# Book of Short Papers
# SIS 2020

SIS2020 *Pisa*

SIS
Società
Italiana di
Statistica

Editors: Alessio Pollice, Nicola Salvati and Francesco Schirripa Spagnolo

# Contents

## Specialized sessions

# Solicited Sessions

VIII

# Contributed papers and Posters

# Finance, business and official statistics ......................................886

## Machine Learning and Data Science........................................................1023

## Models and methods – Sampling ................................................1271

## Models and methods - Theoretical Issues in Statistical Inference ..........1314

## Models and methods - Time Series and Longitudinal Data.....................1350

## Population and society ...................................................................1411

XVI

# Preface

The COVID-19 pandemic is putting our society under incredible health, emotional, and economic stress. Facing its harmful effects and their uncertainty, the Executive Board of the Italian Statistical Society (SIS) and the Local Organizing Committee, to ensure the highest level of safety for members and delegates, deliberated to cancel the 50th Meeting of the Italian Statistical Society originally planned to be held in Pisa in June 2020 and to postpone the conference to June 2021. The Executive Board and the Local Organizing Committee continue to monitor closely the pandemic evolving situation, and keep the members of SIS and the researchers informed about the potential new dates for the next meeting. To give value to the work of those who prepared their presentation for the conference, the Program Committee decided to publish the volume *Book of short papers - SIS 2020* despite the conference cancellation.

The conference program included 4 plenary sessions, 16 specialized sessions, 24 solicited sessions, 32 contributed sessions and the poster exhibition. Plenary sessions concerned with robust statistics, human longevity, statistical models for climate changes and small area estimation for educational poverty. The meeting had to host also 2 round tables on data privacy and innovation in statistics. Activities focused on topics of interest for a wider audience included two round tables on Teaching Statistics and on the SIS journal Statistical Methods & Applications, and the Stats Under the Stars (SUS6) competition for young statisticians. The SUS6 event attracted many sponsors from statistical, financial and editorial firms as well as numerous students. The conference committee had registered 345 accepted submissions, including 143 to be presented in invited plenary, specialized and solicited sessions, and 202 spontaneously submitted for oral and poster sessions.

This book includes most of the scientific contributions that had to be presented at the 50th Meeting of the Italian Statistical Society. It is organized into 49 chapters corresponding to 15 specialized, 23 solicited sessions, and to 11 general topics for contributed papers and posters. All 268 contributions provide a wide overview of the state-of-the-art of the subjects, from methodological and theoretical contributions, to applied works and case studies. The result is a very lively picture of the Italian statistical community with its international connections.

We would like to thank all contributors for having submitted their work to the conference, the members of the Program Committee and the extra reviewers for their efforts in this difficult period. Although the Conference did not take place, the organization went on until cancellation was decided for safety reasons. It would have been impossible without the joint effort of Università di Pisa, Scuola Superiore Sant'Anna and National Research Council of Pisa. Members these three institutions took part actively in the Local Organizing Committee. Finally we wish to express our gratitude to the publisher Pearson Italia for all the support received.

This book is our contribution to encourage the scientific community and the network of the Italian Statistical Society to go on and transform this difficult period into an opportunity of scientific debate for better statistics in a better world.

Alessio Pollice
Università degli Studi di Bari Aldo Moro
Chair of the Program Committee

Nicola Salvati
Università di Pisa
Chair of the Local Organizing Committee

Francesco Schirripa Spagnolo
Università di Pisa

**Program Committee**: Alessio Pollice (Chair), Serena Arima, Marilena Barbieri, Alessandra Brazzale, Eugenio Brentari, Alessia Caponera, Antonio Lepore, Antonella Plaia, Tommaso Proietti, MarcoRiani, Nicola Salvati, Pasquale Sarnacchiaro, Mauro Scanu, Manuela Stranges, Valentina Tocchioni, Simone Vantini, Massimo Ventrucci, Paola Vicard, Donatella Vicari.

**Local Organizing Committee**: Nicola Salvati (Chair), Gaia Bertarelli, Bruno Cheli, Alessandra Coli, Paolo Frumento, Fosca Giannotti, Caterina Giusti, Piero Manfredi, Stefano Marchetti, Lucio Masserini, Vincenzo Mauro, Barbara Pacini, Dino Pedreschi, Francesco Schirripa Spagnolo, Chiara Seghieri.

**Organizers of Specialized and Solicited Sessions**: Giada Adelfio, Bruno Arpino, Emanuele Aliverti, Nicoletta Balbo, Mara Bernardi, Silvia Bozza, Pierpaolo Brutti. Annalisa Busetta, Michela Cameletti, Carlo Cavicchia, Fabrizio Durante, Leonardo Egidi, Pietro D. Falorsi, Francesco Finazzi, Livio Finos, Stefania Galimberti, Michele Gallo, Caterina Giusti, Francesca Greselin, Alessandra Guglielmi, Francesca Ieva, Tiziana Laureti, Achille Lemmi, Brunero Liseo, Fabio Massimo Lo Verde, Daria Mendola, Roberta Pappadà, Lea Petrella, Alessandra Petrucci, Alessia Pini, Sabrina Prati, Maria Giovanna Ranalli, Davide Risso, Fabrizio Ruggeri, Silvana Salvini, Monica Scannapieco, Francesco Stingo, Luca Tardella, Grazia Vicario, Susanna Zaccarin, Maroussa Zagoraiou.

# Specialized sessions

# Aggregating Gaussian mixture components
## *Come aggregare le componenti gaussiane di un miscuglio*

Roberto Rocci[1]

**Abstract** The finite mixture of Gaussians is a well-known model frequently used to classify a sample of observations. It considers the sample as drawn from a heterogeneous population where each subpopulation, cluster, is Gaussian and corresponds to one component of the mixture. Whenever such assumption is false, the model may use two or more Gaussians to describe a single cluster. In this case, the researcher has the problem of how to identify the clusters starting from the estimated components. This work proposes to solve this problem by aggregating the components in clusters by optimizing an appropriate criterion based on their posterior probabilities.

**Abstract** *I modelli miscuglio di gaussiane sono spesso utilizzati nell'analisi dei gruppi. L'idea è quella di considerare la popolazione che ha generato il campione come formata da sottopopolazioni, gruppi, ognuna ben descritta da una componente del miscuglio. Quando questa assunzione risulta falsa, il modello tende ad utilizzare due o più componenti per rappresentare un unico gruppo, creando così il problema di come identificare i gruppi a partire dalle componenti stimate. In questo lavoro proponiamo di risolvere il problema aggregando le componenti in modo da ottimizzare un criterio basato sulle loro probabilità a posteriori.*

**Key words:** Unsupervised Classification, Finite mixtures of Gaussians, Within and Between deviances.

## 1 Introduction

In cluster analysis, or unsupervised classification, quite frequently observations are classified by using a finite mixture of Gaussians (see for example Hennig et al., 2015). The idea is to consider the population as heterogeneous, i.e. formed by sub-populations, clusters, which are well represented by the mixture components.

---

[1]Department of Statistical Science, Sapienza University of Rome; email: roberto.rocci@uniroma1.it

Such approach works well in practice unless one or more subpopulations have a distribution different from, or not well approximated by, a single component of the mixture. In this case, the model may use more than one component, i.e. a sub-mixture of two or more Gaussians, to describe a single cluster destroying the one to one correspondence between clusters and components. This problem is well recognized in practice and several solutions have been proposed.

The first idea is to assume for each component a functional form that is more flexible than the Gaussian (see McNicholas, 2016, for an excellent review about this approach). This solves the problem in many cases in practice. However, it cannot be considered as the definitive solution because even in this case we cannot exclude that more than one component could be necessary to represent a cluster. This derives from the identifiability of the finite mixture model. For example, a cluster that is a finite mixture of two components cannot be represented by only one and vice versa.

A different idea comes up by observing that, if the mixture fits well the data, then the information about the clustering structure is contained in the estimated model and it can be recovered by aggregating, in an opportune way, the components. This allow us to represent a very wide variety of possible distributions for the clusters and to relax the, usually made and sometimes restrictive, assumption of same functional form for the distribution of each cluster. Technically, the model would become a finite mixture of finite mixtures of Gaussians, i.e. a finite mixture where each component is a finite mixture of Gaussians. Unfortunately, such a model is not identified because different aggregations of the components give the same population distribution. The estimation is then possible only by using some constraints on model parameters making the model identified (see for example Di Zio et al. 2007), or by introducing a criterion determining the aggregation. On the latter approach, there are several proposals in the literature where the Gaussian components are hierarchically aggregated into clusters on the basis of a measure of proximity (see Hennig 2010, Comas-Cufi et al. 2017 and references there in for some examples). However, such methods are optimal only locally. They establish what is the best way to merge two components not what is the best way to aggregate, say $G$, components into, say $K$, clusters. It is not specified how the "internal cohesion" and "external isolation" (Cormack, 1971) are measured, related and optimized (e.g. total deviance = within deviance + between deviance in $K$-means). To achieve this goal a partitioning method should be adopted but, as far as we know, only Li (2005) considered a partitioning approach based on the application of the $K$-means on the mean components. This proposal makes clear the aforementioned aspects but it is based on a measure of dissimilarity between components depending only on their locations. Our purpose is to go beyond the limits of this proposal.

In our paper we are going to propose a new method to aggregate the Gaussian components of a finite mixture model making clear how the identified partition optimize the "internal cohesion" and "external isolation" of the clustering. The plan of the paper is the following. Our proposal, based on the use of the Kullback Leibler divergence to measure the dissimilarity among components, will be presented in section 2. In section 3, some insights on how to extend the technique are presented with particular reference to other dissimilarity measures and its hierarchical version.

## 2 Partitioning Gaussian components by the Kullback-Liebler divergence: the KL-components method

The finite mixture of Gaussians (McLachlan & Peel, 2000) is based on the assumption that the probability density of a multivariate observation is of the form

$$f(\mathbf{x}_i;\Theta) = \sum_{g=1}^{G} \pi_g f_g(\mathbf{x}_i;\vartheta_g).$$ (1)

where $\mathbf{x}_i = [x_{i1}, x_{i2},\ldots,x_{iJ}]'$ is a random vector of $J$ variables sampled from a population parametrized by $\Theta = \{\vartheta_1,\vartheta_2,\ldots,\vartheta_G, p_1,p_2,\ldots,p_G\}$, which consists of $G$ groups, or subpopulations, in proportions $\pi_1, \pi_2,\ldots, \pi_G$, where $\pi_g$ is the prior probability to sample one observation from group $g$. The density $f_g(\mathbf{x}_i;\vartheta_g)$ of $\mathbf{x}_i$ in the $g^{th}$ group is multivariate normal (Gaussian). Let $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n$ be a sample of $n$ independent and identically distributed observations, we can use the above model to classify the observations into $G$ classes. First, we compute the posterior probabilities

$$P(g \mid \mathbf{x}_i) = \pi_{g|i} = \pi_g f_g(\mathbf{x}_i;\vartheta_g) / \sum_{h=1}^{G} \pi_h f_h(\mathbf{x}_i;\vartheta_h),$$ (2)

then, we use the MAP rule (Maximum A Posterior probability) to assign the observations to the Gaussian components. Usually, the parameter $\Theta$ and the number of components $G$ are unknown and estimated from the data.

Our method, named the KL-components technique, originates from the observation that two components, say $g$ and $h$, are equal, with respect to the data, if and only if their posterior probabilities are proportional. In formulas

$$g = h \Leftrightarrow \boldsymbol{\pi}_g = P(g \mid \mathbf{x}_i) / \pi_g = P(h \mid \mathbf{x}_i) / \pi_h = \boldsymbol{\pi}_h, \quad g,h = 1,\ldots,G \text{ and } i = 1,\ldots,n. \quad (3)$$

It seems quite natural to measure the dissimilarity between components as a function of the diversity between the normalized posteriors, say profiles. In particular, we investigate the use of the Kullback-Leibler (KL) divergence

$$\mathrm{KL}(\boldsymbol{\pi}_g, \boldsymbol{\pi}_h) = \frac{1}{n}\sum_{i=1}^{n} \frac{\pi_{g|i}}{\pi_g} \log\left(\frac{\pi_{g|i}}{\pi_g} \frac{\pi_h}{\pi_{h|i}}\right) = \frac{1}{n}\sum_{i=1}^{n} \frac{f_g(\mathbf{x}_i)}{f(\mathbf{x}_i)} \log\left(\frac{f_g(\mathbf{x}_i)}{f_h(\mathbf{x}_i)}\right),$$ (4)

which can be considered an estimate of

$$\mathrm{E}\left(\frac{f_g(\mathbf{x})}{f(\mathbf{x})} \log\left(\frac{f_g(\mathbf{x})}{f_h(\mathbf{x})}\right)\right) = \int f_g(\mathbf{x}) \log\left(\frac{f_g(\mathbf{x})}{f_h(\mathbf{x})}\right) d\mathbf{x},$$ (5)

153

i.e. the KL divergence between the two components. It is well known, and evident from formula (5), that the KL divergence is not symmetric. However, this is not a problem for us because we use it to define a sort of deviance rather than to measure the dissimilarity between two components. In particular, we define the Total deviance as the weighted sum of the KL divergences among the profiles of the mixture components and their barycenter. By noting that the latter quantity is

$$\sum_{g=1}^{G} \pi_g \frac{\pi_{g|i}}{\pi_g} = 1, \quad i = 1, 2, ..., n,$$ (6)

the Total deviance results to be

$$D_T = \sum_{g=1}^{G} \pi_g \mathrm{KL}(\boldsymbol{\pi}_g, \mathbf{1}) = \frac{1}{n} \sum_{g=1}^{G} \sum_{i=1}^{n} \pi_{g|i} \log\left(\frac{\pi_{g|i}}{\pi_g}\right).$$ (7)

Given a partition of the $G$ components into $K$ clusters according to the binary row stochastic membership matrix $\mathbf{U} = [u_{gk}]$, where $u_{gk}$ is equal to 1 if component $g$ belong to cluster $k$ and 0 otherwise, we note that the posteriors probabilities of cluster $k$ are

$$p_{k|i} = \sum_{g=1}^{G} u_{gk} \pi_{g|i},$$ (8)

the priors are

$$p_k = \sum_{g=1}^{G} u_{gk} \pi_g = \sum_{g=1}^{G} u_{gk} \frac{1}{n} \sum_{i=1}^{n} \pi_{g|i} = \frac{1}{n} \sum_{i=1}^{n} \sum_{g=1}^{G} u_{gk} \pi_{g|i} = \frac{1}{n} \sum_{i=1}^{n} p_{k|i},$$ (9)

and the barycenter is

$$\frac{1}{\sum_{g=1}^{G} u_{gk} \pi_g} \sum_{g=1}^{G} u_{gk} \pi_g \frac{\pi_{g|i}}{\pi_g} = \frac{1}{p_k} \sum_{g=1}^{G} u_{gk} \pi_{g|i} = \frac{p_{k|i}}{p_k}, \quad i = 1, 2, ..., n.$$ (10)

According to (7), the Within deviance is defined as

$$D_W(\mathbf{U}) = \sum_{k,g} u_{gk} \pi_g \mathrm{KL}(\boldsymbol{\pi}_g, \mathbf{p}_k) = \sum_{k,g} u_{gk} \pi_g \frac{1}{n} \sum_{i=1}^{n} \frac{\pi_{g|i}}{\pi_g} \left[ \log\left(\frac{\pi_{g|i}}{\pi_g}\right) - \log\left(\frac{p_{k|i}}{p_k}\right) \right].$$ (11)

Formula (11) suggests in a very natural way a partitioning method based on its minimization with respect to $\mathbf{U}$. The minimization of (11) guarantees the maximum internal cohesion of the clusters. However, we should ask: what about the external

isolation? To give an answer to this question, we should find a way to measure the external isolation, i.e. the cluster separation. In coherence with (7), we can define the Between deviance as the weighted sum of the KL divergences among the cluster profiles and their barycenter, that is the vector of ones even in this case. In formulas

$$D_B(\mathbf{U}) = \sum_{k=1}^{K} p_k \text{KL}(\mathbf{p}_k, \mathbf{1}) = \frac{1}{n} \sum_{k=1}^{K} \sum_{i=1}^{n} p_{k|i} \log\left(\frac{p_{k|i}}{p_k}\right). \tag{12}$$

Once again, we have a formula suggesting, in a very natural way, a partitioning method corresponding to the maximization of (12) with respect to $\mathbf{U}$. However, it is possible to show that, as in the case of $K$-means, the minimization of $D_W$ is equivalent to the maximization of $D_B$ because the sum of the two is constant and equal to $D_T$.

The Within deviance (11) can be minimized by using a coordinate descent algorithm that minimize (11) with respect to $\mathbf{U}$ and the centroids $\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_K$.



**Figure 1:** Simulated sample from a 6-component mixture along with the estimated classification in 7 components given by a homoscedastic mixture of Gaussians.

In order to check if the proposed method works properly, we considered the simulated dataset analysed in section 4.1 of Baudry et al. (2010). It is a sample from a finite mixture of 6 bivariate Gaussians. In Figure 1, the data is shown along with the classification in 7 components given by a homoscedastic mixture of Gaussians where the number of components have been selected by using BIC and the parameters estimated by maximum likelihood. Looking at the figure, it is not clear if the number of true clusters is 2 or 4 and then we considered both. The algorithm has been run for $K = 2$ and 4, from several different starting points. The technique aggregated the 7 components as {1,2,3,4} and {5,6,7} for $K = 2$ and as {1,2}, {3,4}, {5,6} and {7} for $K = 4$. From Figure 1, it is clear that in both cases, KL-components has been able to find the correct aggregation of the components.

# 3 Final comments and extensions

The method here presented has been extended to the use of divergences different from the KL. In particular, we have proven that all the properties of the KL-components technique shown in the previous section do hold if a Bregman divergence (Bregman, 1967) is used. The proofs are not reported here for the sake of space.

In practical applications, especially when the clusters of components are not well separated, the coordinate descent algorithm quite often remains trapped into a local optimum. Frequently, this problem can be simply solved by starting several times the algorithm from different random partitions. However, there is still the need to have a method able to produce good rational, non random, starting points. To this end, a hierarchical clustering procedure has been proposed, not shown here for the sake of space, where at each step two clusters are merged by minimizing the increment of Within deviance. The hierarchical solution is then used to start the partitioning algorithm.

The aim of our method is not to find the true number of clusters, even if it can help us in this task exploring the possible components aggregations. In this respect, further insights can be obtained by looking at the plot of $D_W$ vs $K$, computing the Calinski-Harabasz index (1974) or any other index based on $D_W$ and/or $D_B$.

We conclude by noting that the results presented here do not use the assumption that the components of the mixture are normally distributed. It follows that the proposed techniques can be used to cluster components that are not Gaussians.

# References

1. Baudry, J.P., Raftery, A., Celeux, A., Lo, K., Gottardo, R.: Combining mixture components for clustering, *Journal of Computational and Graphical Statistics*, 19, 332-353 (2010).
2. Bregman, L. M.: The relaxation method of finding the common points of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics* 7: 200–217 (1967)
3. Calinski, T., Harabasz, J.: A dendrite method for cluster analysis.Com-munications in Statistics, 3, no. 1:1–27 (1974).
4. Comas-Cufi M., Martin-Fernandez J.A., Mateu-Figueras G.: Merging the components of a finite mixture using posterior probabilities, *Statistical Modelling*, 19(2), 1–31 (2017)
5. Cormack, R. M.: A review of classification (with discussion). *Journal of the Royal Statistical Society*, A, 134, 321-67 (1971)
6. Di Zio M., Guarnera U., Rocci R.: A mixture of mixture models for a classification problem: the unity measure error. *Computational Statistics and Data Analysis*, 51, 5, 2573-2585 (2007)
7. Hennig, C.: Methods for merging Gaussian mixture components. *Advances in Data Analysis and Classification* 4, 3–34 (2010)
8. Hennig, C., Meila, M., Murtagh, F., Rocci R.: *Handbook of Cluster Analysis*, Chapman and Hall/CRC (2015)
9. Li J.: Clustering based on a multilayer mixture model. *Journal of Computational and Graphical Statistics* 14:547–568 (2004)
10. McLachlan, G.J., Peel, D.: *Finite Mixture Models*. Wiley, New York (2000)
11. McNicholas, P.D.: *Mixture Model-Based Classification*, Boca Raton FL: Chapman & Hall/CRC Press (2016)