

Clustering using EM and CEM, cluster number selection via the Von Mises-Fisher mixture models

Wafia Parr Bouberima, Mohamed Nadif, Yamina Khemal Bencheikh

LIPADE EA 2517, UFR Maths-Info, Paris Descartes University, France
e-mail:wboub@yahoo.fr

LIPADE EA 2517, UFR Maths-Info, Paris Descartes University, France
e-mail:mohamed.nadif@univ-paris5.fr

LMFN laboratory, Department of Mathematics, UFA University, Setif. Algeria
e-mail :bencheikh.00@yahoo.fr

Abstract

We consider the clustering problem of directional data and specifically the choice of the number of clusters. Setting this problem under the mixture approach, we perform a comparative study of different criteria. Monte Carlo simulations are performed taking into account the overlap degree of clusters and the size of data.

Keywords: *clustering, model selection, mixture models, information criteria, Von Mises-Fisher distribution.*

1 Introduction

Clustering is a key form of scientific research utilized within a variety of different scientific disciplines. Principally the classification method is used to produce g different clusters of wide distinctions. It should be noted that the optimum number of clusters g leading to the greatest separation is not known a priori and must be computed from the data, this is an heuristic problem in the classification topic; and this paper will be mainly concerned with this issue. In its main usage (Mainly), Clustering supports two approaches: a geometric one where the quality of the clustering depends on the chosen distance, and a probabilistic approach which is considered as a standard approach [14]. The latter covers the most widely used clustering methods. In this approach, data is presumed to come from a sampled mixture of g components which are modelled

by a distribution of probability. This approach can support several situations, depending on the parameters of the model, to obtain a best description of a heterogeneous population considering a selected model which is in itself another problem. The clustering problem can be resolved by mixture modelling and we can, for this, consider two approaches: the Maximum Likelihood (ML) and the Classification Maximum Likelihood (CML) approaches. The former is based on the maximization of the Likelihood, and the latter one is based on the maximization of the Classification (or complete data) Likelihood. These maximizations can be performed respectively by the *EM* algorithm and by the Classification *EM* (CEM) [9]. The model selection problem is to find the most appropriate and concise model to express given data.

Here, we merely examine some criteria from a practical point of view and in the context of the directional data utilizing a suitable distribution for mixture of directional data. The von Mises-Fisher distributions (*VMF*) are defined on the hypersphere $S^{(d-1)}$ [2] and appear adapted in this context. We consider Monte Carlo simulations and examine through numerical experiments on "real data" to see the validity of the proposed criteria for our main goal to estimate the number of clusters in a mixture.

This paper is organized as follows. Section 2 is devoted to describe the VMF mixture model. Section 3 begins with a review of the ML and CML approaches and a description of the EM and CEM algorithms. In Section 3, we review several criteria used in the determination of the number of clusters, and we evaluate these criteria. Finally, the last section summarizes the main points of this paper.

Notation Along this work, we assume that the data matrix \mathbf{x} is a contingency table, crossing, for example, n documents (rows) and d words (columns). In this case, each document is represented by $\mathbf{x}_i = (x_i^1, \dots, x_i^d) \in \mathbf{R}^d$, with $\|\mathbf{x}_i\| = 1$ ($\|\cdot\|$ denotes the standard L_2). Each value x_i^j corresponds to the frequency of a word j in a document i . A clustering of n documents provides a partition z into g classes.

2 Clustering via the Von Mises-Fisher mixture models

2.1 Finite mixture model

Finite mixture models underpin a variety of techniques in major areas of statistics including cluster analysis; see for instance [14]. With a mixture model-based approach clustering, it is assumed that the data to be clustered are generated by a mixture of underlying probability distributions in which each

component represents a different cluster. Given observations $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, let $\varphi_k(\mathbf{x}_i; \alpha_k)$ be the density of an observation \mathbf{x}_i from the k th component, where the α_k 's are the corresponding parameters and let g be the number of components in the mixture. The probability density function is

$$f(\mathbf{x}_i; \theta) = \sum_{k=1}^g \pi_k \varphi_k(\mathbf{x}_i; \alpha_k), \quad (1)$$

where π_k is the probability that an observation belongs to the k th component and θ is the vector of the unknown parameters $(\pi_1, \dots, \pi_g; \alpha_1, \dots, \alpha_g)$.

Setting the clustering problem of directional data under the mixture model approach, we assume that x is generated from a von Mises-Fisher mixture of g components. In this case

$$\varphi_k(\mathbf{x}_i; \alpha_k) = c_d(\xi_k) \exp \xi_k^T \mu_k \mathbf{x}_i$$

where $\alpha_k = (\mu_k, \xi_k)$; μ_k is the centre, ξ_k is the concentration of the k th cluster and $c_d(\xi) = \frac{\xi^{\frac{d}{2}-1}}{(2\pi)^{\frac{d}{2}} I_{\frac{d}{2}-1}(\xi)}$ with $I_{\frac{d}{2}}(\xi)$ is the modified Bessel function of the 1st type and of order $\frac{d}{2}$: $I_d(\xi) = \frac{1}{2\pi} \int_0^{2\pi} e^{\xi \cos \theta} \cos(d\theta) d\theta$.

Note that we can consider different parsimonious models by imposing constraints on π_k and ξ_k .

1. the proportions π_k of clusters and the concentrations ξ_k are supposed not equal, this model is noted $[\pi_k, \xi_k]$,
2. the concentrations ξ_k are supposed equal, this model is noted $[\pi_k, \xi]$,
3. the proportions π_k of clusters are supposed equal, this model is noted $[\pi, \xi_k]$,
4. the proportions π_k of clusters and the concentrations ξ_k are supposed equal, this model is noted $[\pi, \xi]$.

Next we focus on the general model $[\pi_k, \xi_k]$.

2.2 ML and CML approaches

The problem of clustering can be studied in the mixture model using the ML approach. This one, by maximizing the likelihood

$$L(\theta) = \prod_i \sum_{k=1}^g \pi_k \varphi_k(\mathbf{x}_i; \alpha_k),$$

has been by far the most commonly used approach to the fitting of mixture distribution and is appropriate to tackle this problem. It estimates the parameters

of the mixture, and the partition of I is derived from these parameters using the maximum a posteriori principle (MAP). Classical optimization techniques such as Newton-Raphson or gradient methods can be used but, in mixture context, the EM algorithm [10] has been successfully applied and is one of the most widely used procedures.

2.2.1 EM and CEM Algorithms

The EM algorithm is a method for maximizing the log-likelihood $L(\theta)$ iteratively, using the maximization of the conditional expectation of the complete-data log-likelihood given a previous current estimate $\theta^{(c)}$ and the observed data \mathbf{x} . In mixture model, we take the complete-data to be the vector (\mathbf{x}, \mathbf{z}) where the unobservable vector \mathbf{z} is the label data; the complete-data log-likelihood $L_c(\theta; \mathbf{x}, \mathbf{z})$ noted also $L_c(\mathbf{z}; \theta)$ is

$$L_c(\mathbf{z}; \theta) = \sum_{i,k} z_{ik} \log \pi_k \varphi_k(\mathbf{x}_i; \alpha_k) \quad (2)$$

and its conditional expectation can be written

$$\begin{aligned} Q(\theta, \theta^{(c)}) &= \sum_{i,k} s_{ik}^{(c)} \log(\pi_k \varphi_k(\mathbf{x}_i; \alpha_k)) \\ &= \sum_{i,k} s_{ik}^{(c)} \log(\pi_k c_d(\xi_k) e^{\xi_k^T \mu_k x_i}) \end{aligned}$$

where $s_{ik}^{(c)} = P(z_{ik} = 1 | \mathbf{x}, \theta^{(c)}) = \frac{\pi_k^{(c)} \varphi_k(\mathbf{x}_i; \alpha_k^{(c)})}{\sum_{k'=1}^g \pi_{k'}^{(c)} \varphi_{k'}(\mathbf{x}_i; \alpha_{k'}^{(c)})}$ denotes the conditional probability, given \mathbf{x} and $\theta^{(c)}$, that \mathbf{x}_i arises from the mixture component with density $\varphi_k(\mathbf{x}_i; \alpha_k)$. Each iteration of EM has two steps: an E-step and a M-step. The $(c+1)$ st E-step finds the conditional expectation of the complete-data log-likelihood. Note that in the mixture case this step reduces to the computation of the conditional density of the $s_{ik}^{(c)}$. The $(c+1)$ st M-step finds $\theta^{(c+1)}$ maximizing $Q(\theta, \theta^{(c)})$.

The characteristics of the EM algorithm are well documented. It leads in general to simple equations, has the nice property of increasing the log-likelihood at each iteration until stationarity, and in many circumstances, it derives sensible parameter estimates and consequently it is a popular tool to obtain maximum likelihood estimation. The EM algorithm can be viewed as a soft algorithm, and the partition can be derived from the parameters by using the MAP.

Note that a hard version CEM [9] can be performed by substituting $Q(\theta, \theta^{(c)})$ by $L_c(\theta)$. The main modifications concern therefore the conditional maximization of complete data log-likelihoods w.r. to \mathbf{z} given θ . in this context, we are not treating the estimation problem but; we are dealing with the problems of the selection of the number of components in a mixture.

2.2.2 The EM steps for a von Mises-Fisher mixture model

The EM algorithm, as explained previously is used to compute the maximum likelihood (ML) estimates of all the parameters through the iterated application of the estimation and maximization of $Q(\theta, \theta^{(c)})$. Starting from an initial situation $\theta^{(0)}$, an iteration ($c > 0$) is defined as follows: After the Estimation step, where the current posterior $s_{ik}^{(c)}$ is computed. The Maximization step compute the ML estimates $\theta^{(c)} = (\mu_k^{(c)}, \pi_k^{(c)}, \xi_k^{(c)})$, as following:

- $\pi_k^{(c)} = \frac{\sum_{i=1}^n s_{ik}^{(c)}}{n}$
- $\mu_k^{(c)} = \frac{\sum_{i=1}^n s_{ik}^{(c)} x_i}{\left\| \sum_{i=1}^n s_{ik}^{(c)} x_i \right\|}$
- $\xi_k^{(c)} = A_d^{-1} \left(\frac{\left\| \sum_{i=1}^n s_{ik}^{(c)} x_i \right\|}{\pi_k^{(c)} \times n} \right)$
with $A_d(\xi) = \frac{I_{\frac{d}{2}}(\xi)}{I_{\frac{d}{2}-1}(\xi)}$

Then, a partition $z = (z_1, \dots, z_k)$ of the data can be directly derived from the ML estimates of the mixture parameters by assigning each x_i to the component which provided the greatest posterior probability.

3 Number of components selection

The determination of the numbers of components g and m can be viewed as a model selection problem which can be solved by different criteria: information model selection criteria, methods based on confidence interval, and empirical criteria [8]. In this paper we will explore some information criteria, which are the most important and popular techniques. They consist to penalize the models with additional parameters. These criteria split into two terms: one for the model fitting fit, which is data likelihood or complete data likelihood, and one for the model complexity. There is also different methods of selection of the number of components in a mixture, in which the iterations of the algorithm delete the empty or less condensed components or more like the hierarchical classification, merging many clusters to one in an agglomerative way to a fixed level of dissimilarity [11]. In the mixture modelling the stochastic version of EM is particularly appealing to estimate the components of a mixture, in its stochastic iteration, the SEM exclude the components in which the cardinal is lower than a fixed initial number; iteratively this algorithm estimates the number of clusters.

3.1 SEM algorithm

The SEM algorithm is a stochastic version of EM incorporating between E and M steps a restoration of the unknown component labels $z_i, i = 1, \dots, n$, by drawing them at random from their current conditional distribution, starting from an initial parameter, consisting the three steps:

- Expectation E : compute the conditional probabilities t_{ik} for the current parameters of the mixture
- Stochastic S : assign each point at random to one of the mixture components according to the multinomial distribution with parameters t_{ik} .
- Maximisation M : update the ML estimates of the parameters of the mixture using the partition result of the step S.

3.2 Information criteria

Let L be the \log -likelihood of observed data, L_c be the complete data \log -likelihood with the parameter $\hat{\theta}$ obtained by the EM algorithm, v be the number of free parameters in the mixture model and $E = \sum_{i,k} s_{ik} \log(s_{ik})$ the entropy criterion. The terms L, L_c, v and E depend on g . In the following, we shall focus on twelve criteria.

- $Bic(g) = -2L(g) + v \ln n$, proposed by Schwarz [18] and Rissanen [17]
- $Aic(g) = -2L(g) + 2v$, proposed by Akaike [1]
- $Aic3(g) = -2L(g) + 3v$, proposed by Bozdogan [7]
- $Aic4(g) = -2L(g) + 4v$, proposed by Bozdogan [7]
- $Aicc(g) = Aic(g) + \frac{2v(v+1)}{n-v-1}$, proposed by Hurvich and Tsai [13]
- $Aicu(g) = Aicc(g) + n \ln n / (n - v - 1)$, proposed by McQuarrie, Schwarz and Tsai [15]
- $CAic(g) = -2L(g) + v(1 + \ln n)$, proposed by Bozdogan [6]
- $Cle(g) = -2L(g) + 2E(g)$, proposed by Biernacki [4]
- $IclBic(g) = Bic(g) + 2E(g)$, proposed by Biernacki, Celeux and Govaert [5]
- $Ll(g) = -L(g) + \frac{v}{2} \sum_k \ln \frac{n\pi_k}{2} + \frac{g}{2} \ln \frac{n}{12} + \frac{g(v+1)}{2}$, proposed by Figueiredo and Jain (2002) [12]

- $Icl(g) = -2L_c(g) + v \ln n$, proposed by Biernacki, Celeux and Govaert [5]
- $Awe(g) = -2L_c(g) + 2v(\frac{3}{2} + \ln n)$, proposed by Banfield and Raftery [3]

3.3 Experimental conditions

In our experiments, we perform a study according to the degree of overlap of clusters and the size of data.

1. The concept of cluster separation is difficult to visualize easily for our model, but the degree of overlap can be measured by the true error rate approximated by comparing the partitions simulated with those we obtained by applying a classification step. From our numerical experiments, we present only 3 situations corresponding to 3 levels of overlap degrees: clusters well separated ($\approx 5\%$), moderately separated ($\approx 15\%$) and poorly separated ($\approx 23\%$).

To explain the different degrees of overlap, the following plots represents six samples simulated regarding the three degrees, the first group in 2 dimensions and the second in 3 dimensions, with $n = 120$, $g = 3$ and different parameters (μ, k, p) .

2. We selected several sizes of data 600×3 , 1800×3 , 6000×3 , 6000×50 and 6000×50 data arising from 3- components mixture model corresponding to the three degrees of overlap.

Before any application lets observe; the model complexity for most of the criteria exposed previously depend directly on the number of the unknown free parameters in the clustering model. Such number considering a model of von Mises-Fisher with the unknown parameters (μ_k, π_k, ξ_k) , is: $v = g(d+2) - 1$, so if $g = 2, \dots, 5$, a quick calculus give us an idea about this quantity and about the $\log(n)$, in the following table:

Table 1: number of free parameters for g clusters.

g/d	3	50	100	n	$\ln(n)$
2	9	103	203	600	6.3969
3	14	155	305	1800	7.4955
4	19	207	407	3000	8.0063
5	24	259	509	6000	8.6995

Furthermore, let us compute the penalty terms for all the criteria, where the term is independent of the iterations of the EM algorithm (table2). According to the increasing values of g all the criteria values increase uniformly.

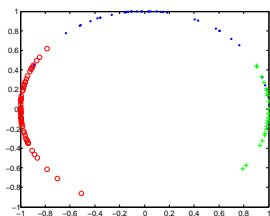


Figure 1: Sample1: 5%
degrees of overlap, $d = 2$

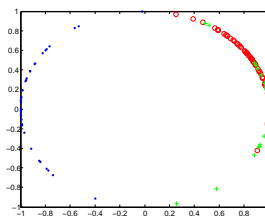


Figure 2: Sample2: 15%
degrees of overlap, $d = 2$

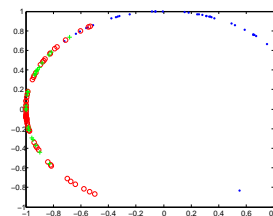


Figure 3: Sample3: 25%
degrees of overlap, $d = 2$

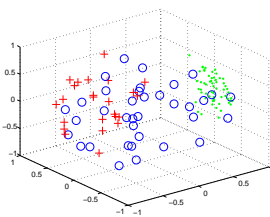


Figure 4: Sample4: 5%
degrees of overlap, $d = 3$

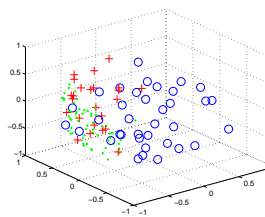


Figure 5: Sample5: 15%
degrees of overlap, $d = 3$

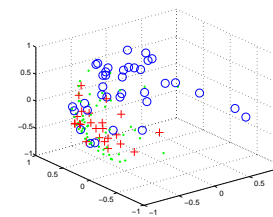


Figure 6: Sample6: 25%
degrees of overlap, $d = 3$

Aic , $Aic3$, $Aic4$ and $Aicc$ are independent of the number of the lines n , for these criteria the penalty term increases modestly when the number of the columns d increases. The Bic , $Icl - Bic$, Icl and $Caic$ penalty term is a composition of both numbers n and d , and it increases according to both of them. By the same way the Awe 's penalty increases, but quicker. It's clear that the quality of any of the above criteria is affected directly by this term.

Now, to evaluate the EM algorithm and the previous criteria, many applications on simulated data was realized. For each θ leading the degree of overlap, we generated 20 samples. For each sample and to avoid local optima in the generated estimation process, the $EM(g)$ algorithm ($g = 2, \dots, 5$) regarding the general model $[\pi_k, \xi_k]$, is repeated 20 times starting from the best partition obtained by the spherical k means [2] which is a CEM applied with the model $[\pi, \xi]$. From the best solution,

1. we compute the percent of documents misclassified by comparing the true partition and the obtained partition with the same number of clusters,
2. we compute all criteria previously cited in function of different values of g ,
3. we count the number of times on 20 that each criterion detects the original number of clusters *fit*, overestimates it *over-fit* or underestimates it *under-fit*. In table 1 are reported all results obtained by all criteria.

From these experiments, the main points arising are the following.

- The EM algorithm gives good results by comparing the true partition and the obtained one by EM(3).
- When the clusters are well or moderately separated $Aic3$, $Aic4$, $Aicu$ and Bic are the more efficient for the studied sizes.
- When the clusters are poorly separated, the quality of these criteria increases with the size of the data n and when $n \gg d$.
- Moreover note that $Aic3$ and $Aicu$ outperform Bic when the number of columns increases and remain interesting in the most situations. In fact, Bic seems very sensitive to the dimension, it underestimates the number of clusters.

In these first experiments, we can consider that $Aic3$ and $Aicu$ are the best criteria. Note that $Aic3$ is also interesting for the Bernoulli mixture model for the binary data [16]. Nevertheless, we have noted that their performances decrease when we are in the high dimension. Then we illustrate the behavior of all criteria by using a well known set of data known as Classic3 as a real data

application. This is a set of documents from three well separated sources. Classic3 contains 3893 documents (vectors) with a total of 4303 features (words). The data matrix consists of 1400 Cranfield documents from aeronautical system papers, 1033 from Medline documents obtained from medical journals, and 1460 Cisi documents obtained from information retrieval papers. Each vector was normalized in order to be used as a unit vector. In order to select a number of clusters in $g = 2, \dots, 5$, we have computed the same criteria as previously, we applied the $EM(g)$ algorithm regarding the general model $[\pi_k, \xi_k]$ and we obtained the following results:

- Bic, Caic, Icl-Bic, Icl overestimate the number of clusters and give 4 clusters.
- Aic, Aic3, Aic4, Aicc, Clc overestimate the number of clusters and give 5 clusters.
- Aicu, Ll, Awe underestimate the number of clusters and give 2 clusters.

4 Conclusion

Setting the clustering of directional data in the mixture approach context, we have performed some experiments in order to evaluate the EM algorithm and to assess the number of clusters. Different information criteria have been tested on different sizes of data according different degree of overlap. We have observed that some of them such as Aic3, Aic, Aicu and Bic are interesting. Moreover we have noted that their performance increases on the size of data and Aic3 and Aicu appear as the best.

5 Open Problem

In future work, it will be interesting :

- 1) to take into account the high dimension in these criteria.
- 2) to tackle simultaneously the problem of assessing of the number of clusters combined to the choice of the parsimonious models $[\pi_k, \xi]$, $[\pi_k, \xi]$ and $[\pi, \xi]$.

References

- [1] Akaike, H., "Information theory and an extension of maximum likelihood principle," *Second International Symposium on Information Theory*, Akademia Kiado, 267-281, 1973.

- [2] Banerjee, A., Dhillon, I. S., Ghosh J., Sra S., "Clustering on the Unit Hypersphere Using Von Mises-Fisher Distributions," *Journal of Machine Learning Research*, 6:1345-1382, 2005.
- [3] J. D. Banfield and A. E. Raftery., "Model-based gaussian and non-gaussian clustering," *Biometrics*, 49:803821, 1993.
- [4] Biernacki, C., "Choix de modèles en Classification," *PhD Thesis*, Compiègne University of Technology, 1997.
- [5] Biernacki, C., Celeux, G. and Govaert, G., "Assessing a Mixture model for Clustering with the integrated Completed Likelihood," *IEEE Transactions on Pattern analysis and Machine Intelligence 22 (7)*, pp. 719-725, 2000.
- [6] Bozdogan, H., "Model Selection and Akaike's Information Criterion (AIC): The General Theory and its Analytical Extensions", *Psycometrika 52 (3)*, pp. 345-370, 1987.
- [7] Bozdogan, H., "Mixture-Model Cluster Analysis Using Model Selection Criteria and a New Information Measure of Complexity" *Proceedings of the first US/Japan conference on the Frontiers of Statistical Modeling: An Informational Approach*, 1 ed. 3 vols. Vol. 1., Dordrecht, Kluwer Academic Publishers, 1994.
- [8] Bubna K., Stewart, C.V., "Model Selection Techniques and Merging Rules for Range Data Segmentation Algorithms," *Computer Vision and Image Understanding*, 80: 215-245, 2000.
- [9] Celeux, G., Govaert, G., "A classification EM Algorithm for clustering and two stochastic versions," *Computational statistics & Data analysis*, 14:315-332, 1992.
- [10] Dempster, A.P., Laird, N.M., Rubin, D., "Maximum Likelihood from Incomplete Data via the EM Algorithm," *J of the Royal Stat Soc*, B 39: 1-38, 1977
- [11] Dumais, S. T., Chen, H., "Hierarchical classification of web content." *Proceedings of the 23rd ACM International Conference on Research and Development in Information Retrieval*, 256263. ACM Press, 2000.
- [12] Figueiredo, M.A.T., and Jain, A.K., "Unsupervised Learning of Finite Mixture Models," *IEEE Transactions on pattern analysis and Machine Intelligence 24 (3)*, pp. 1-16, 2002.

- [13] *Hurvich, C.M., and Tsai, C.-L., "Regression and Time Series Model Selection in Small Samples," Biometrika 76 (2), pp. 297-307, 1989.*
- [14] *McLachlan, G.J., Peel, D., Finite mixture models, Wiley, New York, 2000.*
- [15] *McQuarrie, A., Shumway, R. and Tsai, C.-L., "The model selection criterion AIC_u," Statistics & Probability Letters 34, pp. 285-292, 1997.*
- [16] *Nadif, M., Govaert, G., "Clustering for binary data and mixture models: Choice of the model," Applied Stochastic Models and Data Analysis, 13: 269-278, 1998.*
- [17] *Rissanen, J., "Modelling by shortest data description," Automatica, 14:465-471, 1978.*
- [18] *Schwarz, G., "Estimating the Dimension of a Model," The Annals of Statistics 6 (2), pp. 461-464, 1978.*

Table 2: Values of the penalty term for some information criteria.

<i>Size</i>	<i>Clusters</i>	<i>Bic</i> <i>Icl - Bic</i> <i>Icl</i>	<i>Aic</i>	<i>Aic3</i>	<i>Aic4</i>	<i>Aicc</i>	<i>Aicu</i>	<i>Caic</i>	<i>Awe</i>
600 × 3	2	57.572	18	27	36	18.033	28.118	66.572	142.144
	3	89.557	28	42	56	28.051	43.241	103.557	221.114
	4	121.541	38	57	76	38.068	58.409	140.541	300.083
	5	153.526	48	72	96	48.086	73.622	177.526	379.052
1800 × 3	2	67.459	18	27	36	18.011	28.039	76.459	161.919
	3	104.937	28	42	56	28.016	43.079	118.937	251.875
	4	142.415	38	57	76	38.022	58.134	161.415	341.830
	5	179.893	48	72	96	48.028	73.203	203.893	431.786
6000 × 3	2	78.295	18	27	36	18.003	28.011	87.295	183.591
	3	121.793	28	42	56	28.005	43.023	135.793	285.586
	4	165.290	38	57	76	38.006	58.040	184.290	387.581
	5	208.788	48	72	96	48.008	73.060	232.788	489.576
3000 × 50	2	824.655	206	309	412	206.071	311.917	927.655	1958.311
	3	1240.986	310	465	620	310.109	470.312	1395.986	2946.973
	4	1657.318	414	621	828	414.148	629.711	1864.318	3935.636
	5	2073.649	518	777	1036	518.189	790.152	2332.649	4924.298
6000 × 50	2	896.050	206	309	412	206.035	310.947	999.050	2101.100
	3	1348.424	310	465	620	310.053	468.117	1503.424	3161.849
	4	1800.799	414	621	828	414.071	625.762	2007.799	4222.599
	5	2253.174	518	777	1036	518.090	783.892	2512.174	5283.348
6000 × 100	2	1766.001	406	609	812	406.070	613.619	1969.001	4141.002
	3	2653.351	610	915	1220	610.107	924.186	2958.351	6221.703
	4	3540.702	814	1221	1628	814.145	1236.680	3947.702	8302.405
	5	4428.053	1018	1527	2036	1018.185	1551.173	4937.053	10383.106

Table 3: Evaluation of EM and all information criteria for the model $[\pi_k, \xi_k]$. For each criterion, the numbers of times on 20 that a criterion detects or not the good number of clusters (a).

<i>size</i>	<i>degree</i>	<i>EM(3)</i>	<i>fit</i>	<i>Bic</i>	<i>Aic</i>	<i>Aic3</i>	<i>Aic4</i>	<i>Aicc</i>	<i>Aicu</i>	<i>CAic</i>	<i>Clc</i>	<i>Icl - Bic</i>	<i>Ll</i>	<i>Icl</i>	<i>Awe</i>	
600 × 3	4.88%	5.17%	under-fit	0	0	0	0	0	0	0	0	0	0	0	0	
			fit	20	15	19	20	15	19	20	15	20	20	20	20	
			over-fit	0	5	1	0	5	1	0	5	0	5	0	0	0
1800 × 3	5.16%	4.83%	under-fit	0	0	0	0	0	0	0	0	0	0	0	0	
			fit	20	18	19	20	18	19	20	18	20	20	20	20	
			over-fit	0	2	1	0	2	1	0	2	0	2	0	0	0
3000 × 50	4.74%	6.85%	under-fit	0	0	0	0	0	0	0	0	0	1	7	0	4
			fit	20	1	20	20	1	20	20	6	19	13	20	16	
			over-fit	0	19	0	0	19	0	0	14	0	0	0	0	0
600 × 3	14.63%	16.33%	under-fit	0	0	0	0	0	0	0	9	16	7	7	16	
			fit	20	17	20	20	17	20	20	9	4	13	13	4	
			over-fit	0	3	0	0	3	0	0	2	0	0	0	0	
1800 × 3	15.10%	15.83%	under-fit	0	0	0	0	0	0	0	14	18	0	2	7	
			fit	20	19	20	20	19	20	20	6	2	20	18	13	
			over-fit	0	1	0	0	1	0	0	0	0	0	0	0	
3000 × 50	13.68%	14.10%	under-fit	0	0	0	0	0	0	0	0	3	0	0	20	
			fit	20	10	20	20	10	20	20	18	17	20	20	0	
			over-fit	0	10	0	0	10	0	0	2	0	0	0	0	
6000 × 100	15.11%	18.35%	under-fit	0	0	0	0	0	0	4	20	20	20	20	20	
			fit	20	0	20	20	0	20	16	0	0	0	0	0	
			over-fit	0	20	0	0	20	0	0	0	0	0	0	0	

Table 4: Evaluation of EM and all information criteria for the model $[\pi_k, \xi_k]$.
 For each criterion, the numbers of times on 20 that a criterion detects or not
 the good number of clusters (b).

<i>size</i>	<i>degree</i>	<i>EM(3)</i>	<i>fit</i>	<i>Bic</i>	<i>Aic</i>	<i>Aic3</i>	<i>Aic4</i>	<i>Aicc</i>	<i>Aicu</i>	<i>CAic</i>	<i>Clc</i>	<i>Icl - Bic</i>	<i>Ll</i>	<i>Icl</i>	<i>Awe</i>
600 × 3	24.96%	29.17%	under-fit	20	15	17	20	15	18	20	20	20	20	20	20
			fit	0	3	3	0	3	2	0	0	0	0	0	0
			over-fit	0	2	0	0	2	0	0	0	0	0	0	0
1800 × 3	25.19%	35.94%	under-fit	20	12	17	19	12	17	20	20	20	20	20	20
			fit	0	8	3	1	8	3	0	0	0	0	0	
			over-fit	0	0	0	0	0	0	0	0	0	0	0	
6000 × 3	27.49%	30.95%	under-fit	0	0	0	0	0	0	0	20	20	8	20	20
			fit	20	20	20	20	20	20	0	0	0	12	0	0
			over-fit	0	0	0	0	0	0	0	0	0	0	0	
3000 × 50	24.75%	32.26%	under-fit	18	0	0	0	0	0	20	20	20	20	20	20
			fit	2	8	20	20	8	20	0	0	0	0	0	0
			over-fit	0	12	0	0	12	0	0	0	0	0	0	
6000 × 50	25.61%	30.17%	under-fit	20	0	1	16	0	1	20	20	20	20	20	20
			fit	0	11	19	4	11	19	0	0	0	0	0	0
			over-fit	0	9	0	0	9	0	0	0	0	0	0	
6000 × 100	26.74%	42.91%	under-fit	20	8	20	20	8	20	20	19	20	20	20	20
			fit	0	3	0	0	3	0	0	1	0	0	0	
			over-fit	0	9	0	0	9	0	0	0	0	0	0	