (\pi_1, \dots, \pi_K, \omega_1, \dots, \omega_K) \in \Theta_K \right\},$$

$$\text{with } \begin{cases} \Theta_K \subset \Pi_K \times (\mathbb{R}^d \times \mathbb{S}_+^d)^K \\ \Pi_K = \{(\pi_1, \dots, \pi_K) \in [0, 1]^K : \sum_{k=1}^K \pi_k = 1\}. \end{cases}$$

Let us denote:

- ▶ $f(\cdot; \theta) = \sum_{k=1}^K \pi_k \phi(\cdot; \omega_k)$, for all $\theta \in \Theta_K$.
- ▶ $D_K = \dim(\Theta_K)$, “number of free parameters”.

- ▶ Mixture form \leftrightarrow Choice of constraints on Θ_K .
- ▶ One model \leftrightarrow One number of components K .

Model-Based Clustering

- Fit a mixture model to the data.

- Design classes according to the rule

“One Gaussian component = One class”

.

Model-Based Clustering

- Fit a mixture model to the data. Usually:

$$\hat{\theta}_K^{\text{MLE}} \in \underset{\theta \in \Theta_K}{\operatorname{argmax}} \underbrace{\sum_{i=1}^n \log \sum_{k=1}^K \pi_k \phi(X_i; \omega_k)}_{\log L(\theta)}.$$

Study based on a good estimation of the sample distribution.

- Design classes according to the rule

“One Gaussian component = One class”

based on the Maximum A Posteriori:

$$\forall x, \forall k, \forall \theta \in \Theta_K, \quad \tau_k(x; \theta) = \frac{\pi_k \phi(x; \omega_k)}{\sum_{k'=1}^K \pi_{k'} \phi(x; \omega_{k'})}.$$

$$\hat{z}^{\text{MAP}}(\hat{\theta}_K^{\text{MLE}}) = \underset{1 \leq k \leq K}{\operatorname{argmax}} \tau_k(x; \hat{\theta}_K^{\text{MLE}}).$$

Choosing the Number of Classes: Model Selection

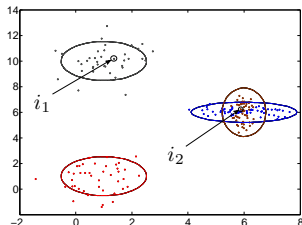
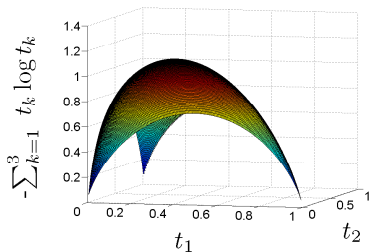
Penalized Likelihood Criteria.

- Efficiency: minimize $d_{KL}(f^{\wp}, f(\cdot; \hat{\theta}_K^{\text{MLE}}))$.
 - ▶ AIC : $\hat{K}^{\text{AIC}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + D_K\}$;
 - ▶ Slope heuristics (Birgé and Massart, 2006).
- Identification: minimize $\min_{\theta \in \Theta_K} d_{KL}(f^{\wp}, f(\cdot; \theta))$.
 - ▶ BIC : $\hat{K}^{\text{BIC}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + \frac{D_K}{2} \log n\}$.
- A criterion adapted to clustering:
 - ▶ ICL (Biernacki, Celeux, Govaert, 2000) :
$$\hat{K}^{\text{ICL}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + \text{ENT}(\hat{\theta}_K^{\text{MLE}}) + \frac{D_K}{2} \log n\}.$$

Entropy: Measure of the Assignment Confidence

$$\text{ENT}(\theta; x) = - \sum_{k=1}^K \tau_k(x; \theta) \log \tau_k(x; \theta) \in [0, \log K].$$

$$\text{ENT}(\theta) = \sum_{i=1}^n \text{ENT}(\theta; x_i).$$



$\text{ENT}(\hat{\theta}_4^{\text{MLE}}; x_{i_1})$ close to 0.
 $\text{ENT}(\hat{\theta}_4^{\text{MLE}}; x_{i_2})$ close to $\log 2$.

Choosing the Number of Classes: Model Selection

Penalized Likelihood Criteria.

- Efficiency: minimize $d_{KL}(f^\varphi, f(\cdot; \hat{\theta}_K^{\text{MLE}}))$.
 - ▶ AIC : $\hat{K}^{\text{AIC}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + D_K\}$;
 - ▶ Slope heuristics (Birgé and Massart, 2006).
- Identification: minimize $\min_{\theta \in \Theta_K} d_{KL}(f^\varphi, f(\cdot; \theta))$.
 - ▶ BIC : $\hat{K}^{\text{BIC}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + \frac{D_K}{2} \log n\}$.
- A criterion adapted to clustering:
 - ▶ ICL (Biernacki, Celeux, Govaert, 2000) :
$$\hat{K}^{\text{ICL}} = \operatorname{argmin}_{1 \leq K \leq K_M} \{-\log L(\hat{\theta}_K^{\text{MLE}}) + \text{ENT}(\hat{\theta}_K^{\text{MLE}}) + \frac{D_K}{2} \log n\}.$$

Conditional Classification Likelihood

The classification log-likelihood for the complete data $(\underline{X}, \underline{Z})$ in model \mathcal{M}_K :

$$\log L_c(\theta; (\underline{X}, \underline{Z})) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \pi_k \phi(X_i; \omega_k).$$

A key relation:

$$\log L_c(\theta) = \log L(\theta) + \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \tau_k(X_i; \theta).$$

Considering the conditional expectation of this relation:

Definition

The conditional classification log-likelihood is

$$\log L_{cc}(\theta) = \log L(\theta) - \text{ENT}(\theta).$$

Contrast Minimization for Clustering

Contrast:

$$-\log L_{\text{cc}}(\theta) = -\log L(\theta) + \text{ENT}(\theta).$$

Associated loss:

$$\begin{aligned} \mathbb{E}_{f^\varphi} [-\log L_{\text{cc}}(\theta)] - \min_{\theta} \mathbb{E}_{f^\varphi} [-\log L_{\text{cc}}(\theta)] \\ \longleftrightarrow d_{\text{KL}}(f^\varphi, f(\cdot; \theta)) + \mathbb{E}_{f^\varphi} [\text{ENT}(\theta)]. \end{aligned}$$

Approximation in the model \mathcal{M}_K :

$$\Theta_K^0 = \operatorname{argmin}_{\theta \in \Theta_K} \left\{ d_{\text{KL}}(f^\varphi, f(\cdot; \theta)) + \mathbb{E}_{f^\varphi} [\text{ENT}(\theta)] \right\}.$$

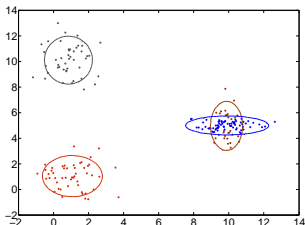
Even if $f^\varphi \in \mathcal{M}_K$, there is no reason that $f^\varphi \in \Theta_K^0$.

Estimation: MLccE

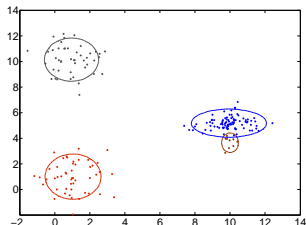
Definition

The minimum empirical contrast estimator, called “Maximum conditional classification Likelihood Estimator”, is defined by

$$\hat{\theta}^{\text{MLccE}} \in \underset{\theta \in \Theta_K}{\operatorname{argmin}} \{ -\log L_{\text{cc}}(\theta) \}.$$



MLE



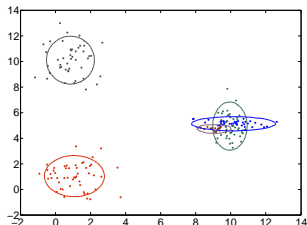
MLccE

Estimation: MLccE

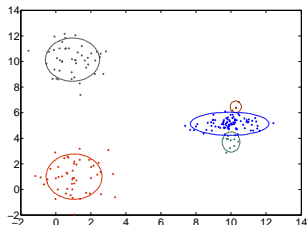
Definition

The minimum empirical contrast estimator, called “Maximum conditional classification Likelihood Estimator”, is defined by

$$\hat{\theta}^{\text{MLccE}} \in \underset{\theta \in \Theta_K}{\operatorname{argmin}} \{ -\log L_{\text{cc}}(\theta) \}.$$



MLE



MLccE

MLccE Properties

Theorem

Let $K \in \mathbb{N}^*$. Assume

- Θ_K compact and convex;
- $H'_K(x) = \sup_{\theta \in \Theta_K} \left\| \left(\frac{\partial \log L_{cc}}{\partial \theta} \right)_{(\theta; x)} \right\|_{\infty} < \infty$ a.s. and $\|H'_K\|_1 < \infty$.
- $\hat{\theta}^{\text{MLccE}} \in \Theta_K$ such that for all $\theta_K^0 \in \Theta_K^0$,

$$-\log L_{cc}(\hat{\theta}_K^{\text{MLccE}}) \leq -\log L_{cc}(\theta_K^0) + o_{\mathbb{P}}(1).$$

Then

$$d(\hat{\theta}_K^{\text{MLccE}}, \Theta_K^0) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0.$$

With $d(\theta, \tilde{\Theta}) = \inf_{\tilde{\theta} \in \tilde{\Theta}} \|\theta - \tilde{\theta}\|_{\infty}$.

Computing MLccE: L_{cc}-EM

The L_{cc}-EM algorithm for MLccE:

$$\hat{\theta}_K^{\text{MLccE}} = \operatorname{argmax}_{\theta \in \Theta_K} \left\{ \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \phi(X_i; \omega_k) - \text{ENT}(\theta) \right\}.$$

Initialization: Random, Small_L_{cc}-EM, CEM, Km1...

Iteration $\theta^j \rightarrow \theta^{j+1}$:

E Step Compute $Q(\theta, \theta^j) = \mathbb{E}_{\theta^j} [\log L_c(\theta; \underline{X}, \underline{Z}) | \underline{X} = \underline{x}]$.

This amounts to computing $\tau_k(x_i; \theta^j)$.

M Step Maximization of $Q(\theta, \theta^j) - \text{ENT}(\theta)$ with respect to $\theta \in \Theta_K$:

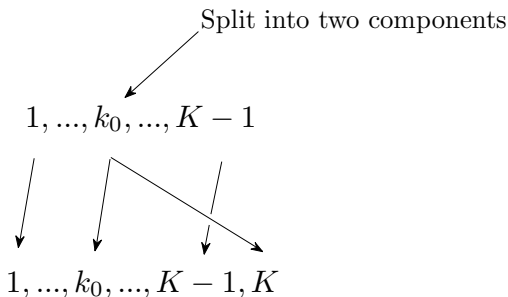
$$\theta^{j+1} \in \operatorname{argmax}_{\theta \in \Theta_K} \left\{ \underbrace{\mathbb{E}_{\theta^j} [\log L_c(\theta; \underline{X}, \underline{Z}) | \underline{X} = \underline{x}]}_{\log L(\theta) + \sum_{i=1}^n \sum_{k=1}^K (\tau_k(x_i; \theta^j) + \tau_k(x_i; \theta)) \log \tau_k(x_i; \theta)} - \text{ENT}(\theta) \right\}$$

Computing MLccE: L_{CC}-EM

The L_{CC}-EM algorithm for MLccE:

$$\hat{\theta}_K^{\text{MLccE}} = \operatorname{argmax}_{\theta \in \Theta_K} \left\{ \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \phi(X_i; \omega_k) - \text{ENT}(\theta) \right\}.$$

Initialization: Random, Small_L_{CC}-EM, CEM, Km1...



Consistent Model Selection

Identification point of view:

$$K_0 = \min_{1 \leq K \leq K_M} \operatorname{argmin} \mathbb{E}_{f^\circ} [-\log L_{\text{cc}}(\Theta_K^0)].$$

Procedures are considered such that

$$\hat{K} = \operatorname{argmin}_{1 \leq K \leq K_M} \left\{ -\log L_{\text{cc}}(\hat{\theta}_K^{\text{MLccE}}) + \text{pen}(K) \right\}.$$

Consistent Model Selection

Theorem

Let us consider the model family $(\mathcal{M}_K)_{K \in \{1, \dots, K_M\}}$. Let us assume:

• $\forall K$, Θ_K is compact and convex.

• $\forall K$, $\forall \theta \in \Theta_K$, $\forall \theta_{K_0}^0 \in \Theta_{K_0}^0$,

$$\mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)] = \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta_{K_0}^0)] \iff -\log L_{cc}(\theta; x) = -\log L_{cc}(\theta_{K_0}^0; x) \text{ a.s.}$$

• $\forall K$, $H_K(x) = \sup_{\theta \in \Theta_K} |\log L_{cc}(\theta; x)| < \infty$ a.s. and $\|H_K\|_\infty < \infty$.

• $\forall K$, $H'_K(x) = \sup_{\theta \in \Theta_K} \left\| \left(\frac{\partial \log L_{cc}}{\partial \theta} \right)_{(\theta; x)} \right\|_\infty < \infty$ a.s. and $\|H'_K\|_2 < \infty$.

• $\forall K$, $\forall \theta_K^0 \in \Theta_K^0$, $\frac{\partial^2}{\partial \theta^2} \left(\mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)] \right)_{|\theta_K^0}$ is nonsingular.

Let $\text{pen} : \{1, \dots, K_M\} \rightarrow \mathbb{R}^+$ such that $\begin{cases} \text{pen}(K) = o_{\mathbb{P}}(n) \text{ as } n \rightarrow \infty \\ (\text{pen}(K) - \text{pen}(K')) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \infty \text{ if } K' < K. \end{cases}$

Then

$$\mathbb{P}[\hat{K} \neq K_0] \xrightarrow[n \rightarrow \infty]{} 0.$$

A New Light on ICL

- Analogy with model selection criteria in the usual observed likelihood framework.
- A good identification criterion, by analogy with BIC:

$$\hat{K}^{\text{L}_{\text{cc}}\text{-ICL}} = \underset{K \in \{1, \dots, K_M\}}{\operatorname{argmin}} \left\{ -\log L_{\text{cc}}(\hat{\theta}_K^{\text{ML}_{\text{cc}}\text{E}}) + \frac{D_K}{2} \log n \right\}.$$

- ICL is an approximation of $L_{\text{cc}}\text{-ICL}$.

A New Light on ICL

- Analogy with model selection criteria in the usual observed likelihood framework.
- A good identification criterion, by analogy with BIC:

$$\hat{K}^{\text{L}_{\text{cc}}\text{-ICL}} = \underset{K \in \{1, \dots, K_M\}}{\operatorname{argmin}} \left\{ -\log L_{\text{cc}}(\hat{\theta}_K^{\text{ML}_{\text{cc}}\text{E}}) + \frac{D_K}{2} \log n \right\}.$$

- ICL is an approximation of $L_{\text{cc}}\text{-ICL}$:

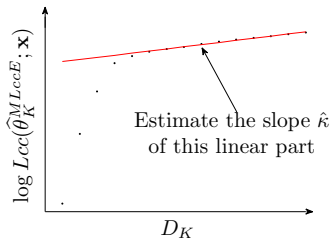
$$\hat{K}^{\text{ICL}} = \underset{K \in \{1, \dots, K_M\}}{\operatorname{argmin}} \left\{ -\log L_{\text{cc}}(\hat{\theta}_K^{\text{MLE}}) + \frac{D_K}{2} \log n \right\}.$$

Slope Heuristics

(Birgé and Massart, 2006)

“Data-driven Slope Estimation”
Assume an “optimal” penalty is known up to a constant κ_{opt} :

$$\text{pen}(K) = \kappa_{\text{opt}} D_K.$$

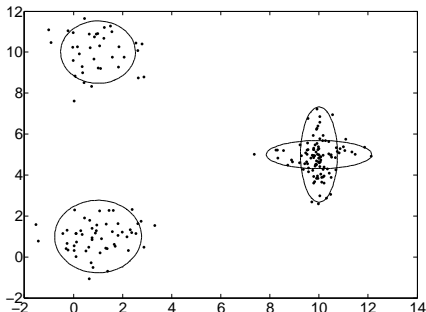


$$\hat{K}^{\text{SHLcc}} = \underset{K \in \{1, \dots, K_M\}}{\text{argmin}} \left\{ -\log L_{\text{cc}}(\hat{\theta}_K^{\text{MLccE}}) + 2\hat{\kappa} D_K \right\}$$

This data-driven procedure may be applied:

- to the usual observed likelihood contrast;
- to the conditional classification likelihood contrast.

“Cross” Dataset



- Simulated data in \mathbb{R}^2 .
- Sample size: 200.
- Number of components: 4.
- Diagonal mixture models fitted: $f^\varnothing \in \mathcal{M}_4$.

“Cross” Dataset: Results

Selected number of components	2	3	4	5	6	7	8	9	10–20
AIC	0	0	1	1	2	2	3	3	88
BIC	0	4	91	5	0	0	0	0	0
SHL	0	2	84	10	3	0	0	0	1
ICL	0	96	3	1	0	0	0	0	0
L_{cc} -ICL	0	99	1	0	0	0	0		
SHL_{cc}	2	79	8	8	3	0	0		

Results for 100 experiments

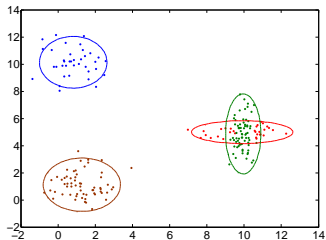
L	Risk $\times 10^3$	$\frac{\text{Risk}}{\text{Oracle Risk}}$
Oracle	59	1
AIC	506	8.03
BIC	65	1.10
(ICL)	156	2.62
SHL	69	1.17

“Oracle” number of components: 4

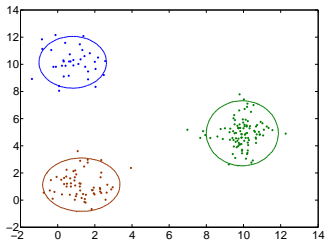
L_{cc}	“Risk” $\times 10^3$
Oracle	3618
ICL	3622
L_{cc} -ICL	3623
SHL_{cc}	3632

“Oracle” number of components : 3

Mixtures of Mixtures

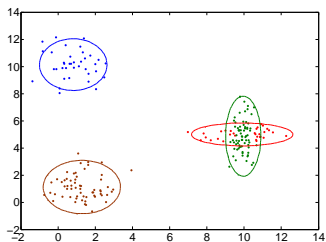


BIC Solution

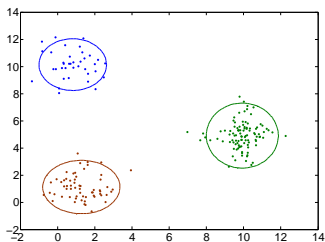


ICL Solution

Mixtures of Mixtures

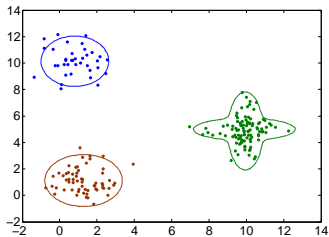


BIC Solution



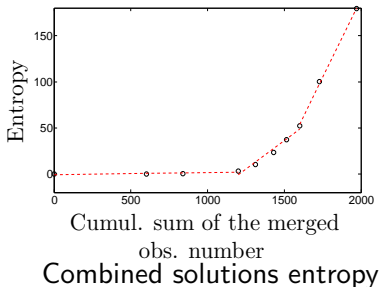
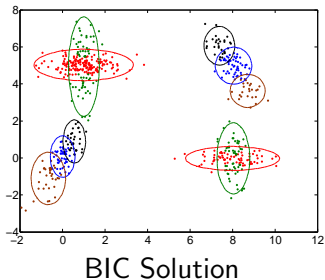
ICL Solution

Combine
two
components



Mixtures of Mixtures

- Hierarchical classes combining:
 - ▶ From the BIC solution with \hat{K}^{BIC} classes;
 - ▶ By minimizing the entropy of the combined solution at each step;
 - ▶ Until there is only one class left.
- Choosing the number of classes:
 - ▶ May be based on substantive ground;
 - ▶ The whole hierarchy may be of interest to the user;
 - ▶ The plot of the entropy against the number of classes may be helpful for the analysis;
 - ▶ Link with the works about penalized criteria?



Conclusions and Perspectives

- The theoretical study of ICL led to the definition of a contrast adapted to the clustering objective, and thus to the corresponding estimator and model selection procedures.
- Solutions are proposed to put these into practice. They may also be applied with benefit when computing the usual MLE through the usual EM algorithm.
- A new light is thrown on ICL, viewed as an approximation of L_{CC} -ICL. This is a contribution to the study of the “class” notion in model-based clustering.
- This “class” notion may be further studied.
- The MLE and $ML_{CC}E$ estimators on the one hand; the ICL and L_{CC} -ICL model selection criteria on the other hand, may be further compared, notably from a practical point of view.
- The slope heuristics, in this mixture models framework, may be further studied.