

# Estimation and model selection for model-based clustering with the conditional classification likelihood

Jean-Patrick Baudry

► **To cite this version:**

Jean-Patrick Baudry. Estimation and model selection for model-based clustering with the conditional classification likelihood. *Electronic journal of statistics*, Shaker Heights, OH: Institute of Mathematical Statistics, 2015, 9 (1), pp.1041-1077. 10.1214/15-EJS1026 . hal-00699578v2

**HAL Id: hal-00699578**

**<https://hal.sorbonne-universite.fr/hal-00699578v2>**

Submitted on 31 May 2012

**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.

# Estimation and Model Selection for Model-Based Clustering with the Conditional Classification Likelihood

Jean-Patrick Baudry\*  
*LSTA*

*Université Pierre et Marie Curie - Paris VI*

May 31, 2012

## Abstract

The Integrated Completed Likelihood (ICL) criterion has been proposed by Biernacki et al. (2000) in the model-based clustering framework to select a relevant number of classes and has been used by statisticians in various application areas.

A theoretical study of this criterion is proposed. A contrast related to the clustering objective is introduced: the *conditional classification likelihood*. This yields an estimator and a model selection criteria class. The properties of these new procedures are studied and ICL is proved to be an approximation of one of these criteria.

We oppose these results to the current leading point of view about ICL, that it would not be consistent. Moreover these results give insights into the class notion underlying ICL and feed a reflection on the class notion in clustering.

General results on penalized minimum contrast criteria and on mixture models are derived, which are interesting in their own right.

**Keywords** BIC; Bracketing entropy; Classification likelihood; Contrast minimization; ICL; Model-based clustering; Model selection; Mixture models; Number of classes; Penalized criteria

---

\*This research was partly supported by Université Paris XI, Laboratoire de Mathématiques d'Orsay UMR 8628; INRIA Saclay Île-de-France, Projet SELECT and Université Paris Descartes, MAP5 UMR 8145.

# 1 Introduction

Model-based clustering is introduced in Sections 1.1 and 1.2. Our purpose is to better understand the behavior of the ICL model selection criterion of Biernacki et al. (2000), which is presented in Section 1.3.

The main topic of this work is the choice of the number of classes in a model-based clustering framework, and then the choice of the number of components of a Gaussian mixture. The interested reader may refer to Titterton et al. (1985) or McLachlan and Peel (2000) for comprehensive studies on Gaussian mixture models. The last also provides an overview on the approaches for assessing the number of components, and particularly on the standard and widely used penalized likelihood criteria, such as AIC (Akaike, 1973) or BIC (Schwarz, 1978).

The ICL criterion studied here is an alternative to BIC. It was up to now widely presented as a penalized likelihood criterion, which penalty involves an “entropy” term. Here, however, we prove that it is actually a penalized contrast criterion with a criterion which is different from the standard likelihood: this justifies why this is not surprising, nor a drawback, that ICL does not asymptotically select the “true” number of components, even when the “true” model is considered. Even for data arising from a mixture distribution, a relevant number of classes may differ from the true number of components of the mixture.

The reason why we introduce this new contrast  $L_{cc}$  (Section 2.1) is not that we believe it *a priori* to be the better one for a clustering purpose, but rather that it enables to theoretically study and understand ICL. We prove (Section 4.3) that ICL is an approximation of a criterion linked to this contrast: studying further ICL then amounts to studying  $L_{cc}$ . The notion of class underlying ICL is proved to be a compromise between Gaussian mixture density estimation and a strictly “cluster” point of view (Section 5).

Let  $X$  be a random variable in  $\mathbb{R}^d$  with distribution  $f^{\varphi \cdot \lambda}$  and  $X_1, \dots, X_n$  an i.i.d. sample of the same distribution. Let us denote  $\mathbf{X} = (X_1, \dots, X_n)$ .

All proofs are gathered in Section 6.

## 1.1 Gaussian Mixture Models

$\mathcal{M}_K$  is the Gaussian mixture model with  $K$  components:

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

where  $\phi$  is the Gaussian density and  $\Theta_K \subset \Pi_K \times (\mathbb{R}^d \times \mathbb{S}_+^d)^K$  with  $\Pi_K = \{(\pi_1, \dots, \pi_K) \in [0, 1]^K : \sum_{k=1}^K \pi_k = 1\}$  and  $\mathbb{S}_+^d$  the set of positive definite  $d \times d$  real matrices. Constraints on the model may be imposed by restricting  $\Theta_K$ . We typically have in mind the decomposition suggested by Celeux and Govaert (1995). “General” (no constraint) and “diagonal” (diagonal covariance matrices) models will be considered here as examples.

Those are studied here as parametric models. It is then assumed the existence of a parametrization  $\varphi : \Theta_K \subset \mathbb{R}^{D_K} \rightarrow \mathcal{M}_K$ . It is assumed that  $\Theta_K$  and  $\varphi$  are “optimal”, in the sense that  $D_K$  is minimal.  $D_K$  is the number of *free parameters* in the model  $\mathcal{M}_K$  and is called the *dimension* of  $\mathcal{M}_K$ . For example, at most  $(K - 1)$  mixing proportions need to be parametrized.

It shall not be needed to assume the parametrization to be identifiable, i.e. that  $\varphi$  is injective. Indeed our purpose is twofold: identifying a relevant number of classes to be designed; and actually designing those classes. Theorem 4.2 justifies that the first task can be achieved under a weaker “identifiability” assumption. Theorem 3.2 then guarantees that our estimator converges to the best parameters set, any of which is as good as the others. There will be no “true parameter” assumption. The classes can finally be defined through the MAP rule (see Section 1.2). Practically, the parameters themselves are never the quantities of interest here. They only stand as a convenient notation and this is also why we expect that the assumption about the Fisher information (see Theorem 4.2) is technical and could maybe be avoided with other techniques. Please refer to Baudry (2009, Chapter 4) for a more comprehensive discussion about the identifiability question.

## 1.2 Model-Based Clustering

Although the results are stated first for much more general situations, this paper is devoted to the question of clustering through Gaussian mixture models.

The process is standard (see Fraley and Raftery, 2002):

- fit each considered mixture model;
- select a model and a number of components based on the first step;
- classify the observations through the MAP rule (recalled below) with respect to the mixture distribution fitted in the selected model.

Notably, the usual choice is made here, to identify a class with each fitted Gaussian component. The number of classes to be designed is then chosen

at the second step. See for example Hennig (2010) and Baudry et al. (2010) for alternative approaches.

Let us recall the MAP classification rule. It involves the *conditional probabilities* of the components

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

$\tau_k(x; \theta)$  is the probability that  $X$  arises from the  $k^{\text{th}}$  component, conditionally to  $X = x$ , under the distribution defined by  $\theta$ . Let us also denote  $\tau_{ik}(\theta) = \tau_k(X_i; \theta)$ . The MAP classification rule for  $x$  is then

$$\hat{z}^{\text{MAP}}(\theta) = \operatorname{argmax}_{k \in \{1, \dots, K\}} \tau_k(x; \theta).$$

Let us denote by  $L$  the *observed likelihood* associated to  $\mathbf{X}$ :

$$\forall \theta \in \Theta_K, L(\theta; \mathbf{X}) = \prod_{i=1}^n \left( \sum_{k=1}^K \pi_k \phi(X_i; \omega_k) \right).$$

The maximum likelihood estimator in the model  $\mathcal{M}_K$  is denoted by  $\hat{\theta}_K^{\text{MLE}}$ .

### 1.3 ICL

Our motivation is to better understand the ICL (Integrated Completed Likelihood) criterion. Let us introduce the *classification likelihood* associated to the complete data sample  $(\mathbf{X}, \mathbf{Z})$  ( $Z \in \{0, 1\}^K$  is the unobserved label of  $X$ :  $Z_k = 1 \Leftrightarrow X$  arises from component  $k$ ):

$$\forall \theta \in \Theta_K, L_c(\theta; (\mathbf{X}, \mathbf{Z})) = \prod_{i=1}^n \prod_{k=1}^K (\pi_k \phi(X_i; \omega_k))^{Z_{ik}}. \quad (1)$$

To mimic the derivation of the BIC criterion (Schwarz, 1978) in a clustering framework, Biernacki et al. (2000) approximate the integrated classification likelihood through a Laplace's approximation. Then they assume that the classification likelihood mode can be identified with  $\hat{\theta}_K^{\text{MLE}}$  as  $n$  is large enough and replace the unobserved  $Z_{ik}$ 's by their MAP estimators under  $\hat{\theta}_K^{\text{MLE}}$ . This is questionable, notably when the components of  $\hat{\theta}_K^{\text{MLE}}$  are not well separated. They derive the ICL criterion:

$$\text{crit}_{\text{ICL}}(K) = \log L(\hat{\theta}_K^{\text{MLE}}) + \sum_{i=1}^n \sum_{k=1}^K \hat{Z}_{ik}^{\text{MAP}}(\hat{\theta}_K^{\text{MLE}}) \log \tau_{ik}(\hat{\theta}_K^{\text{MLE}}) - \frac{\log n}{2} D_K.$$

McLachlan and Peel (2000) replace the  $Z_{ik}$ 's by their conditional expectations  $\tau_{ik}(\hat{\theta}_K^{\text{MLE}})$ :

$$\text{crit}_{\text{ICL}}(K) = \log L(\hat{\theta}_K^{\text{MLE}}) + \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}(\hat{\theta}_K^{\text{MLE}}) \log \tau_{ik}(\hat{\theta}_K^{\text{MLE}}) - \frac{\log n}{2} D_K. \quad (2)$$

Both versions of the ICL appear to behave analogously, and the latter is considered from now on.

The ICL differs from the standard and widely used BIC criterion of Schwarz (1978) through the *entropy* term (see Section 2.2):

$$\forall \theta \in \Theta_K, \text{ENT}(\theta; \mathbf{X}) = - \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}(\theta) \log \tau_{ik}(\theta). \quad (3)$$

The BIC is known to be consistent, in the sense that it asymptotically selects the true number of components, at least when the true distribution actually lies in one of the considered models (Keribin, 2000; Nishii, 1988). This nice property may however not suit a clustering purpose. In many applications, there is no reason to assume that the distribution conditional on the (unobserved) labels  $Z$  is Gaussian. The BIC in this case tends to overestimate the number of components since several Gaussian components are needed to approximate each non-Gaussian component of the true mixture distribution  $f^\varphi$ . And the user may rather be interested in a *cluster* notion — as opposed to this strictly *component* approach — which also includes a separation notion and which be robust to non-Gaussian components. Of course, it depends on the application, and on what a *class* should be. It may be of interest to discriminate into two different classes a group of observations which the best fit is reached with a mixture of two Gaussian components having quite different parameters (we particularly think of the covariance matrices parameters). BIC tends to do so. But it may also be more relevant and may conform to an intuitive notion of cluster, to identify two very close — or largely overlapping — Gaussian components as a single non-Gaussian shaped cluster (see for example Figure 3)...

ICL has been derived with this viewpoint. It is widely understood and explained (for instance in Biernacki et al., 2000) as the BIC criterion with a supplemental penalty, which is the entropy (Section 2.2). Since the last penalizes models which maximum likelihood estimator yields an uncertain MAP classification, ICL is more robust than BIC to non-Gaussian components. However we do not think that the entropy should be considered as a penalty term and an other point of view will be developed in this paper.

The references here were found by browsing the result obtained from Google Scholar citations about Biernacki et al. (2000). Only 3 pages of 16 have been studied... The behavior of ICL has been studied through simulations and real data studies by Biernacki et al. (2000), McLachlan and Peel (2000, Section 6.11), Steele and Raftery (2010) and in several simulation studies (See Baudry, 2009, Chapter 4). Besides several authors chose to use it for the mentioned reasons in various applications area: Goutte et al. (2001) (fMRI images); Pigeau and Gelgon (2005) (image collection automatic sorting); Hamelryck et al. (2006) (protein structure prediction); De Granville et al. (2006) (robots learning); Mariadassou et al. (2010) (uncovering groups of nodes in valued graphs and application to host-parasite interaction networks in forest ecosystems analysis); Rigaiil et al. (2012) (comparative genomic hybridization profile); etc.

This practical interest for ICL lets us think that it meets an interesting notion of cluster, corresponding to what some users expect. But no theoretical study is available. Our main motivation is to go further in this direction. This leads to considering new estimation and model selection procedures for clustering, similar to ICL but for which the development of the underlying logic is driven to its conclusion, from the estimation step to the model selection step, instead of introducing the MLE. It is proved that ICL is an approximation of a criterion which is consistent for a particular loss function.

## 2 A New Contrast: Conditional Classification Likelihood

The contrast minimization framework turns out to be a fruitful approach. It enables to fully understand that ICL is not a penalized likelihood criterion, as opposed to the usual point of view. It should rather be linked to an other contrast: the *conditional classification likelihood*.

### 2.1 Definition, Origin

In a clustering context, the classification likelihood (see (1)) is an interesting quantity but neither the labels  $\mathbf{Z}$  are observed, nor we assume that they even exist (think of the case several models with different number of components are fitted: then at most one can correspond to the true number of classes). Beside the first-mentioned works of Biernacki et al. (2000), Biernacki and Govaert (1997), for example, already proposed to directly involve

the classification likelihood to select the number of classes, by estimating the unobserved data. We propose here to consider its expectation conditional on the observed sample  $\mathbf{X}$ . In case there exists a true classification and a model with the true number of classes is considered, this conditional expectation can be interpreted as the quantity the closest to the classification likelihood, which can be considered given the available information.

Let us report the following algebraic relation between  $L$  and  $L_c$ :

$$\forall \theta \in \Theta_K, \log L_c(\theta) = \log L(\theta) + \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \tau_{ik}(\theta). \quad (4)$$

Then, denoting the conditional expectation of  $\log L_c(\theta)$  by  $\log L_{cc}(\theta)$  (for Conditional Classification log Likelihood),

$$\begin{aligned} \log L_{cc}(\theta) &= \mathbb{E}_\theta [\log L_c(\theta) | \mathbf{X}] \\ &= \log L(\theta) + \underbrace{\sum_{i=1}^n \sum_{k=1}^K \tau_{ik}(\theta) \log \tau_{ik}(\theta)}_{-\text{ENT}(\theta; \mathbf{X})}, \end{aligned}$$

which is obviously linked to the clustering objective. We consider in the following  $-\log L_{cc}$  as an empirical contrast to be minimized.

## 2.2 Entropy

$\log L_{cc}$  differs from  $\log L$  through the *entropy* (see (3)).

The behavior of the entropy is based on the properties of the function  $h : t \in [0, 1] \mapsto (-t \log t)$  (with  $h(0) = 0$ ). This nonnegative function (see Figure 1) takes zero value if and only if  $t = 0$  or  $t = 1$ . It is continuous but not differentiable at 0, and in particular it is not Lipschitz over  $[0, 1]$ , which will be a cause of analysis difficulties. Let us also introduce the function  $h_K : (t_1, \dots, t_K) \in \Pi_K \mapsto \sum_{k=1}^K h(t_k)$ . This nonnegative function (see Figure 2) then takes zero value if and only if there exists  $k_0 \in \{1, \dots, K\}$  such that  $t_{k_0} = 1$  and  $t_k = 0$  for  $k \neq k_0$ . It reaches its maximum value  $\log K$  at  $(t_1, \dots, t_K) = (\frac{1}{K}, \dots, \frac{1}{K})$  (proof in Section 6).

Now, the contribution  $\text{ENT}(\theta; x_i)$  of a single observation to the total entropy  $\text{ENT}(\theta; \mathbf{x})$  is considered. Figure 3 represents a dataset simulated from a four-component Gaussian mixture. Let  $\theta$  be such that  $f(\cdot; \theta) = f^\varphi$ . First,  $\text{ENT}(\theta; x_i) \approx 0$  if and only if there exists  $k_0$  such that  $\tau_{ik_0} \approx 1$  and  $\tau_{ik} \approx 0$  for  $k \neq k_0$ . There is no difficulty to classify  $x_i$  in such a case (for example  $x_{i_1}$ ). Second,  $\text{ENT}(\theta; x_i)$  is all the greater that  $(\tau_{i1}, \dots, \tau_{iK})$  is



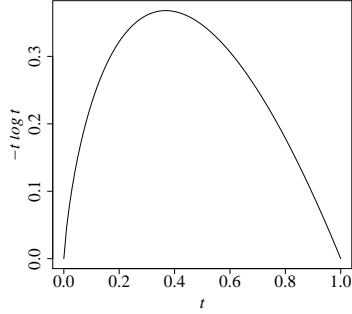


Figure 1:  $h : t \mapsto -t \log t$

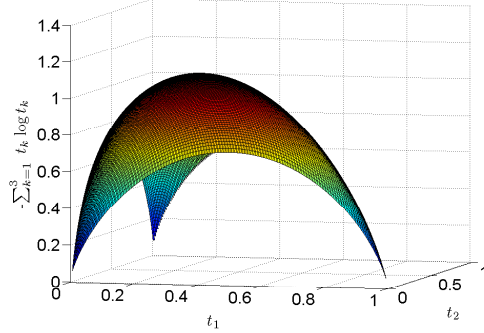


Figure 2:  $h_3 : (t_1, t_2, t_3) \in \Pi_3 \mapsto \sum_{k=1}^3 h(t_k)$

closer to  $(\frac{1}{K}, \dots, \frac{1}{K})$ , i.e. that the classification through the MAP rule is uncertain. The worst case is reached as the conditional distribution over of the components  $1, \dots, K$  is uniform. The observation  $x_{i_2}$  for example has about the same posterior probability  $\frac{1}{2}$  to arise from each one of the components surrounding it. Its individual entropy is about  $\log 2$ .

In conclusion the individual entropy is a measure of the *assignment confidence* of the considered observation through the MAP classification rule. The total entropy  $\text{ENT}(\theta; \mathbf{x})$  is the empirical mean assignment confidence, and then measures the MAP classification quality for the whole sample.

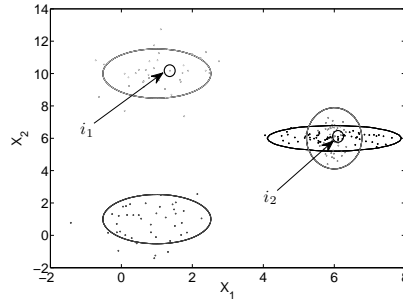


Figure 3: A dataset example

Involving this quantity in a clustering study means that one expects the classification to be confident. The class notion underlying the choice of the conditional classification likelihood as a contrast is then a compromise between the fit (and then the idea of Gaussian-shaped classes) because of the likelihood term on the one hand, and the assignment confidence because

of the entropy term on the other hand (which is rather a *cluster* point of view).

### 2.3 $\log L_{cc}$ as a Contrast

See for example Massart (2007) for an introduction to contrast minimization. Let us consider the best distribution from the  $L_{cc}$  point of view in a model  $\mathcal{M}_m = \{f(\cdot; \theta) : \theta \in \Theta_m\}$ , namely the distribution minimizing the corresponding *loss function*

$$\theta_m \in \underbrace{\operatorname{argmin}_{\theta \in \Theta_m} \{d_{\text{KL}}(f^\varphi, f(\cdot; \theta)) + \mathbb{E}_{f^\varphi} [\text{ENT}(\theta; X)]\}}_{\substack{\operatorname{argmin}_{\theta \in \Theta_m} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)] \\ \text{this set is denoted by } \Theta_m^0}}$$

The existence of  $\mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)]$  is a very mild assumption. The non-emptiness of  $\Theta_m^0$  may be guaranteed for example by assuming  $\Theta_m$  to be compact. Let  $K$  be fixed and consider the minimization of the loss function at hand in the model  $\mathcal{M}_K$  (Section 1.1). First of all, remark that  $\log L_{cc} = \log L$  if  $K = 1$ :  $\Theta_1^0$  is the set of parameters of the distributions which minimize the Kullback-Leibler divergence to  $f^\varphi$ . Now, if  $K > 1$ ,  $\theta_K^0 \in \Theta_K^0$  may be close to minimizing the Kullback-Leibler divergence if the corresponding components do not overlap since then, the entropy is about zero. But if those components overlap, this is not the case anymore (Example 2.1).

To completely define the loss function, and to fully understand this framework, it is necessary to consider the *best element of the universe*  $\mathcal{U}$ :

$$\operatorname{argmin}_{\theta \in \mathcal{U}} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)].$$

The *universe*  $\mathcal{U}$  must be chosen with care. There is no natural relevant choice, on the contrary to the density estimation framework where the set of all densities may be chosen. First the considered contrast is well-defined in a parametric mixture setup, and not necessarily over any mixture densities set because of the definition of the entropy term involving the definition of each component. However, this would still enable to consider mixtures much more general than mixtures of Gaussian components. The ideas developed in Baudry et al. (2010) may for example suggest to involve mixtures which components are Gaussian mixtures. But this would not make sense. The mixture with one component which is a mixture of  $K$  Gaussian components,

and which then yields a single non-Gaussian-shaped class, always has a smaller  $-\log L_{cc}$  value than the corresponding Gaussian mixture yielding  $K$  classes. This illustrates how carefully the components involved in the study must be chosen: involving for example any mixture of Gaussian mixtures means that one considers that a class may be almost anything and may notably contain two Gaussian-shaped clusters very far from each other! The components should in any case be chosen with respect to the corresponding cluster shape. The most natural is then to involve in the universe only Gaussian mixtures:  $\mathcal{U}$  may be chosen as  $\cup_{1 \leq K \leq K_M} \mathcal{M}_K$ .

**Example 2.1.**  $f^\varphi$  is the normal density  $\mathcal{N}(0, 1)$  ( $d = 1$ ). The model  $\mathcal{M}_2 = \{\frac{1}{2}\phi(\cdot; -\mu, \sigma^2) + \frac{1}{2}\phi(\cdot; \mu, \sigma^2); \mu \in \mathbb{R}, \sigma^2 > 0\}$  is considered.

Let us consider  $\Theta_2^0$  in this most simple situation. We numerically obtain that  $\Theta_2^0 = \{(-\mu_0, \sigma_0^2), (\mu_0, \sigma_0^2)\}$ , so that, up to a label switch, there exists a unique minimizer of  $\mathbb{E}_{f^\varphi} [-\log L_{cc}(\mu, \sigma^2)]$  in  $\Theta_2$  in this case (see Figure 4), with  $\mu_0 \approx 0.83$  and  $\sigma_0^2 \approx 0.31$ . This solution is obviously not the same as the one minimizing the Kullback-Leibler divergence (see Figure 5). This illustrates that the objective with the  $-\log L_{cc}$  contrast is not to recover the true distribution, even when it is available in the considered model.

The necessity of choosing a relevant model is striking in this example: this two-component model should obviously not be used for a clustering purpose, at least for datasets with great enough size.

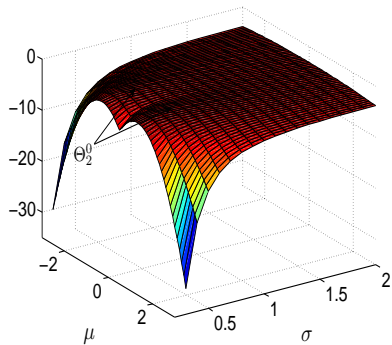


Figure 4:  $\mathbb{E}_{f^\varphi} [-\log L_{cc}(\mu, \sigma^2)]$  w.r.t.  $\mu$  and  $\sigma$ , and  $\Theta_2^0$ , for Example 2.1

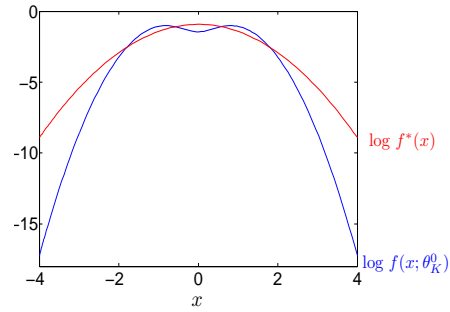


Figure 5:  $\log f^\varphi$  (red, which is also  $\log f(\cdot; \theta_2^{KL})$ ) and  $\log f(\cdot; \theta_2^0)$  (blue) for Example 2.1

The estimator associated to the  $-\log L_{cc}$  contrast is now considered.

### 3 Estimation: MLccE

Let us fix the number of components  $K$  and the model  $\mathcal{M}_K$ . The subscript  $K$  is omitted in the notation of this section. A new minimum contrast estimator is considered. Results are stated in a general parametric model setting with a general contrast  $\gamma$  and a model  $\mathcal{M}$  with parameter space  $\Theta \subset \mathbb{R}^D$ , and then the conditions they involve are discussed in our framework. General conditions ensuring the consistency of such an estimator are given in Theorem 3.1. They notably involve the Glivenko-Cantelli property of the class of functions  $\{\gamma(\theta) : \theta \in \Theta\}$ . Sufficient conditions in terms of bracketing entropy for this property to hold are recalled and verified in the considered context in Section 3.2. Those results are also useful in the study of the model selection step (Section 4). Brought together, they provide the consistency of the estimator in Gaussian mixture models: this is Theorem 3.2.

Here and hereafter, all expectations  $\mathbb{E}$  and probabilities  $\mathbb{P}$  are taken with respect to  $f^\varphi \cdot \lambda$ . For a general contrast  $\gamma$ , we write its empirical version:  $\gamma_n(\theta) = \frac{1}{n} \sum_{i=1}^n \gamma(\theta; X_i)$ .  $\mathbb{R}^D$  is equipped with the infinite norm:  $\forall \theta \in \mathbb{R}^D, \|\theta\|_\infty = \max_{1 \leq i \leq D} |\theta_i|$ . For any  $r \in \mathbb{N}^* \cup \{\infty\}$  and for any  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\|g\|_r = \mathbb{E}_{f^\varphi} [|g(X)|^r]^{\frac{1}{r}}$  if  $r < \infty$  and  $\|g\|_\infty = \text{ess sup}_{X \sim f^\varphi} |g(X)|$  (recall that  $\text{ess sup}_{Z \sim \mathbb{P}} Z = \inf\{z : \mathbb{P}[Z \leq z] = 1\}$  and thus:  $\|g\|_\infty \leq \sup_{x \in \text{supp } f^\varphi} |g(x)|$ ). For any linear form  $l : \mathbb{R}^D \rightarrow \mathbb{R}$ ,  $\|l\|_\infty = \max_{\|\theta\|_\infty=1} l(\theta)$ .

#### 3.1 Definition, Consistency

The minimum contrast estimator is named MLccE (Maximum conditional classification Likelihood Estimator):

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

To ensure its existence, we assume that  $\Theta$  is compact. This is a heavy assumption, but it will be natural and necessary for the following results to hold. That the covariance matrices are bounded from below is a reasonable and necessary assumption in the Gaussian mixture framework: without this assumption, neither the log likelihood, nor the conditional classification likelihood would be bounded (for  $K \geq 2$ ). Insights to choose lower bounds on the proportions and the covariance matrices are suggested in Baudry (2009, Section 5.1). The upper bound on the covariance matrices and the compactness condition on the means, although not necessary in the standard likelihood framework, do not seem to be avoidable here (see Section 3.2).

This is a consequence of the behavior of the entropy term as a component goes to zero.

The following theorem, which is directly adapted from van der Vaart (1998, Section 5.2), gives sufficient conditions for the consistency of a minimum contrast estimator  $\hat{\theta}$ . We write  $\forall \theta \in \Theta, \forall \tilde{\Theta} \subset \Theta, d(\theta, \tilde{\Theta}) = \inf_{\tilde{\theta} \in \tilde{\Theta}} \|\theta - \tilde{\theta}\|_\infty$ .

**Theorem 3.1.** *Let  $\Theta \subset \mathbb{R}^D$  and  $\gamma : \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}$ . Assume:*

$$\exists \theta^0 \in \Theta \text{ such that } \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] = \min_{\theta \in \Theta} \mathbb{E}_{f^\varphi} [\gamma(\theta)] \text{ (}\Leftrightarrow \Theta^0 \text{ is not empty)} \quad (\text{A1})$$

$$\forall \varepsilon > 0, \inf_{\{\theta : d(\theta, \Theta^0) > \varepsilon\}} \mathbb{E}_{f^\varphi} [\gamma(\theta)] > \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] \quad (\text{A2})$$

$$\sup_{\theta \in \Theta} \left| \gamma_n(\theta) - \mathbb{E}_{f^\varphi} [\gamma(\theta)] \right| \xrightarrow{\mathbb{P}} 0 \quad (\text{A3})$$

Define  $\forall n, \hat{\theta} = \hat{\theta}(X_1, \dots, X_n) \in \Theta$  such that  $\gamma_n(\hat{\theta}) \leq \gamma_n(\theta^0) + o_{\mathbb{P}}(1)$ .

Then  $d(\hat{\theta}, \Theta^0) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$ .

The strong consistency holds if (A3) is replaced by an almost sure convergence (this is the case under the conditions we are to define) and if the inequality in the definition of  $\hat{\theta}$  holds almost surely.

Assumption (A1) is the least that can be expected. It is guaranteed if the parameter space is compact.

Assumption (A2) holds, too, under this compactness assumption: since  $\theta \in \Theta_K \mapsto \mathbb{E}_{f^\varphi} [\gamma(\theta)]$  reaches its minimum value on the compact  $\Theta_K \setminus \{\theta \in \Theta_K : d(\theta, \Theta_K^0) > \varepsilon\}$ , it is necessarily strictly greater than  $\mathbb{E}_{f^\varphi} [\gamma(\theta^0)]$ .

Assumption (A3) is a bit strong but it will be guaranteed under the compactness assumption through bracketing entropy arguments in Section 3.2.

*Sketch of proof.* The assumptions guarantee a convenient situation. With great probability as  $n$  grows, from (A3),  $\gamma_n(\theta)$  is uniformly close to  $\mathbb{E}_{f^\varphi} [\gamma(\theta)]$ : this holds for  $\hat{\theta}$  and  $\theta^0$ . Then, from the definition of  $\hat{\theta}$ ,  $\mathbb{E}_{f^\varphi} [\gamma(\hat{\theta})]$  cannot be much larger than  $\mathbb{E}_{f^\varphi} [\gamma(\theta^0)]$  which reaches the minimal value. By (A2), this implies that  $\hat{\theta}$  cannot be far from  $\Theta^0$ . □

Let us apply Theorem 3.1 to Gaussian mixtures, with  $\gamma = -\log L_{\text{cc}}$  and  $\Theta = \Theta_K$ . The two following hypotheses will be involved:

$$\|M\|_r < \infty \text{ with } M(x) = \sup_{\theta \in \Theta} |\gamma(\theta; x)| < \infty \text{ } f^\varphi d\lambda\text{-a.e.} \quad (H_{\gamma, \Theta, r}^M)$$

$$\|M'\|_r < \infty \text{ with } M'(x) = \sup_{\theta \in \Theta} \left\| \left( \frac{\partial \gamma}{\partial \theta} \right)_{(\theta; x)} \right\|_\infty < \infty \text{ } f^\varphi d\lambda\text{-a.e.} \quad (H_{\gamma, \Theta, r}^{M'})$$

**Theorem 3.2** (Weak Consistency of MLccE, compact case). *Let  $\mathcal{M}$  be a Gaussian mixture model with compact parameter space  $\Theta \subset \mathbb{R}^D$ . Let  $\Theta^0 = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)]$ . Let  $\Theta^\circ \subset \mathbb{R}^D$  open over which  $\log L_{cc}$  is defined, such that  $\Theta \subset \Theta^\circ$  and assume that  $H_{\log L_{cc}, \Theta^\circ, 1}^{M'}$  holds. Let  $\hat{\theta}^{MLccE} \in \Theta$  be an estimator (almost) maximizing  $\log L_{cc}$ :*

$$\forall \theta^0 \in \Theta^0, \forall n \in \mathbb{N}^*, -\log L_{cc}(\hat{\theta}^{MLccE}) \leq -\log L_{cc}(\theta^0) + o_{\mathbb{P}}(n).$$

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

$H_{\log L_{cc}, \Theta^\circ, 1}^{M'}$  results from lemma 3.2 and shall be discussed in Section 3.2.

Under the compactness assumption,  $\hat{\theta}^{MLccE}$  is then consistent. It is even strongly consistent if it minimizes the empirical contrast almost surely. Let us highlight that it then converges to the set of parameters minimizing the loss function, which has no reason to contain the true distribution — except for  $K = 1$  — even if the last lies in  $\mathcal{M}$ .

## 3.2 Bracketing Entropy and Glivenko-Cantelli Property

Recall a class of functions over  $\mathbb{R}^d$  is  $\mathbb{P}$ -*Glivenko-Cantelli*, with  $\mathbb{P}$  a probability measure over  $\mathbb{R}^d$ , if it fulfills a uniform law of large numbers for the distribution  $\mathbb{P}$ . A sufficient condition for a family  $\mathcal{G}$  to be  $\mathbb{P}$ -Glivenko-Cantelli is that it is not too complex, which can be measured through its *entropy with bracketing*:

**Definition 3.1** ( $L_r(\mathbb{P})$ -entropy with bracketing). *Let  $r \in \mathbb{N}^*$  and  $l, u \in L_r(\mathbb{P})$ . The bracket  $[l, u]$  is the set of all functions  $g \in \mathcal{G}$  with  $l \leq g \leq u$ .  $[l, u]$  is an  $\varepsilon$ -bracket if  $\|l - u\|_r \leq \varepsilon$ . The bracketing number  $N_{[]}(\varepsilon, \mathcal{G}, L_r(\mathbb{P}))$  is the minimum number of  $\varepsilon$ -brackets needed to cover  $\mathcal{G}$ . The entropy with bracketing  $\mathcal{E}_{[]}(\varepsilon, \mathcal{G}, L_r(\mathbb{P}))$  of  $\mathcal{G}$  with respect to  $\mathbb{P}$  is the logarithm of  $N_{[]}(\varepsilon, \mathcal{G}, L_r(\mathbb{P}))$ .*

It is quite natural that the behavior of all functions lying inside an  $\varepsilon$ -bracket can be uniformly controlled by the behavior of the extrema of the bracket. If those endpoints belong to  $L_1(\mathbb{P})$ , they fulfill a law of large numbers, and if the number of them needed to cover  $\mathcal{G}$  is finite, then this is no surprise that  $\mathcal{G}$  can be proved to fulfill a uniform law of large numbers:

**Theorem 3.3.** *Every class  $\mathcal{G}$  of measurable functions such that  $\mathcal{E}_{[]}(\varepsilon, \mathcal{G}, L_1(\mathbb{P})) < \infty$  for every  $\varepsilon > 0$  is  $\mathbb{P}$ -Glivenko-Cantelli.*

The reader is referred to van der Vaart (1998, Chapter 19) for accurate definitions and a proof of this result. This is a generalization of the

usual Glivenko-Cantelli theorem. We shall prove that the class of functions  $\{\gamma(\cdot; \theta) : \theta \in \Theta_K\}$  has finite  $\varepsilon$ -bracketing entropy for any  $\varepsilon > 0$  and the assumption (A3) will be ensured.

From now on, since  $\Theta$  is typically assumed to be compact, it is assumed that  $\Theta \subset \Theta^\mathcal{O} \subset \mathbb{R}^D$  with  $\Theta^\mathcal{O}$  open over which  $\gamma$  is defined and  $C^1$  for  $f^\varphi d\lambda$ -almost all  $x$ . This is no problem for Gaussian mixture models with  $\log L_{cc}$  (or the standard likelihood by the way), for example with the general or diagonal model. But this requires (with the  $\log L_{cc}$  contrast) the proportions to be positive. Actually, this could be avoided here, but we will need this assumption for the definition of  $M'$  (Hypothesis  $H_{\gamma, \Theta, r}^{M'}$ ). As already mentioned, components going to zero must be avoided. For the same technical reason, we have to assume the mean parameters to be bounded.

Lemma 3.1 guarantees that the bracketing entropy of  $\{\gamma(\cdot; \theta) : \theta \in \Theta\}$  is finite for any  $\varepsilon$ , if  $\Theta$  is convex and bounded. The assumption about the differential of the contrast is not a difficulty in our framework, provided that non-zero lower bounds over  $\Theta$  on the proportions and the covariance matrices are imposed. The lemma is written for any  $\tilde{\Theta}$  bounded and included in  $\Theta$  (which is not assumed to be bounded itself) since it will be applied locally around  $\theta^0$  in the Section 4.

For any bounded  $\tilde{\Theta} \subset \mathbb{R}^D$ ,  $\text{diam } \tilde{\Theta} = \sup\{\|\theta_1 - \theta_2\|_\infty : \theta_1, \theta_2 \in \tilde{\Theta}\}$ .

**Lemma 3.1** (Bracketing Entropy, Convex Case). *Let  $r \in \mathbb{N}^*$ ,  $D \in \mathbb{N}^*$  and  $\Theta \subset \mathbb{R}^D$  assumed to be convex. Let  $\Theta^\mathcal{O} \subset \mathbb{R}^D$  open such that  $\Theta \subset \Theta^\mathcal{O}$  and  $\gamma : \Theta^\mathcal{O} \times \mathbb{R}^d \rightarrow \mathbb{R}$ .  $\theta \in \Theta^\mathcal{O} \mapsto \gamma(\theta; x)$  is assumed to be  $C^1$  for  $f^\varphi d\lambda$ -almost all  $x$ . Assume that  $H_{\gamma, \Theta, r}^{M'}$  holds. Then*

$$\forall \tilde{\Theta} \subset \Theta, \forall \varepsilon > 0, N_{[\cdot]}(\varepsilon, \{\gamma(\theta) : \theta \in \tilde{\Theta}\}, \|\cdot\|_r) \leq \left( \frac{\|M'\|_r \text{diam } \tilde{\Theta}}{\varepsilon} \right)^D \vee 1.$$

Remark that  $\Theta$  does not have to be compact. Its proof is a calculation which relies on the mean value theorem, hence the convexity assumption. The natural parameter space of diagonal Gaussian mixture models, with equal volumes (if  $d > 1$ ) or not, for instance, is convex (see Examples 6.1 and 6.2, p. 26). General mixture models have a convex natural parameter space, too, since the set of definite positive matrices is convex. However, there is no reason that the parameter space  $\Theta$  should be convex in general.

Lemma 3.1 can then be generalized at the price of assuming  $\Theta$  to be compact, and included in an open set  $\Theta^\mathcal{O}$  such that  $H_{\gamma, \Theta^\mathcal{O}, r}^{M'}$  holds. This is no difficulty for the mixture models we consider, under the same lower bounds constraints as before (since  $\Theta^\mathcal{O}$  itself can be chosen to be included

in a compact subset of the set of possible parameters). The entropy is then increased by a multiplying factor  $Q$ , which only depends on  $\Theta$  and roughly measures its “nonconvexity”. Since the exponential behavior of the entropy with respect to  $\varepsilon$  is of concern, this does not make the result really weaker.

**Lemma 3.2** (Bracketing Entropy, Compact Case). *Let  $r \in \mathbb{N}^*$ ,  $D \in \mathbb{N}^*$  and  $\Theta \subset \mathbb{R}^D$  assumed to be compact. Let  $\Theta^{\mathcal{O}} \subset \mathbb{R}^D$  open such that  $\Theta \subset \Theta^{\mathcal{O}}$  and  $\gamma : \Theta^{\mathcal{O}} \times \mathbb{R}^d \rightarrow \mathbb{R}$ .  $\theta \in \Theta^{\mathcal{O}} \mapsto \gamma(\theta; x)$  is assumed to be  $C^1$  for  $f^{\varphi} d\lambda$ -almost all  $x$ . Assume that  $H_{\gamma, \Theta^{\mathcal{O}}, r}^{M'}$  holds.*

*Then*

$$\exists Q \in \mathbb{N}^*, \forall \tilde{\Theta} \subset \Theta, \forall \varepsilon > 0,$$

$$N_{[\cdot]}(\varepsilon, \{\gamma(\theta) : \theta \in \tilde{\Theta}\}, \|\cdot\|_r) \leq Q \left( \frac{\|M'\|_r \text{diam} \tilde{\Theta}}{\varepsilon} \right)^D \vee 1.$$

$Q$  is a constant which depends on the geometry of  $\Theta$  ( $Q = 1$  if  $\Theta$  is convex).

This lemma is proved by applying Lemma 3.1 since  $\Theta$  is still locally convex. Since it is compact, it can be covered with a finite number  $Q$  of open balls, which are convex. Lemma 3.1 then applies to the convex hull of the intersection of  $\Theta$  with each one of them. The supremum of  $M'$  is taken over  $\Theta^{\mathcal{O}}$  — instead of  $\Theta$  — to make sure that the assumptions of Lemma 3.1 are fulfilled over those entire balls, which may not be included in  $\Theta$ .

The result we need for Section 4 is Lemma 3.3, obtained from Lemma 3.1 by a slight modification. Since it is applied locally there, the convexity assumption is no problem. A supplementary and strong assumption  $H_{\gamma, \Theta, \infty}^M$  is made. This is not fulfilled in the general Gaussian mixtures framework. A sufficient condition is that the support of  $f^{\varphi}$  is bounded. This is false of course for most usual distributions we may have in mind, but this is a reasonable modeling assumption: most modeled phenomena are bounded. Another sufficient condition to guarantee this assumption is that the contrast is bounded from above. This is actually not the case of the contrast  $-\log L_{cc}$ , but this can be imposed: replace  $-\log L_{cc}$  by  $(-\log L_{cc} \wedge C)$  and, provided that  $C$  is large enough, this new contrast behaves like  $\log L_{cc}$ . This is a supplemental difficulty in practice to choose a relevant  $C$  value, though.

**Lemma 3.3** (Bracketing Entropy, Convex Case). *Let  $r \geq 2$ ,  $D \in \mathbb{N}^*$  and  $\Theta \subset \mathbb{R}^D$  assumed to be convex. Let  $\Theta^{\mathcal{O}} \subset \mathbb{R}^D$  open such that  $\Theta \subset \Theta^{\mathcal{O}}$  and  $\gamma : \mathbb{R}^D \times \Theta^{\mathcal{O}} \rightarrow \mathbb{R}$ .  $\theta \in \Theta^{\mathcal{O}} \mapsto \gamma(\theta; x)$  is assumed to be  $C^1$  for  $f^{\varphi} d\lambda$ -almost*



all  $x$ . Assume that  $H_{\gamma, \Theta, \infty}^M$  and  $H_{\gamma, \Theta, 2}^{M'}$  hold. Then

$$\forall \tilde{\Theta} \subset \Theta, \forall \varepsilon > 0,$$

$$N_{[\cdot]}(\varepsilon, \{\gamma(\theta) : \theta \in \tilde{\Theta}\}, \|\cdot\|_r) \leq \left( \frac{2^{r-2} \|M\|_{\infty}^{\frac{r-2}{2}} \|M'\|_2 \text{diam} \tilde{\Theta}}{\varepsilon^{\frac{r}{2}}} \right)^D \vee 1.$$

Let us remark that those results are quite general. We are interested here in their application to the conditional classification likelihood, but they hold all the same in the standard likelihood framework. Maugis and Michel (2011) already provide bracketing entropy results in this framework. Our results cannot be directly compared to theirs since they consider the Hellinger distance. The dependency they get on the parameter space bounds and the variable space dimension  $d$  is explicit. This is helpful to derive an oracle inequality. But they could not derive a local control of the entropy, hence an unpleasant logarithm term in the expression of the optimal penalty they get. Their results also suggest the necessity of assuming the contrast to be bounded: see the discussion after the Theorem 4.2. The results we propose achieve the same rate with respect to  $\varepsilon$ . They depend on more opaque quantities ( $\|M\|_{\infty}$  and  $\|M'\|_2$ ). This notably implies, from this first step already, the assumption that the contrast is bounded — over the true distribution support. However, it could be expected to control those quantities with respect to the parameter space bounds. Moreover, beside their simplicity, they straightforwardly enable to derive a local control of the entropy.

## 4 Model Selection

As illustrated by Example 2.1, model selection is a crucial step. The number of classes may even be the target of the study. Anyhow, a relevant number of classes must obviously be chosen so as to design a good classification.

Model selection procedures introduced here are penalized conditional classification likelihood criteria:

$$\text{crit}(K) = -\log L_{\text{cc}}(\hat{\theta}_K^{\text{MLccE}}) + \text{pen}(K).$$

Most results are stated for a general contrast  $\gamma$  and any family of models  $\{\mathcal{M}_K\}_{1 \leq K \leq K_M}$  and then applied to  $-\log L_{\text{cc}}$  and the Gaussian mixtures family of models  $\{\mathcal{M}_K\}_{1 \leq K \leq K_M}$  introduced in Section 1.1.

In Section 4.1, the consistency of such a model selection procedure (“identification” point of view) is proved for a class of penalties. Sufficient

conditions are given in the general Theorem 4.1, which is applied to the framework we are interested in in Theorem 4.2. The heaviest condition of Theorem 4.1 (B4) may be guaranteed under regularity and (weak) identifiability assumptions, and is discussed in Section 4.2. Our approach is inspired from works of Massart (2007) and is the first step to reach non-asymptotic results.

#### 4.1 Consistent Penalized Criteria

Assume that  $K_0$  exists such that

$$\begin{aligned} & \forall K < K_0, \inf_{\theta \in \Theta_{K_0}} \mathbb{E}_{f^\varphi} [\gamma(\theta)] < \inf_{\theta \in \Theta_K} \mathbb{E}_{f^\varphi} [\gamma(\theta)] \\ \text{and} & \forall K \geq K_0, \inf_{\theta \in \Theta_{K_0}} \mathbb{E}_{f^\varphi} [\gamma(\theta)] \leq \inf_{\theta \in \Theta_K} \mathbb{E}_{f^\varphi} [\gamma(\theta)] \end{aligned}$$

which means that the bias of the models is stationary from the model  $\mathcal{M}_{K_0}$ : it is the “best” model. Remark that the last property should hold mostly in the mixtures framework, and notably if the models were not constrained, and then were nested. Under this assumption, a model selection procedure is expected to asymptotically recover  $K_0$ , i.e. to be *consistent*. This is an *identification* aim (see McQuarrie and Tsai, 1998, Chapter 1). It would be disastrous to select a model which does not (almost) minimize the bias. And it is besides assumed that the model  $\mathcal{M}_{K_0}$  contains all the interesting information (typically, the structure of the classes).

Let us stress that the “true” number of components of  $f^\varphi$  is not directly of concern: it is in particular not assumed that it equals  $K_0$ , and is not even assumed to be defined ( $f^\varphi$  does not have to be a Gaussian mixture).  $K_0$  is the best choice from the particular point of view introduced by using the  $\log L_{cc}$  contrast, which is not density estimation, neither is it identification of the “true” number of components.

**Theorem 4.1.**  $\{\Theta_K\}_{1 \leq K \leq K_M}$  a collection of models with  $\Theta_K \subset \mathbb{R}^{D_K}$  ( $D_1 \leq \dots \leq D_{K_M}$ ) and let  $\theta_K^0 \in \Theta_K^0$ , with  $\Theta_K^0 = \operatorname{argmin}_{\theta \in \Theta_K} \mathbb{E}_{f^\varphi} [\gamma(\theta)]$ . Assume

$$K_0 = \min_{1 \leq K \leq K_M} \operatorname{argmin}_{\theta \in \Theta_K} \mathbb{E}_{f^\varphi} [\gamma(\theta)] \quad (\text{B1})$$

$$\begin{aligned} \forall K, \hat{\theta}_K \in \Theta_K \text{ defined such that } \gamma_n(\hat{\theta}_K) &\leq \gamma_n(\theta_K^0) + o_{\mathbb{P}}(1) \\ \text{fulfills } \gamma_n(\hat{\theta}_K) &\xrightarrow{\mathbb{P}} \mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] \end{aligned} \quad (\text{B2})$$

$$\forall K, \begin{cases} \text{pen}(K) > 0 \text{ and } \text{pen}(K) = o_{\mathbb{P}}(1) & \text{when } n \rightarrow +\infty \\ n(\text{pen}(K) - \text{pen}(K')) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \infty & \text{when } K > K' \end{cases} \quad (\text{B3})$$

$$n(\gamma_n(\hat{\theta}_{K_0}) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1) \text{ for any } K \in \underset{1 \leq K \leq K_M}{\text{argmin}} \mathbb{E}_{f^\varphi} [\gamma(\Theta_K^0)]. \quad (\text{B4})$$

Define  $\hat{K}$  such that  $\hat{K} = \min_{1 \leq K \leq K_M} \underset{\text{crit}(K)}{\text{argmin}} \left\{ \underbrace{\gamma_n(\hat{\theta}_K) + \text{pen}(K)}_{\text{crit}(K)} \right\}$ .

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

*Sketch of proof.* First prove that  $\hat{K}$  cannot asymptotically “underestimate”  $K_0$ . Suppose  $\mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] > \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)]$ . From (B2),  $(\gamma_n(\hat{\theta}_K) - \gamma_n(\hat{\theta}_{K_0}))$  is asymptotically of order  $\mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] - \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)] > 0$ . Since the penalty is  $o_{\mathbb{P}}(1)$  from (B3),  $\text{crit}(K_0) < \text{crit}(K)$  asymptotically and  $\hat{K} > K$ .

That  $\hat{K}$  does not asymptotically “overestimate”  $K_0$ , involves the heaviest assumption (B4). It is more involved since then  $(\mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] - \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)])$  is zero. The fluctuations of  $(\gamma_n(\hat{\theta}_K) - \gamma_n(\hat{\theta}_{K_0}))$  around zero then have to be evaluated and canceled by a penalty large enough. According to (B4), a penalty larger than  $\frac{1}{n}$  should suffice. (B3) guarantees this condition.  $\square$

Assumption (B3) defines the range of possible penalties; Assumption (B2) is guaranteed under assumption (A3) of Theorem 3.1:

**Lemma 4.1.** *For a fixed  $K$ , assume (A3). Then (B2) holds.*

Indeed, asymptotically, minimizing  $\theta \mapsto \gamma_n(\theta)$  cannot differ much from minimizing  $\theta \mapsto \mathbb{E}_{f^\varphi} [\gamma(\theta)]$  if they are uniformly close to each other (A3).

Assumption (B4) is the heaviest assumption. Section 4.2 is devoted to deriving sufficient conditions so that it holds. This will justify the

**Theorem 4.2.**  $(\mathcal{M}_K)_{1 \leq K \leq K_M}$  Gaussian mixture models with compact parameter space  $\Theta_K$  and  $\Theta_K^0 = \underset{\theta \in \Theta_K}{\text{argmin}} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)]$  for any  $K$ . Let  $K_0 = \min_{1 \leq K \leq K_M} \underset{\theta \in \Theta_K^0}{\text{argmin}} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)]$ . Assume  $\forall K, \forall \theta \in \Theta_K, \forall \theta_{K_0}^0 \in \Theta_{K_0}^0$ ,

$$\begin{aligned} \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)] &= \mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta_{K_0}^0)] \\ \iff -\log L_{cc}(\theta; x) &= -\log L_{cc}(\theta_{K_0}^0; x) \quad f^\varphi d\lambda - a.e. \quad (\text{C1}) \end{aligned}$$

For any  $K$ , let  $\Theta_K^{\mathcal{O}} \subset \mathbb{R}^{D_K}$  open over which  $\log L_{cc}$  is defined, such that  $\Theta_K \subset \Theta_K^{\mathcal{O}}$ ; assume that  $H_{\log L_{cc}, \Theta_K^{\mathcal{O}}, \infty}^M$  and  $H_{\log L_{cc}, \Theta_K^{\mathcal{O}}, 2}^{M'}$  hold and that  $\forall \theta_K^0 \in \Theta_K^{\mathcal{O}}$ ,  $I_{\theta_K^0} = \frac{\partial^2}{\partial \theta^2} (\mathbb{E}_{f^\varphi} [-\log L_{cc}(\theta)])|_{\theta_K^0}$  is nonsingular; let  $\hat{\theta}_K^{MLccE} \in \Theta_K$  with

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

Let  $\text{pen} : \{1, \dots, K_M\} \rightarrow \mathbb{R}^+$  (which may depend on  $n$ ,  $(\Theta_K)_{1 \leq K \leq K_M}$  and the data) such that

$$\forall K \in \{1, \dots, K_M\}, \begin{cases} \text{pen}(K) > 0 \text{ and } \text{pen}(K) = o_{\mathbb{P}}(n) & \text{when } n \rightarrow +\infty \\ (\text{pen}(K) - \text{pen}(K')) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \infty & \text{for any } K' < K. \end{cases}$$

Select  $\hat{K}$  such that  $\hat{K} = \min_{1 \leq K \leq K_M} \text{argmin} \{-\log L_{cc}(\hat{\theta}_K^{MLccE}) + \text{pen}(K)\}$ .

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

If  $\Theta_K$  is convex,  $M$  and  $M'$  can be defined as suprema over  $\Theta_K$  instead of  $\Theta_K^{\mathcal{O}}$  and there is no need to introduce the sets  $\Theta_K^{\mathcal{O}}$ . The new ‘‘identifiability’’ assumption (C1) introduced is reasonable: as expected the label switching phenomenon is no problem here. But it is necessary for the identification point of view to make sense, that a single value of the *contrast function*  $x \mapsto \gamma(\theta; x)$  minimizes the loss. Remark that in the standard likelihood framework, this holds at least if any model contains the sample distribution, since it is the unique Kullback-Leibler divergence minimizer. Obviously, several parameter values, perhaps in different models, may represent it, besides the label switching. We do not know any such result with the  $-\log L_{cc}$  contrast and hypothesize that the assumption holds.

The assumption about the nonsingularity of  $I_{\theta^0}$  is unpleasant, since it is hard to be guaranteed. Hopefully, it could be weakened. The result of Massart (2007) (Theorem 7.11) which inspires this, and is available in a standard likelihood context, does not require such an assumption since it does not rely on the study of this link between the contrast and the parameters but on a clever choice of the involved distances (Hellinger distances), and on particular properties of the log function. However, this is a usual assumption (see Redner and Walker, 1984, or below).

Massart (2007) moreover does not require the contrast (i.e. the likelihood) to be bounded, as we have to. Remark however that the application of his Lemma 7.23 to obtain a genuine oracle inequality involves an assumption similar to the boundedness of the contrast. So that it seems reasonable

that the assumptions about  $M$  and  $M'$  (the last is much milder than the former) be necessary. They are typically ensured if either the contrast is bounded or if the support of  $f^\varphi$  is bounded.

The conditions about the penalty form are analogous to that of Nishii (1988) or Keribin (2000), which are both derived in the standard maximum likelihood framework. As those of Keribin (2000), they can be regarded as generalizing those of Nishii (1988) when the considered models are Gaussian mixture models. Indeed, Nishii (1988) considers penalties of the form  $c_n D_K$  and proves the model selection procedure to be weakly consistent if  $\frac{c_n}{n} \rightarrow 0$  and  $c_n \rightarrow \infty$ . Note that Nishii (1988) assumes the parameter space to be convex. He moreover notably assumes that  $\Theta_K^0 = \{\theta_K^0\}$  and that the counterpart of  $I_{\theta_K^0}$  is nonsingular, together with other regularity assumptions. Those results are not particularly designed for mixture models. Instead, as we do, Keribin (2000) considers general penalty forms and proves the procedure to be consistent if  $\frac{\text{pen}(K)}{n} \rightarrow 0$ ,  $\text{pen}(K) \rightarrow \infty$  and  $\liminf \frac{\text{pen}(K)}{\text{pen}(K')} > 1$  if  $K > K'$ . These conditions are equivalent to Nishii's if  $\text{pen}(K) = c_n D_K$ . In a general mixture model framework, she assumes the model family to be well-specified, the same notion of identifiability as we do, and a condition which does not seem to be directly comparable to ours about  $I_{\theta_K^0}$  but which tastes roughly the same. It might be milder. Those assumptions are proved to hold with the standard likelihood contrast for Gaussian mixture models with lower bounded, spherical covariance matrices which are the same for all components, and if the means belong to a compact. Our conditions about the penalty are a little weaker than Keribin's, but they still are quite analogous. Moreover, as compared to those results, we notably have to keep the proportions away from zero. This is necessary because the entropy term must be handled. It does not seem easy to extend the methods used by Keribin (2000) to our framework.

The strong version of Theorem 4.2, which would state the almost sure consistency of  $\hat{K}$  to  $K_0$ , would then probably involve penalties a little heavier, as Nishii (1988) and Keribin (2000) proved in their respective frameworks: both had to assume that  $\frac{\text{pen}(K)}{\log \log n} \rightarrow \infty$ .

Theorem 4.2 is a direct consequence of Theorem 4.1, Lemma 4.1, Theorem 3.2, which can be applied under those assumptions, and of Corollary 4.2 below and the discussion about its assumptions along the lines of Section 4.2.

## 4.2 Sufficient Conditions to Ensure Assumption (B4)

Let us introduce the notation  $S_n \gamma(\theta) = n(\gamma_n(\theta) - \mathbb{E}_{f^\varphi}[\gamma(\theta)])$ . The main result of this section is Lemma 4.2. Some intermediate results which en-

able to link Lemma 4.2 to Theorem 4.1 via Assumption (B4) are stated as corollaries and proved subsequently. Lemma 4.2 provides a control of  $\sup_{\theta \in \Theta} \frac{S_n(\gamma(\theta^0) - \gamma(\theta))}{\|\theta^0 - \theta\|_\infty^2 + \beta^2}$  (with respect to  $\beta$ ) and then of  $\frac{S_n(\gamma(\theta^0) - \gamma(\hat{\theta}))}{\|\theta^0 - \hat{\theta}\|_\infty^2 + \beta^2}$ . With a good choice of  $\beta$ , and if  $S_n(\gamma(\theta^0) - \gamma(\hat{\theta}))$  can be linked to  $\|\theta^0 - \hat{\theta}\|_\infty^2$ , it is proved in Corollary 4.1 that it may then be assessed that  $n\|\hat{\theta} - \theta^0\|_\infty^2 = O_{\mathbb{P}}(1)$ .

Plugging this last property back into the result of Lemma 4.2 yields (Corollary 4.2)  $n(\gamma_n(\theta_K^0) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1)$  for any model  $K \in \operatorname{argmin}_{1 \leq K \leq K_M} \mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)]$  and then, under mild identifiability condition,  $n(\gamma_n(\theta_{K_0}^0) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1)$ , which is Assumption (B4).

**Lemma 4.2.** *Let  $D \in \mathbb{N}^*$  and  $\Theta \subset \mathbb{R}^D$  convex. Let  $\Theta^\mathcal{O} \subset \mathbb{R}^D$  open such that  $\Theta \subset \Theta^\mathcal{O}$  and  $\gamma : \Theta^\mathcal{O} \times \mathbb{R}^d \rightarrow \mathbb{R}$ .  $\theta \in \Theta^\mathcal{O} \mapsto \gamma(\theta; x)$  is assumed to be  $C^1$  over  $\Theta^\mathcal{O}$  for  $f^\varphi d\lambda$ -almost all  $x$ . Let  $\theta^0 \in \Theta$  such that  $\mathbb{E}_{f^\varphi} [\gamma(\theta^0)] = \inf_{\theta \in \Theta} \mathbb{E}_{f^\varphi} [\gamma(\theta)]$ .*

*Assume that  $H_{\gamma, \Theta, \infty}^M$  and  $H_{\gamma, \Theta, 2}^{M'}$  hold.*

*Then  $\exists \alpha > 0 / \forall n, \forall \beta > 0, \forall \eta > 0$ , with probability larger than  $(1 - e^{-\eta})$ ,*

$$\sup_{\theta \in \Theta} \frac{S_n(\gamma(\theta^0) - \gamma(\theta))}{\|\theta^0 - \theta\|_\infty^2 + \beta^2} \leq \frac{\alpha}{\beta^2} \left( \|M'\|_2 \beta \sqrt{nD} + (\|M\|_\infty + \|M'\|_2 \beta) D \right. \\ \left. + \|M'\|_2 \sqrt{n\eta} \beta + \|M\|_\infty \eta \right)$$

*Note that  $\alpha$  is an absolute constant which notably does not depend on  $\theta^0$ .*

*Sketch of proof.* The proof relies on results of Massart (2007) and on the evaluation of the bracketing entropy of the class of functions at hand. Lemma 3.3 provides a local control of the entropy and hence, through Theorem 6.8 in Massart (2007), a control of the supremum of  $S_n(\gamma(\theta^0) - \gamma(\theta))$  as  $\|\theta - \theta^0\|_\infty^2 < \sigma$ , with respect to  $\sigma$ . The ‘‘peeling’’ Lemma 4.23 in Massart (2007) then enables to take advantage of this local control to derive a fine global control of  $\sup_{\theta \in \Theta} \frac{S_n(\gamma(\theta^0) - \gamma(\theta))}{\|\theta - \theta^0\|_\infty^2 + \beta^2}$ , for any  $\beta$ . This control in expectation, which can be derived conditionally to any event A, yields a control in probability thanks to Lemma 2.4 in Massart (2007), which can be thought of as an application of Markov’s inequality.  $\square$

**Corollary 4.1.** *Same assumptions as Lemma 4.2, but the convexity of  $\Theta$ . Besides assume that  $I_{\theta^0} = \frac{\partial^2}{\partial \theta^2} (\mathbb{E}_{f^\varphi} [\gamma(\theta)])|_{\theta^0}$  is nonsingular. Let  $(\hat{\theta}_n)_{n \geq 1}$  such that  $\hat{\theta}_n \in \Theta$ ,  $\gamma_n(\hat{\theta}_n) \leq \gamma_n(\theta^0) + O_{\mathbb{P}}(\frac{1}{n})$  and  $\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \theta^0$ . Then*

$$n\|\hat{\theta}_n - \theta^0\|_\infty^2 = O_{\mathbb{P}}(1).$$

The constant involved in  $O_{\mathbb{P}}(1)$  depends on  $D$ ,  $\|M\|_{\infty}$ ,  $\|M'\|_2$  and  $I_{\theta^0}$ .

This is a direct consequence of Lemma 4.2: it suffices to choose  $\beta$  well. The dependency of  $O_{\mathbb{P}}(1)$  in  $D$ ,  $\|M\|_{\infty}$ ,  $\|M'\|_2$  and  $I_{\theta^0}$  is not a problem since we aim at deriving an asymptotic result: the order of  $\|\theta - \theta^0\|_{\infty}^2$  with respect to  $n$  when the model is fixed is of concern.

The assumption that  $I_{\theta^0}$  is nonsingular plays an analogous role as Assumption (A2) in Theorem 3.1: this ensures that  $\mathbb{E}_{f^{\varphi}}[\gamma(\theta)]$  cannot be close to  $\mathbb{E}_{f^{\varphi}}[\gamma(\theta^0)]$  if  $\theta$  is not close to  $\theta^0$ . But this stronger assumption is necessary to strengthen the conclusion: the rate of the relation between  $\mathbb{E}_{f^{\varphi}}[\gamma(\theta)] - \mathbb{E}_{f^{\varphi}}[\gamma(\theta^0)]$  and  $\|\theta - \theta^0\|$  can then be controlled...

Should this assumption fail,  $\exists \tilde{\theta} \in \Theta / \tilde{\theta}' I_{\theta^0} \tilde{\theta} = 0 \Rightarrow \mathbb{E}_{f^{\varphi}}[\gamma(\theta^0 + \lambda \tilde{\theta})] = \mathbb{E}_{f^{\varphi}}[\gamma(\theta^0)] + o(\lambda^2)$  and then there is no hope to have  $\alpha > 0$  such that  $\mathbb{E}_{f^{\varphi}}[\gamma(\theta)] - \mathbb{E}_{f^{\varphi}}[\gamma(\theta^0)] > \alpha \|\theta - \theta^0\|^2$ : this approach cannot be applied without this — admittedly unpleasant — assumption. Perhaps an other approach (with distances not involving the parameters but directly the contrast values) might enable to avoid it, as Massart (2007) did in the likelihood framework.

**Corollary 4.2.** *Let  $(\Theta_K)_{1 \leq K \leq K_M}$  be models with, for any  $K$ ,  $\Theta_K \subset \mathbb{R}^{D_K}$ . Assume that  $D_1 \leq \dots \leq D_{K_M}$ . For any  $K$ , assume there exists an open set  $\Theta_K^{\mathcal{O}} \subset \mathbb{R}^{D_K}$  such that  $\Theta_K \subset \Theta_K^{\mathcal{O}}$  and such that with  $\Theta^{\mathcal{O}} = \Theta_1^{\mathcal{O}} \cup \dots \cup \Theta_{K_M}^{\mathcal{O}}$ ,  $\gamma : \Theta^{\mathcal{O}} \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined and  $C^1$  for  $f^{\varphi} d\lambda$ -almost all  $x$ . Assume that  $H_{\gamma, \Theta^{\mathcal{O}}, \infty}^M$  and  $H_{\gamma, \Theta^{\mathcal{O}}, 2}^{M'}$  hold. Let, for any  $K$ ,  $\Theta_K^0 = \operatorname{argmin}_{\theta \in \Theta_K} \mathbb{E}_{f^{\varphi}}[\gamma(\theta)]$  and  $\theta_K^0 \in \Theta_K^0$ .*

*Let  $K_0 = \min \operatorname{argmin}_{1 \leq K \leq K_M} \mathbb{E}_{f^{\varphi}}[\gamma(\Theta_K^0)]$  and assume  $\forall K, \forall \theta \in \Theta_K$ ,*

$$\mathbb{E}_{f^{\varphi}}[\gamma(\theta)] = \mathbb{E}_{f^{\varphi}}[\gamma(\theta_{K_0}^0)] \iff \gamma(\theta) = \gamma(\theta_{K_0}^0) \quad f^{\varphi} d\lambda - a.e.$$

*Let  $\mathcal{K} = \left\{ K \in \{1, \dots, K_M\} : \mathbb{E}_{f^{\varphi}}[\gamma(\theta_K^0)] = \mathbb{E}_{f^{\varphi}}[\gamma(\theta_{K_0}^0)] \right\}$ .*

*For any  $K \in \mathcal{K}$ , let  $\hat{\theta}_K \in \Theta_K$  such that*

$$\gamma_n(\hat{\theta}_K) \leq \gamma_n(\theta_K^0) + O_{\mathbb{P}}\left(\frac{1}{n}\right) \quad \text{and} \quad \hat{\theta}_K \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \theta_K^0.$$

*Assume that  $I_{\theta_K^0} = \frac{\partial^2}{\partial \theta^2} \left( \mathbb{E}_{f^{\varphi}}[\gamma(\theta)] \right)_{|\theta = \theta_K^0}$  is nonsingular for any  $K \in \mathcal{K}$ .*

*Then  $\forall K \in \mathcal{K}$ ,  $n(\gamma_n(\hat{\theta}_{K_0}) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1)$ .*

This last corollary states conditions under which assumption (B4) of Theorem 4.1 is ensured.

### 4.3 A New Light on ICL

The previous section suggests links between model selection penalized criteria with the standard likelihood on the one hand and with the conditional classification likelihood we defined on the other hand. Indeed penalties with the same form as those given by Nishii (1988) or Keribin (2000) with the standard likelihood are proved to be “consistent” in our framework. Therefore, by analogy with the standard likelihood framework, it is expected that penalties proportional to  $D_K$  conform an efficiency point of view (think of AIC), and that penalties proportional to  $D_K \log n$  are optimal for an identification purpose (think of BIC). This possibility to derive an identification procedure from an efficient procedure by a  $\log n$  factor is notified for example by Arlot (2007).

Let us then consider by analogy with BIC the penalized criterion

$$\text{crit}_{L_{cc}\text{-ICL}}(K) = \log L_{cc}(\hat{\theta}_K^{\text{MLccE}}) - \frac{\log n}{2} D_K.$$

The point is that we almost recover ICL (replace  $\hat{\theta}_K^{\text{MLE}}$  by  $\hat{\theta}_K^{\text{MLccE}}$  in (2)), which may then be regarded as an approximation of this  $L_{cc}$ -ICL criterion. The corresponding penalty is  $\frac{\log n}{2} D_K$ , and the derivation of  $L_{cc}$ -ICL illustrates that the entropy should not be considered as a part of the penalty. This notably justifies why ICL does not select the same number of components as BIC or any consistent criterion in the standard likelihood framework, even asymptotically. Actually, it should not be expected to do so.

When  $\hat{\theta}_K^{\text{MLccE}}$  differs from  $\hat{\theta}_K^{\text{MLE}}$ , the former provides more separated clusters. The compromise between the Gaussian component and the cluster viewpoint is achieved with  $\hat{\theta}_K^{\text{MLccE}}$  from the very estimation step. The user is provided a solution which aims at this compromise for each number of classes  $K$ . However, the number of classes selected through  $L_{cc}$ -ICL differs seldom from the one selected by ICL in simulations (See Baudry, 2009, Chapter 4).

Finally,  $L_{cc}$ -ICL is quite close to ICL and enables to better understand the concepts underlying ICL. ICL remains attractive though, notably because it is easier to implement than  $L_{cc}$ -ICL.

## 5 Discussion

Two families of criteria, in the clustering framework, are distinguished: it is shown that ICL’s purpose is of different nature than that of BIC or AIC. The identification theory for the criteria based on the conditional classification likelihood is — not surprisingly — very similar to the one for the



standard likelihood. A major interest of the newly introduced estimator and criteria is to better understand the ICL criterion and the underlying notion of class. This is not a simple notion of cluster — as for example for the k-means procedure — neither a pure notion of “component” — as underlying the MLE/BIC approach — but a compromise between both. ICL leads to discovering classes matching a subtle combination of the notions of well separated, compact, clusters, and (Gaussian) mixture components. It then enjoys the flexibility and modeling possibilities of the model-based clustering approach, but does not break the expected notion of cluster. Better understanding of the ICL criterion now means better understanding the newly involved contrast  $L_{cc}$ .

The choice of the involved mixture components must be handled with care in this framework since it leads the cluster shape underlying the study. Several forms of Gaussian mixtures may be involved: for example, spherical and general models may be compared, or models with free proportions may be compared with models with equal proportions.

Besides it should be further studied how the complexity of the models should be measured when several model kinds are compared. The dimension of the model as a parametric space works for the reported theoretical results. But we are not completely convinced that it is the finest measure of the complexity of Gaussian mixture models. As a matter of fact this simple parametric point of view amounts to considering that all parameters play an analogous role. This is not really natural.

A further theoretical step would be to derive non-asymptotic results and oracle inequalities. This may give more precise insights about the best penalty shape to use, and then justify the use of the slope heuristics of Birgé and Massart (2007) (see also Baudry et al., 2011 or Baudry, 2009 for simulations and discussions on this topic).

A practical challenge is to provide efficient optimization algorithms. Some work has been done in this direction already: see Baudry et al. (2008) and Baudry (2009, Section 5.1). But they need be improved to be more reliable, and above all to run much faster, which would obviously be a condition for a spread practical use of the new contrast.

A possibility to make this contrast more flexible would be to assign different weights to the log likelihood and the entropy:  $\log L_{cc_\alpha} = \alpha \log L + (1 - \alpha) \text{ENT}$ , with  $\alpha \in [0; 1]$ . This would enable to tune how important the assignment confidence is with respect to the Gaussian fit... The difficulty would then be to choose  $\alpha$ . A first insight which comes in mind is to calibrate  $\alpha$  from simulations of situations in which the user knows what solution he expects.

## 6 Proofs

*Proof (Max of  $h_K$ , p. 7).* If  $h_K$  reaches a max. value at  $(t_1^0, \dots, t_K^0)$  under the constraint  $\sum_{k=1}^K t_k^0 = 1$  then, with  $S : (t_1, \dots, t_K) \mapsto \sum_{k=1}^K t_k$ ,

$$\exists \lambda \in \mathbb{R} / dh_K(t_1^0, \dots, t_K^0) = \lambda dS(t_1^0, \dots, t_K^0).$$

This is equivalent to  $\forall k, \log t_k^0 + 1 = \lambda$ . Then,  $\forall k, \forall k', t_k^0 = t_{k'}^0$  and since  $\sum_{k=1}^K t_k^0 = 1$ , this yields  $t_k^0 = \frac{1}{K}$ .  $\square$

*Proof of Theorem 3.1.* Let  $\varepsilon > 0$  and  $\eta = \inf_{d(\theta, \Theta^0) > \varepsilon} \mathbb{E}_{f^\varphi} [\gamma(\theta)] - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] > 0$  (from assumption (A2)). For  $n$  large enough and with large probability, from assumption (A3) and the definition of  $\hat{\theta}$ ,

$$\sup_{\theta \in \Theta} |\gamma_n(\theta) - \mathbb{E}_{f^\varphi} [\gamma(\theta)]| < \frac{\eta}{3} \quad \text{and} \quad \gamma_n(\hat{\theta}) \leq \gamma_n(\theta^0) + \frac{\eta}{3}.$$

Then

$$\begin{aligned} \mathbb{E}_{f^\varphi} [\gamma(\hat{\theta})] - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] &\leq \mathbb{E}_{f^\varphi} [\gamma(\hat{\theta})] - \gamma_n(\hat{\theta}) + \gamma_n(\hat{\theta}) - \gamma_n(\theta^0) \\ &\quad + \gamma_n(\theta^0) - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] \\ &< \eta. \end{aligned}$$

And  $d(\hat{\theta}, \Theta^0) < \varepsilon$  with great probability, as  $n$  is large enough.  $\square$

*Proof of Lemma 3.1.* Let  $\varepsilon > 0$ , and  $\tilde{\Theta} \subset \Theta$ , with  $\tilde{\Theta}$  bounded. Let  $\tilde{\Theta}_\varepsilon$  be a grid in  $\Theta$  which “ $\varepsilon$ -covers”  $\tilde{\Theta}$  in any dimension with step  $\varepsilon$ .  $\tilde{\Theta}_\varepsilon$  is for example  $\tilde{\Theta}_\varepsilon^1 \times \dots \times \tilde{\Theta}_\varepsilon^D$  with

$$\forall i \in \{1, \dots, D\}, \tilde{\Theta}_\varepsilon^i = \left\{ \tilde{\theta}_{\min}^i, \tilde{\theta}_{\min}^i + \varepsilon, \dots, \tilde{\theta}_{\max}^i \right\},$$

where

$$\forall i \in \{1, \dots, D\}, \left\{ \theta^i : \theta \in \tilde{\Theta} \right\} \subset \left[ \tilde{\theta}_{\min}^i - \frac{\varepsilon}{2}, \tilde{\theta}_{\max}^i + \frac{\varepsilon}{2} \right].$$

This is always possible since  $\Theta$  is convex. For the sake of simplicity, it is assumed without loss of generality, that  $\tilde{\Theta}_\varepsilon \subset \tilde{\Theta}$ . With the  $\|\cdot\|_\infty$  norm, the step of the grid  $\tilde{\Theta}_\varepsilon$  is the same as the step over each dimension,  $\varepsilon$ :

$$\forall \tilde{\theta} \in \tilde{\Theta}, \exists \tilde{\theta}_\varepsilon \in \tilde{\Theta}_\varepsilon / \|\tilde{\theta} - \tilde{\theta}_\varepsilon\|_\infty \leq \frac{\varepsilon}{2}.$$

And the cardinal of  $\tilde{\Theta}_\varepsilon$  is at most

$$\prod_{i=1}^D \frac{(\sup_{\theta \in \tilde{\Theta}} \theta^i - \inf_{\theta \in \tilde{\Theta}} \theta^i)}{\varepsilon} \vee 1 \leq \left( \frac{\text{diam } \tilde{\Theta}}{\varepsilon} \right)^D \vee 1.$$

Now, let  $\theta_1$  and  $\theta_2$  in  $\Theta$  and  $x \in \mathbb{R}^d$ .

$$\begin{aligned} |\gamma(\theta_1; x) - \gamma(\theta_2; x)| &\leq \sup_{\theta \in [\theta_1; \theta_2]} \left\| \left( \frac{\partial \gamma}{\partial \theta} \right)_{(\theta; x)} \right\|_{\infty} \|\theta_1 - \theta_2\|_{\infty} \\ &\leq \underbrace{\sup_{\theta \in \Theta} \left\| \left( \frac{\partial \gamma}{\partial \theta} \right)_{(\theta; x)} \right\|_{\infty}}_{M'(x)} \|\theta_1 - \theta_2\|_{\infty}, \end{aligned}$$

since  $\Theta$  is convex. Let  $\tilde{\theta} \in \tilde{\Theta}$  and choose  $\tilde{\theta}_{\varepsilon} \in \tilde{\Theta}_{\varepsilon}$  such that  $\|\tilde{\theta} - \tilde{\theta}_{\varepsilon}\|_{\infty} \leq \frac{\varepsilon}{2}$ . Then

$$\begin{aligned} \forall x \in \mathbb{R}^d, |\gamma(\tilde{\theta}_{\varepsilon}; x) - \gamma(\tilde{\theta}; x)| &\leq M'(x) \frac{\varepsilon}{2} \\ \text{and} \\ \gamma(\tilde{\theta}_{\varepsilon}; x) - \frac{\varepsilon}{2} M'(x) &\leq \gamma(\tilde{\theta}; x) \leq \gamma(\tilde{\theta}_{\varepsilon}; x) + \frac{\varepsilon}{2} M'(x). \end{aligned}$$

The set of  $\varepsilon \|M'\|_r$ -brackets (for the  $\|\cdot\|_r$ -norm)

$$\left\{ \left[ \gamma(\tilde{\theta}_{\varepsilon}) - \frac{\varepsilon}{2} M'; \gamma(\tilde{\theta}_{\varepsilon}) + \frac{\varepsilon}{2} M' \right] : \tilde{\theta}_{\varepsilon} \in \tilde{\Theta}_{\varepsilon} \right\}$$

then has cardinal at most  $\left( \frac{\text{diam} \tilde{\Theta}}{\varepsilon} \right)^D \vee 1$  and covers  $\left\{ \gamma(\tilde{\theta}) : \tilde{\theta} \in \tilde{\Theta} \right\}$ .  $\square$

**Example 6.1** (Diagonal Gaussian Mixture Model Parameter Space is Convex). *Following Celeux and Govaert (1995), we write  $[p\lambda_k B_k]$  for the model of Gaussian mixtures with diagonal covariance matrices and equal mixing proportions. To keep simple notation, let us consider the case  $d = 2$  and  $K = 2$  ( $d = 1$  or  $K = 1$  are obviously particular cases!). A natural parametrization of this model (which dimension is 8) is*

$$\theta \in \mathbb{R}^4 \times \mathbb{R}^{+*4} \xrightarrow{\varphi} \frac{1}{2} \phi \left( \cdot; \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}, \begin{pmatrix} \theta_5 & 0 \\ 0 & \theta_6 \end{pmatrix} \right) + \frac{1}{2} \phi \left( \cdot; \begin{pmatrix} \theta_3 \\ \theta_4 \end{pmatrix}, \begin{pmatrix} \theta_7 & 0 \\ 0 & \theta_8 \end{pmatrix} \right)$$

Then  $[p\lambda_k B_k] = \varphi(\mathbb{R}^4 \times \mathbb{R}^{+*4})$  and the parameter space  $\mathbb{R}^4 \times \mathbb{R}^{+*4}$  is convex.

**Example 6.2** (The Same Model with Equal Volumes is Convex, too...).  *$[p\lambda B_k]$  is the same model as in the previous example, but the covariance matrices determinants have to be equal. With  $d = 2$  and  $K = 2$ , a natural*

parametrization of this model of dimension 7 is

$$\begin{aligned} \theta \in \mathbb{R}^4 \times \mathbb{R}^{+*3} \xrightarrow{\varphi} & \frac{1}{2} \phi \left( \cdot; \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}, \sqrt{\theta_7} \begin{pmatrix} \theta_5 & 0 \\ 0 & \frac{1}{\theta_5} \end{pmatrix} \right) \\ & + \frac{1}{2} \phi \left( \cdot; \begin{pmatrix} \theta_3 \\ \theta_4 \end{pmatrix}, \sqrt{\theta_7} \begin{pmatrix} \theta_6 & 0 \\ 0 & \frac{1}{\theta_6} \end{pmatrix} \right) \end{aligned}$$

Then  $[p\lambda B_k] = \varphi(\mathbb{R}^4 \times \mathbb{R}^{+*3})$  and the parameter space  $\mathbb{R}^4 \times \mathbb{R}^{+*3}$  is convex.

*Proof of Lemma 3.2.* Let  $O_1, \dots, O_Q$  be a finite covering of  $\Theta$  consisting of open balls such that  $\cup_{q=1}^Q O_q \subset \Theta^\circ$ . Such a covering always exists since  $\Theta$  is assumed to be compact. Remark that

$$\Theta = \cup_{q=1}^Q (O_q \cap \Theta) \subset \cup_{q=1}^Q \text{conv}(O_q \cap \Theta).$$

Now, for any  $q$ ,  $\text{conv}(O_q \cap \Theta)$  is convex and  $\sup_{\theta \in \text{conv}(O_q \cap \Theta)} \left\| \left( \frac{\partial \gamma}{\partial \theta} \right)_{(\theta; x)} \right\|_\infty \leq M'(x)$  since  $\text{conv}(O_q \cap \Theta) \subset O_q \subset \Theta^\circ$ . Therefore, Lemma 3.1 applies to  $O_q \cap \tilde{\Theta} \subset \text{conv}(O_q \cap \Theta)$ :

$$N_{[\cdot]}(\varepsilon, \{\gamma(\theta) : \theta \in \tilde{\Theta} \cap O_q\}, \|\cdot\|_r) \leq \left( \frac{\|M'\|_r \text{diam } \tilde{\Theta}}{\varepsilon} \right)^D \vee 1.$$

Since  $N_{[\cdot]}(\varepsilon, \{\gamma(\theta) : \theta \in \tilde{\Theta}\}, \|\cdot\|_r) \leq N_{[\cdot]}(\varepsilon, \cup_{q=1}^Q \{\gamma(\theta) : \theta \in \tilde{\Theta} \cap O_q\}, \|\cdot\|_r)$ , the result follows.  $\square$

*Proof of Lemma 3.3.* Consider the grid  $\tilde{\Theta}_\varepsilon$  of the proof of Lemma 3.1. Let  $\theta_1$  and  $\theta_2$  in  $\Theta$  and  $x \in \mathbb{R}^d$ . Since  $\Theta$  is convex,

$$\begin{aligned} \left| \gamma(\theta_1; x) - \gamma(\theta_2; x) \right|^r & \leq \sup_{\theta \in [\theta_1; \theta_2]} \left\| \left( \frac{\partial \gamma}{\partial \theta} \right)_{(\theta; x)} \right\|_\infty^2 \|\theta_1 - \theta_2\|_\infty^2 \\ & \quad \times \left( 2 \sup_{\theta \in \{\theta_1, \theta_2\}} |\gamma(\theta; x)| \right)^{r-2} \\ & \leq M'(x)^2 \|\theta_1 - \theta_2\|_\infty^2 (2\|M\|_\infty)^{r-2} f^\varphi d\lambda\text{-a.e.} \end{aligned}$$

Let  $\tilde{\theta} \in \tilde{\Theta}$  and choose  $\tilde{\theta}_\varepsilon \in \tilde{\Theta}_\varepsilon$  such that  $\|\tilde{\theta} - \tilde{\theta}_\varepsilon\|_\infty \leq \frac{\varepsilon}{2}$ . Then

$$\left| \gamma(\tilde{\theta}_\varepsilon; x) - \gamma(\tilde{\theta}; x) \right| \leq M'(x)^{\frac{2}{r}} \left( \frac{\varepsilon}{2} \right)^{\frac{2}{r}} (2\|M\|_\infty)^{\frac{r-2}{r}} f^\varphi\text{-a.e.}$$

and the set of brackets

$$\left\{ \left[ \gamma(\tilde{\theta}_\varepsilon; x) - \varepsilon^{\frac{2}{r}} M'(x)^{\frac{2}{r}} \|M\|_\infty^{\frac{r-2}{r}} 2^{1-\frac{4}{r}}; \gamma(\tilde{\theta}_\varepsilon; x) + \varepsilon^{\frac{2}{r}} M'(x)^{\frac{2}{r}} \|M\|_\infty^{\frac{r-2}{r}} 2^{1-\frac{4}{r}} \right] : \tilde{\theta} \in \tilde{\Theta}_\varepsilon \right\}$$

(of  $\|\cdot\|_r$ -norm length  $(2^{2-\frac{4}{r}} \|M\|_\infty^{\frac{r-2}{r}} \|M'\|_2^{\frac{2}{r}} \varepsilon^{\frac{2}{r}})$  has cardinal at most  $\left(\frac{\text{diam } \tilde{\Theta}}{\varepsilon}\right)^D \vee 1$  and covers  $\{\gamma(\tilde{\theta}) : \tilde{\theta} \in \tilde{\Theta}\}$ , which yields Lemma 3.3.  $\square$

*Proof of Theorem 4.1.* Let  $\mathcal{K} = \text{argmin}_{1 \leq K \leq K_M} \mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)]$ . By assumption,  $K_0 = \min \mathcal{K}$ .

It is first proved that  $\hat{K}$  does not asymptotically “underestimate”  $K_0$ . Let  $K \notin \mathcal{K}$ . Let  $\varepsilon = \frac{1}{2} (\mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] - \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)]) > 0$ . From (B2) and (B3) ( $\text{pen}(K) = o_{\mathbb{P}}(1)$ ), with large probability and for  $n$  large enough:

$$|\gamma_n(\hat{\theta}_K) - \mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)]| \leq \frac{\varepsilon}{3} \quad |\gamma_n(\hat{\theta}_{K_0}) - \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)]| \leq \frac{\varepsilon}{3} \quad \text{pen}(K_0) \leq \frac{\varepsilon}{3}.$$

Then

$$\begin{aligned} \text{crit}(K) &= \gamma_n(\hat{\theta}_K) + \text{pen}(K) \geq \mathbb{E}_{f^\varphi} [\gamma(\theta_K^0)] - \frac{\varepsilon}{3} + 0 \\ &= \mathbb{E}_{f^\varphi} [\gamma(\theta_{K_0}^0)] + \frac{5\varepsilon}{3} \geq \underbrace{\gamma_n(\hat{\theta}_{K_0}) + \text{pen}(K_0)}_{\text{crit}(K_0)} + \varepsilon. \end{aligned}$$

Then, with large probability and for  $n$  large enough,  $\hat{K} \neq K$ .

Let now  $K \in \mathcal{K}$  (hence  $K > K_0$ ). This part of the result is more involved than the first one but at this stage, it is not more difficult to derive: all the difficulty is hidden in the strong assumption (B4)... Indeed, it implies that  $\exists V > 0$ , such that for  $n$  large enough and with large probability,

$$n(\gamma_n(\hat{\theta}_{K_0}) - \gamma_n(\hat{\theta}_K)) \leq V.$$

Increase  $n$  enough so that  $n(\text{pen}(K) - \text{pen}(K_0)) > V$  with large probability (which is possible from assumption (B4)). Then, for  $n$  large enough and with large probability,

$$\text{crit}(K) = \gamma_n(\hat{\theta}_K) + \text{pen}(K) \geq \gamma_n(\hat{\theta}_{K_0}) - \frac{V}{n} + \text{pen}(K) > \text{crit}(K_0).$$

And then, with large probability and for  $n$  large enough,  $\hat{K} \neq K$ .

Finally, since  $\mathbb{P}[\hat{K} \neq K_0] = \sum_{K \notin \mathcal{K}} \mathbb{P}[\hat{K} = K] + \sum_{K \in \mathcal{K}, K \neq K_0} \mathbb{P}[\hat{K} = K]$ , the result follows.  $\square$

*Proof of Lemma 4.1.* For any  $\varepsilon > 0$ , with large probability and for  $n$  large enough:

$$\underbrace{\gamma_n(\hat{\theta}) - \mathbb{E}_{f^\varphi} [\gamma(\hat{\theta})]}_{\geq -\varepsilon} + \underbrace{\mathbb{E}_{f^\varphi} [\gamma(\hat{\theta})] - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)]}_{\geq 0} = \gamma_n(\hat{\theta}) - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)] \\ = \underbrace{\gamma_n(\hat{\theta}) - \gamma_n(\theta^0)}_{\geq -\varepsilon} + \underbrace{\gamma_n(\theta^0) - \mathbb{E}_{f^\varphi} [\gamma(\theta^0)]}_{\geq 0}. \quad \square$$

*Proof of Lemma 4.2.* Actually, the proof as it is written below holds for an at most countable model (because this assumption is necessary for Lemma 4.23 and Theorem 6.8 in Massart (2007) to hold). But it can be checked that both those results may be applied to a dense subset of  $\{\gamma(\theta) : \theta \in \Theta\}$  containing  $\theta^0$  and their respective conclusions generalized to the entire set: choose  $\Theta^{\text{count}}$  a countable dense subset of  $\Theta$ . Then, for any  $\theta \in \Theta$ , let  $\theta_n \in \Theta^{\text{count}} \xrightarrow[n \rightarrow \infty]{} \theta$ . Then,  $\gamma(\theta_n; X) \xrightarrow[n \rightarrow \infty]{a.s.} \gamma(\theta; X)$ . Now, whatever  $g : \mathbb{R}^D \times (\mathbb{R}^d)^n \rightarrow \mathbb{R}$  such that  $\theta \in \mathbb{R}^D \mapsto g(\theta, \mathbf{X})$  continue a.s.,  $\sup_{\theta \in \Theta} g(\theta; \mathbf{X}) = \sup_{\theta \in \Theta^{\text{count}}} g(\theta; \mathbf{X})$  a.s. Hence,  $\mathbb{E}_{f^\varphi} [\sup_{\theta \in \Theta} g(\theta; \mathbf{X})] = \mathbb{E}_{f^\varphi} [\sup_{\theta \in \Theta^{\text{count}}} g(\theta; \mathbf{X})]$ . Remark however that the models which are actually considered are discrete, because of the computation limitations.

Let us introduce the centered empirical process

$$S_n \gamma(\theta) = n\gamma_n(\theta) - n\mathbb{E}_{f^\varphi} [\gamma(\theta; X)].$$

Here and hereafter,  $\alpha$  stands for a generic absolute constant, which may differ from a line to an other. Let  $\theta^0 \in \Theta$  such that  $\mathbb{E}_{f^\varphi} [\gamma(\theta^0)] = \inf_{\theta \in \Theta} \mathbb{E}_{f^\varphi} [\gamma(\theta)]$ .

Let us define

$$\forall \sigma > 0, \Theta(\sigma) = \{\theta \in \Theta : \|\theta - \theta^0\|_\infty \leq \sigma\}.$$

On the one hand, for all  $r \in \mathbb{N}^* \setminus \{1\}$ ,

$$\forall \theta \in \Theta(\sigma), |\gamma(\theta^0; x) - \gamma(\theta; x)|^r \leq M'(x)^2 \|\theta^0 - \theta\|_\infty^2 (2M(x))^{r-2}$$

since  $\Theta(\sigma) \subset \Theta$  is convex. And thus,

$$\forall \theta \in \Theta(\sigma), \mathbb{E}_{f^\varphi} [|\gamma(\theta^0) - \gamma(\theta)|^r] \leq \|M'\|_2^2 \|\theta^0 - \theta\|_\infty^2 (2\|M\|_\infty)^{r-2} \\ \leq \frac{r!}{2} (\|M'\|_2 \sigma)^2 \left( \frac{2\|M\|_\infty}{2} \right)^{r-2}. \quad (5)$$

On the other hand, from Lemma 3.3, for any  $r \in \mathbb{N}^* \setminus \{1\}$ , for any  $\delta > 0$ , there exists  $C_\delta$  a set of brackets which cover  $\{(\gamma(\theta^0) - \gamma(\theta)) : \theta \in \Theta(\sigma)\}$  (deduced from a set of brackets which cover  $\{\gamma(\theta) : \theta \in \Theta(\sigma)\}$ ...) such that:

$$\forall r \in \mathbb{N}^* \setminus \{1\}, \forall [g_l, g_u] \in C_\delta, \|g_u - g_l\|_r \leq \left( \frac{r!}{2} \right)^{\frac{1}{r}} \delta^{\frac{2}{r}} \left( \frac{4\|M\|_\infty}{3} \right)^{\frac{r-2}{r}}$$

and such that, writing  $e^{H(\delta, \Theta(\sigma))}$  the minimal cardinal of such a  $C_\delta$ ,

$$e^{H(\delta, \Theta(\sigma))} \leq \left( \frac{\overbrace{\text{diam } \Theta(\sigma)}^{\leq 2\sigma} \|M'\|_2}{\delta} \right)^D \vee 1. \quad (6)$$

Then, according to Theorem 6.8 in Massart (2007),  
 $\exists \alpha, \forall \varepsilon \in ]0, 1], \forall A$  measurable such that  $\mathbb{P}[A] > 0$ ,

$$\begin{aligned} \mathbb{E}^A \left[ \sup_{\theta \in \Theta(\sigma)} S_n(\gamma(\theta^0) - \gamma(\theta)) \right] &\leq \frac{\alpha}{\varepsilon} \sqrt{n} \int_0^{\varepsilon \|M'\|_{2\sigma}} \sqrt{H(u, \Theta(\sigma))} du \\ &\quad + 2 \left( \frac{4}{3} \|M\|_\infty + \|M'\|_{2\sigma} \right) H(\|M'\|_{2\sigma}, \Theta(\sigma)) \\ &\quad + (1 + 6\varepsilon) \|M'\|_{2\sigma} \sqrt{2n \log \frac{1}{\mathbb{P}[A]}} + \frac{8}{3} \|M\|_\infty \log \frac{1}{\mathbb{P}[A]}. \end{aligned} \quad (7)$$

Now, we have

$$\begin{aligned} \forall t \in \mathbb{R}^+, \int_0^t \sqrt{\log \frac{1}{u} \vee 0} du &= \int_0^{t \wedge 1} \sqrt{\log \frac{1}{u}} du \\ &\leq \sqrt{t \wedge 1} \sqrt{\int_0^{t \wedge 1} \log \frac{1}{u} du} = (t \wedge 1) \sqrt{\log \frac{e}{t \wedge 1}}, \end{aligned}$$

by the Cauchy-Schwarz inequality. Together with (6), this yields

$$\begin{aligned} \forall t \in \mathbb{R}^+, \int_0^t \sqrt{H(u, \Theta(\sigma))} du &\leq \sqrt{D} \int_0^t \sqrt{\log \frac{2\|M'\|_{2\sigma}}{u} \vee 0} du \\ &\leq \sqrt{D} (t \wedge 2\|M'\|_{2\sigma}) \sqrt{\log \frac{e}{\frac{t}{2\|M'\|_{2\sigma}} \wedge 1}}, \end{aligned} \quad (8)$$

after a simple substitution.

Next, let us apply Lemma 4.23 in Massart (2007): From (6), (7) and (8),

$$\forall \sigma > 0, \mathbb{E}_{f^\varphi} \left[ \sup_{\theta \in \Theta(\sigma)} S_n(\gamma(\theta^0) - \gamma(\theta)) \right] \leq \varphi(\sigma),$$

$$\begin{aligned} \text{with } \varphi(t) = & \frac{\alpha}{\beta} \sqrt{n} \sqrt{D} \beta \|M'\|_2 t \sqrt{\log \frac{2e}{\varepsilon}} + 2 \left( \frac{4}{3} \|M\|_\infty + \|M'\|_2 t \right) D \log 2 \\ & + (1 + 6\varepsilon) \|M'\|_2 t \sqrt{2n \log \frac{1}{\mathbb{P}[A]}} + \frac{8}{3} \|M\|_\infty \log \frac{1}{\mathbb{P}[A]}. \end{aligned}$$

As required for Lemma 4.23 in Massart (2007) to hold,  $\frac{\varphi(t)}{t}$  is nonincreasing. It follows

$$\forall \beta > 0, \mathbb{E}^A \left[ \sup_{\theta \in \Theta} \frac{S_n(\gamma(\theta^0) - \gamma(\theta))}{\|\theta^0 - \theta\|_\infty + \beta^2} \right] \leq 4\beta^{-2} \varphi(\beta).$$

We then choose  $\varepsilon = 1$  and apply Lemma 2.4 in Massart (2007): for any  $\eta > 0$  and any  $\beta > 0$ , with probability larger than  $1 - e^{-\eta}$ ,

$$\begin{aligned} \sup_{\theta \in \Theta} \frac{S_n(\gamma(\theta^0) - \gamma(\theta))}{\|\theta^0 - \theta\|_\infty^2 + \beta^2} \leq & \frac{\alpha}{\beta^2} \left( \sqrt{nD} \|M'\|_2 \beta \sqrt{\log 2e} \right. \\ & \left. + (\|M\|_\infty + \|M'\|_2 \beta) D \log 2 + \|M'\|_2 \beta \sqrt{n\eta} + \|M\|_\infty \eta \right). \square \end{aligned}$$

*Proof of Corollary 4.1.* Let  $\varepsilon > 0$  such that  $B(\theta^0, \varepsilon) \subset \Theta^\mathcal{O}$ . Then, since  $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta^0$ , there exists  $n_0 \in \mathbb{N}^*$  such that, with large probability, for  $n \geq n_0$ ,  $\hat{\theta}_n \in B(\theta^0, \varepsilon)$ . Now,  $B(\theta^0, \varepsilon)$  is convex and the assumptions of the corollary guarantee that Lemma 4.2 applies. Let us apply it to  $\hat{\theta}_n$ :  $\forall n \geq n_0, \forall \beta > 0$ , with great probability as  $\eta$  is large,

$$\begin{aligned} \frac{S_n(\gamma(\theta^0) - \gamma(\hat{\theta}_n))}{\|\theta^0 - \hat{\theta}_n\|_\infty^2 + \beta^2} \leq & \frac{\alpha}{\beta^2} \left( \sqrt{nD} \|M'\|_2 \beta + (\|M\|_\infty + \|M'\|_2 \beta) D \right. \\ & \left. + \|M'\|_2 \beta \sqrt{n\eta} + \|M\|_\infty \eta \right). \quad (9) \end{aligned}$$

But since  $I_{\theta^0}$  is supposed to be nonsingular,  $\forall \theta \in B(\theta^0, \varepsilon)$ ,

$$\begin{aligned} \mathbb{E}_{f^\varphi}[\theta] - \mathbb{E}_{f^\varphi}[\theta^0] &= (\theta - \theta^0)' I_{\theta^0} (\theta - \theta^0) + r(\|\theta - \theta^0\|_\infty) \|\theta - \theta^0\|_\infty^2 \\ &\geq (2\alpha' + r(\|\theta - \theta^0\|_\infty)) \|\theta - \theta^0\|_\infty^2 \end{aligned}$$

where  $\alpha' > 0$  depends on  $I_{\theta^0}$  and  $r : \mathbb{R}^+ \rightarrow \mathbb{R}$  fulfills  $r(x) \xrightarrow{x \rightarrow 0} 0$ . Then, for  $\|\theta - \theta^0\|_\infty$  small enough ( $\varepsilon$  may be decreased...),

$$\forall \theta \in B(\theta^0, \varepsilon), \mathbb{E}_{f^\varphi}[\theta] - \mathbb{E}_{f^\varphi}[\theta^0] \geq \alpha' \|\theta - \theta^0\|_\infty^2. \quad (10)$$



$$\begin{aligned} S_n(\gamma(\theta^0) - \gamma(\hat{\theta}_n)) &= n(\gamma_n(\theta^0) - \gamma_n(\hat{\theta}_n)) + n\mathbb{E}_{f^\varphi}[\gamma(\hat{\theta}_n) - \gamma(\theta^0)] \quad \text{Since} \\ &\geq O_{\mathbb{P}}(1) + n\mathbb{E}_{f^\varphi}[\gamma(\hat{\theta}_n) - \gamma(\theta^0)], \end{aligned}$$

(9) together with (10) leads (with great probability) to

$$n\|\hat{\theta}_n - \theta^0\|_\infty^2 \leq \frac{\|M'\|_2(\sqrt{nD} + \sqrt{\eta n} + D)\beta + \|M\|_\infty(D + \eta) + O_{\mathbb{P}}(1)}{\frac{\alpha'}{\alpha} - \frac{1}{n\beta^2} \left( \|M'\|_2(\sqrt{nD} + \sqrt{\eta n} + D)\beta + \|M\|_\infty(D + \eta) \right)},$$

as soon as the denominator of the right-hand side is positive. It then suffices to choose  $\beta$  such that this condition is fulfilled and such that the right-hand side is upper-bounded by a quantity which does not depend on  $n$  to get the result. Let us try  $\beta = \frac{\beta_0}{\sqrt{n}}$  with  $\beta_0$  independent of  $n$ :

$$n\|\hat{\theta}_n - \theta^0\|_\infty^2 \leq \frac{\|M'\|_2(\sqrt{D} + \sqrt{\eta} + D)\beta_0 + \|M\|_\infty(D + \eta) + O_{\mathbb{P}}(1)}{\frac{\alpha'}{\alpha} - \frac{1}{\beta_0^2} \left( \|M'\|_2(\sqrt{D} + \sqrt{\eta} + D)\beta_0 + \|M\|_\infty(D + \eta) \right)}.$$

This only holds if the denominator is positive. Choose  $\beta_0$  large enough so as to guarantee this, which is always possible. The result follows: with large probability and for  $n$  larger than  $n_0$ , we have  $n\|\hat{\theta}_n - \theta^0\|_\infty^2 = CO_{\mathbb{P}}(1)$  with  $C$  depending on  $D$ ,  $\|M\|_\infty$ ,  $\|M'\|_2$ ,  $I_{\theta^0}$  and  $\eta$ .  $\square$

*Proof of Corollary 4.2.* This is a direct application of Corollary 4.1. Let  $K \in \mathcal{K}$ :  $\mathbb{E}_{f^\varphi}[\gamma(\theta_K^0)] = \mathbb{E}_{f^\varphi}[\gamma(\theta_{K_0}^0)]$ .  $\Theta_K$  can be assumed to be convex: if it is not,  $\hat{\theta}_K$  lies in  $B(\theta_{K_0}^0, \varepsilon) \subset \Theta^O$  with large probability for large  $n$  and  $\Theta_K$  may be replaced by  $B(\theta_{K_0}^0, \varepsilon)$ . According to Lemma 4.2, with probability larger than  $(1 - e^{-\eta})$  for  $n$  large, with  $\beta = \frac{\beta_0}{\sqrt{n}}$  for any  $\beta_0 > 0$ :

$$\begin{aligned} S_n(\gamma(\theta_K^0) - \gamma(\hat{\theta}_K)) &\leq \alpha \frac{n\|\theta_K^0 - \hat{\theta}_K\|_\infty^2 + \beta_0^2}{\beta_0^2} \left( \|M'\|_2 \left( \sqrt{D_K} + \sqrt{\eta} + \overbrace{\frac{D_K}{\sqrt{n}}}^{\leq D_K} \right) \beta_0 \right. \\ &\quad \left. + \|M\|_\infty(D_K + \eta) \right). \end{aligned}$$

But, according to Corollary 4.1,  $n\|\theta_K^0 - \hat{\theta}_K\|_\infty^2 = O_{\mathbb{P}}(1)$ . Moreover, by definition,

$$S_n(\gamma(\theta_K^0) - \gamma(\hat{\theta}_K)) = n(\gamma_n(\theta_K^0) - \gamma_n(\hat{\theta}_K)) + n \underbrace{\left( \mathbb{E}_{f^\varphi}[\gamma(\hat{\theta}_K)] - \mathbb{E}_{f^\varphi}[\gamma(\theta_K^0)] \right)}_{\geq 0}$$

Thus,  $n(\gamma_n(\theta_K^0) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1)$ . This holds for any  $K \in \mathcal{K}$  and then in particular for  $K_0$  and  $K$ . Besides,  $\gamma_n(\theta_K^0) = \gamma_n(\theta_{K_0}^0)$  since, by assumption,  $\gamma(\theta_K^0) = \gamma(\theta_{K_0}^0)$   $f^\varphi d\lambda$ -a.e. Hence  $n(\gamma_n(\hat{\theta}_{K_0}) - \gamma_n(\hat{\theta}_K)) = O_{\mathbb{P}}(1)$ .  $\square$

## References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Proceedings, 2nd Internat. Symp. on Information Theory*, pages 267–281.
- Arlot, S. (2007). *Resampling and model selection*. PhD thesis, Univ. Paris-Sud.
- Baudry, J., Raftery, A., Celeux, G., Lo, K., and Gottardo, R. (2010). Combining mixture components for clustering. *J. Comput. Graph. Statist.*, 19(2):332–353.
- Baudry, J.-P. (2009). *Model Selection for Clustering. Choosing the Number of Classes*. PhD thesis, Univ. Paris-Sud. <http://tel.archives-ouvertes.fr/tel-00461550/fr/>.
- Baudry, J.-P., Celeux, G., and Marin, J.-M. (2008). Selecting models focussing on the modeler’s purpose. In *COMPSTAT 2008: Proceedings in Computational Statistics*, pages 337–348, Heidelberg. Physica-Verlag.
- Baudry, J.-P., Maugis, C., and Michel, B. (2011). Slope heuristics: overview and implementation. *Statist. Comput.*, 22(2):455–470.
- Biernacki, C., Celeux, G., and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Trans. PAMI*, 22:719–725.
- Biernacki, C. and Govaert, G. (1997). Using the classification likelihood to choose the number of clusters. *Computing Science and Statistics*, 29:451–457.
- Birgé, L. and Massart, P. (2007). Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields*, 138(1-2):33–73.
- Celeux, G. and Govaert, G. (1995). Gaussian parsimonious clustering models. *Pattern Recognition*, 28(5):781 – 793.
- De Granville, C., Southerland, J., and Fagg, A. (2006). Learning grasp affordances through human demonstration. In *Proceedings of the International Conference on Development and Learning, electronically published*.
- Fraley, C. and Raftery, A. (2002). Model-based clustering, discriminant analysis, and density estimation. *J. Amer. Statist. Assoc.*, 97:611–631.

- Goutte, C., Hansen, L., Liptrot, M., and Rostrup, E. (2001). Feature-space clustering for fMRI meta-analysis. *Human Brain Mapping*, 13(3):165–183.
- Hamelryck, T., Kent, J. T., and Krogh, A. (2006). Sampling realistic protein conformations using local structural bias. *PLoS Comput. Biol.*, 2:e131.
- Hennig, C. (2010). Methods for merging Gaussian mixture components. *Adv. Data Anal. Classif.*, 4(1):3–34.
- Keribin (2000). Consistent estimation of the order of mixture models. *Sankhya A*, 62(1):49–66.
- Mariadassou, M., Robin, S., and Vacher, C. (2010). Uncovering latent structure in valued graphs: a variational approach. *Ann. Appl. Stat.*, 4(2):715–742.
- Massart, P. (2007). *Concentration Inequalities and Model Selection*. Lecture Notes in Math. Springer.
- Maugis, C. and Michel, B. (2011). A non asymptotic penalized criterion for Gaussian mixture model selection. *ESAIM Probab. Stat.*, 15:41–68.
- McLachlan, G. and Peel, D. (2000). *Finite Mixture Models*. New York : Wiley.
- McQuarrie, A. and Tsai, C. (1998). *Regression and time series model selection*. World Scientific.
- Nishii, R. (1988). Maximum likelihood principle and model selection when the true model is unspecified. *J. Multivariate Anal.*, 27:392–403.
- Pigeau, A. and Gelgon, M. (2005). Building and tracking hierarchical geographical & temporal partitions for image collection management on mobile devices. In *Proceedings 13th annual ACM internat. conf. on Multimedia*, pages 141–150. ACM New York, NY, USA.
- Redner, R. and Walker, H. (1984). Mixture densities, maximum likelihood and the EM algorithm. *SIAM Rev.*, 26(2):195 – 239.
- Rigail, G., Lebarbier, E., and Robin, S. (2012). Exact posterior distributions and model selection criteria for multiple change-point detection problems. *Statist. Comput.*, pages 1–13.
- Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, 6:461–464.

Steele, R. J. and Raftery, A. (2010). Performance of bayesian model selection criteria for gaussian mixture models. In *Frontiers of Statistical Decision Making and Bayesian Analysis*, pages 113–130. Springer.

Titterington, D., Smith, A., and Makov, U. (1985). *Statistical Analysis of Finite mixture Distributions*. New York : Wiley.

van der Vaart, A. (1998). *Asymptotic Statistics*. Cambridge University Press.

**Acknowledgements** I am deeply grateful to G. Celeux (INRIA Saclay Île-de-France and Université Paris-Sud), J.-M. Marin (Université Montpellier II) and P. Massart (Université Paris-Sud) for their essential help.

*Jean-Patrick Baudry*  
Université Pierre et Marie Curie - Paris VI  
Boîte 158, Tour 15-25, 2<sup>e</sup> étage  
4 place Jussieu, 75252 Paris Cedex 05  
France.  
Jean-Patrick.Baudry@upmc.fr  
<http://www.lsta.upmc.fr/Baudry>