

Variational Bayesian Approximation method for Classification and Clustering with a mixture of Student-t model

Ali Mohammad-Djafari

Laboratoire des signaux et systèmes (L2S),
CNRS-CentraleSupélec-UNIV PARIS SUD, Gif-sur-Yvette, FRANCE,
djafari@lss.supelec.fr,
WWW home page: <http://djafari.free.fr>

Abstract. Clustering, classification and Pattern Recognition in a set of data are between the most important tasks in statistical researches and in many applications. In this paper, we propose to use a mixture of Student-t distribution model for the data via a hierarchical graphical model and the Bayesian framework to do these tasks. The main advantages of this model is that the model accounts for the uncertainties of variances and covariances and we can use the Variational Bayesian Approximation (VBA) methods to obtain fast algorithms to be able to handle large data sets.

1 Introduction

Clustering and classification of a set of data are not trivial problems. In fact, we can consider them as ill-posed inverse problems in which the solutions are not unique. Mixture models are natural ones for classification and clustering [1–8]. The Mixture of Gaussians (MoG) models have been used very extensively [9–11]. In this paper, we propose to use a mixture of Student-t model and a Bayesian framework for these tasks. The main advantages of this model is that the model accounts for the uncertainties of variances and covariances and we can use the Variational Bayesian Approximation (VBA) methods to obtain fast algorithms as well. Even if this model may have been used before [12–26], here we propose a novel unifying presentation for all the steps: training, supervised or semi-supervised classification and clustering (non-supervised). We also use VBA framework and some simplifications to develop fast algorithms to be able to handle big data sets.

2 Mixture models for classification and clustering

A mixture model is generally given as:

$$p(\mathbf{x}|\mathbf{a}, \boldsymbol{\Theta}, K) = \sum_{k=1}^K a_k p_k(\mathbf{x}_k|\boldsymbol{\theta}_k), \quad (1)$$

where K is the number of classes, $\mathbf{a} = \{a_k, k = 1, \dots, K\}$ the proportion parameters and $\Theta = \{\theta_k, k = 1, \dots, K\}$ all the other parameters of the model. If we assume different classes can be modelled by the same family $p_k(\mathbf{x}_k|\theta_k) = p(\mathbf{x}_k|\theta_k)$ and introduce a hidden class variable $c_n \in \{1, \dots, K\}$, then for a given sample \mathbf{x}_n in class k we can write:

$$p(\mathbf{x}_n|c_n = k, \theta_k) = p(\mathbf{x}_n|\theta_k) \quad (2)$$

or

$$p(\mathbf{x}_n, c_n = k|a_k, \theta_k, K) = a_k p(\mathbf{x}_n|\theta_k). \quad (3)$$

The Mixture of Gaussians (MoG) corresponds to the case where $p(\mathbf{x}_n|c_n = k, \theta_k) = \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$ with $\theta_k = (\mu_k, \Sigma_k)$.

Now, imagine a set of data $\mathbf{X} = \{\mathbf{x}_n, n = 1, \dots, N\}$ where each element \mathbf{x}_n can be in one of these classes. Then, we can write:

$$p(\mathbf{X}_n, c_n = k|\mathbf{a}, \theta) = \prod_{n=1}^N p(\mathbf{x}_n, c_n = k|\mathbf{a}, \theta). \quad (4)$$

Noting by $\mathbf{c} = \{c_n, n = 1, N\}$ with $c_n \in \{1, \dots, K\}$, $\mathbf{a} = \{a_k, k = 1, \dots, K\}$ and $\Theta = \{\theta_k, k = 1, \dots, K\}$, we have:

$$\begin{aligned} p(\mathbf{X}_n, \mathbf{c}|\mathbf{a}, \Theta, K) &= \prod_{n=1}^N \prod_{k=1}^K p(c_n = k) p(\mathbf{x}_n|\theta_k) \\ &= \prod_{n=1}^N \prod_{k=1}^K a_k p(\mathbf{x}_n|\theta_k). \end{aligned} \quad (5)$$

The classification problems can then be summarized as follows:

Training:

Given a set of (training) data \mathbf{X} and classes \mathbf{c} , estimate the parameters \mathbf{a} and Θ . The classical frequentist method is the Maximum Likelihood (ML) which defines the solution as

$$(\hat{\mathbf{a}}, \hat{\Theta}) = \arg \max_{(\mathbf{a}, \Theta)} \{p(\mathbf{X}, \mathbf{c}|\mathbf{a}, \Theta, K)\}. \quad (6)$$

The Bayesian way is to assign priors $p(\mathbf{a}|K)$ and $p(\Theta|K) = \prod_{k=1}^K p(\theta_k)$, then the joint posterior laws is given by:

$$p(\mathbf{a}, \Theta|\mathbf{X}, \mathbf{c}, K) = \frac{p(\mathbf{X}, \mathbf{c}|\mathbf{a}, \Theta, K) p(\mathbf{a}|K) p(\Theta|K)}{p(\mathbf{X}, \mathbf{c}|K)} \quad (7)$$

where

$$p(\mathbf{X}, \mathbf{c}|K) = \iint p(\mathbf{X}, \mathbf{c}|\mathbf{a}, \Theta|K) p(\mathbf{a}|K) p(\Theta|K) d\mathbf{a} d\Theta \quad (8)$$

from which we can deduce $\hat{\mathbf{a}}$ and $\{\hat{\theta}_k, k = 1, \dots, K\}$ either as the Maximum A Posteriori (MAP) or Posterior Mean (PM).

Supervised classification:

For a given sample \mathbf{x}_m and given the parameters K , \mathbf{a} and Θ determine

$$p(c_m = k | \mathbf{x}_m, \mathbf{a}, \Theta, K) = \frac{p(\mathbf{x}_m, c_m = k | \mathbf{a}, \Theta, K)}{p(\mathbf{x}_m | \mathbf{a}, \Theta, K)} \quad (9)$$

where $p(\mathbf{x}_m, c_m = k | \mathbf{a}, \Theta, K) = a_k p(\mathbf{x}_m | \theta_k)$ and

$$p(\mathbf{x}_m | \mathbf{a}, \Theta, K) = \sum_{k=1}^K a_k p(\mathbf{x}_m | \theta_k) \quad (10)$$

and its best class k^* , for example the MAP solution:

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \mathbf{a}, \Theta, K)\}. \quad (11)$$

Semi-supervised classification:

For a given sample \mathbf{x}_m and given the parameters K and Θ , determine the probabilities

$$p(c_m = k | \mathbf{x}_m, \Theta, K) = \frac{p(\mathbf{x}_m, c_m = k | \Theta, K)}{p(\mathbf{x}_m | \Theta, K)} \quad (12)$$

where

$$p(\mathbf{x}_m, c_m = k | \Theta, K) = \int p(\mathbf{x}_m, c_m = k | \mathbf{a}, \Theta, K) p(\mathbf{a} | K) d\mathbf{a} \quad (13)$$

and

$$p(\mathbf{x}_m | \Theta, K) = \sum_{k=1}^K p(\mathbf{x}_m, c_m = k | \Theta, K) \quad (14)$$

and its best class k^* , for example the MAP solution:

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \Theta, K)\}. \quad (15)$$

Clustering or non-supervised classification:

Given a set of data \mathbf{X} , determine K and \mathbf{c} . When these are determined, we can also determine the characteristics of those classes \mathbf{a} and Θ . To do this we need the following relations:

$$p(K = L | \mathbf{X}) = \frac{p(\mathbf{X}, K = L)}{p(\mathbf{X})} = \frac{p(\mathbf{X} | K = L) p(K = L)}{p(\mathbf{X})} \quad (16)$$

and

$$p(\mathbf{X}) = \sum_{L=1}^{L_0} p(K = L) p(\mathbf{X} | K = L), \quad (17)$$

where L_0 is the a priori maximum number of classes and

$$p(\mathbf{X} | K = L) = \int \int \prod_n \prod_{k=1}^L a_k p(\mathbf{x}_n, c_n = k | \theta_k) p(\mathbf{a} | K) p(\Theta | K) d\mathbf{a} d\Theta. \quad (18)$$

As we will see later the main difficulty is the computation of these two last equations. The Variational Bayesian Approximation technics try to find upper and lower bounds for them.

3 Mixture of Student-t model

Let us consider the following representation of the Student-t probability density function (pdf):

$$\mathcal{T}(\mathbf{x}|\nu, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, z^{-1}\boldsymbol{\Sigma}) \mathcal{G}(z|\frac{\nu}{2}, \frac{\nu}{2}) dz, \quad (19)$$

where

$$\begin{aligned} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right] \\ &= |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\text{Tr}\{(\mathbf{x} - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})'\}\right], \end{aligned} \quad (20)$$

and

$$\mathcal{G}(z|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} z^{\alpha-1} \exp[-\beta z]. \quad (21)$$

Let us also consider the finite mixture of Student-t model:

$$p(\mathbf{x}|\{\nu_k, \mathbf{a}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, k = 1, \dots, K\}, K) = \sum_{k=1}^K a_k \mathcal{T}(\mathbf{x}_n|\nu_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad (22)$$

Introducing the hidden variables z_{nk} this model can be written via:

$$p(\mathbf{x}_n, c_n = k, z_{nk}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, K) = a_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, z_{nk}^{-1}\boldsymbol{\Sigma}_k) \mathcal{G}(z_{nk}|\frac{\nu_k}{2}, \frac{\nu_k}{2}). \quad (23)$$

Noting by: $\mathbf{Z} = \{z_{nk}\}$, $\mathbf{z}_k = \{z_{nk}, n = 1, \dots, N\}$, $\mathbf{c} = \{c_n, n = 1, \dots, N\}$, $\boldsymbol{\theta}_k = \{\nu_k, \mathbf{a}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$, $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_k, k = 1, \dots, K\}$ and assigning the priors $p(\boldsymbol{\Theta}) = \prod_k p(\boldsymbol{\theta}_k)$, we can write:

$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|K) = \prod_n \prod_k a_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, z_{n,k}^{-1}\boldsymbol{\Sigma}_k) \mathcal{G}(z_{n,k}|\frac{\nu_k}{2}, \frac{\nu_k}{2}) p(\boldsymbol{\theta}_k) \quad (24)$$

Then, the joint posterior law of all the unknowns $(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})$ given the data \mathbf{X} and K can be written as

$$p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, K) = \frac{p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|K)}{p(\mathbf{X}|K)}. \quad (25)$$

The main task now is to propose some approximations to it in such a way that we can use it easily in all the above mentioned tasks of classification or clustering. The main idea behind the VBA technics is exactly this.

4 Variational Bayesian Approximation (VBA)

4.1 Main idea

The main idea behind the VBA is to propose an approximation $q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})$ for $p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, K)$. This approximation can be in such a way that $\text{KL}(q : p)$ be minimized. Interestingly, by noting that $p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, K) = p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|K)/p(\mathbf{X}|K)$, it is easy to showed that

$$\text{KL}(q : p) = -\mathcal{F}(q) + \ln p(\mathbf{X}|K) \quad (26)$$

where

$$\mathcal{F}(q) = \langle -\ln p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) \rangle_q \quad (27)$$

is called free energy of q and we have the following properties:

- Maximizing $\mathcal{F}(q)$ or minimizing $\text{KL}(q : p)$ are equivalent and both give an upper bound to the evidence of the model $\ln p(\mathbf{X} | K)$.
- When the optimum q^* is obtained, $\mathcal{F}(q^*)$ can be used as a criterion for model selection.
- If p is in the exponential family, then choosing appropriate conjugate priors, the structure of q will be the same and we can obtain appropriate fast optimization algorithms.

In our case, noting that

$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) = \prod_n \prod_k p(\mathbf{x}_n, c_n, z_{nk} | a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \nu_k) \prod_k [p(\alpha_k) p(\beta_k) p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) p(\boldsymbol{\Sigma}_k)] \quad (28)$$

with

$$p(\mathbf{x}_n, c_n, z_{nk} | a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \nu_k) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, z_{n,k}^{-1} \boldsymbol{\Sigma}_k) \mathcal{G}(z_{nk} | \alpha_k, \beta_k) \quad (29)$$

is separable, in one side for $[\mathbf{c}, \mathbf{Z}]$ and in other size in components of $\boldsymbol{\Theta}$, we propose to use

$$q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}) = q(\mathbf{c}, \mathbf{Z}) q(\boldsymbol{\Theta}). \quad (30)$$

With this decomposition, the expression of the Kullback-Leibler divergence becomes:

$$\text{KL}(q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta}) : p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | \mathbf{X}, K)) = \sum_{\mathbf{c}} \int \int q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta}) \ln \frac{q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta})}{p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | \mathbf{X}, K)} d\boldsymbol{\Theta} d\mathbf{Z} \quad (31)$$

and the expression of the Free energy becomes:

$$\mathcal{F}(q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta})) = \sum_{\mathbf{c}} \int \int q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta}) \ln \frac{p(\mathbf{X}, \mathbf{c}, \mathbf{Z} | \boldsymbol{\Theta}, K) p(\boldsymbol{\Theta} | K)}{q_1(\mathbf{c}, \mathbf{Z}) q_2(\boldsymbol{\Theta})} d\boldsymbol{\Theta} d\mathbf{Z}. \quad (32)$$

In the following we propose appropriate priors and obtain the expressions of q and appropriate fast algorithms.

5 Proposed VBA for Mixture of Student-t priors model

As we discussed in previous section, here we consider the Mixture of Student-t priors model and propose appropriate conjugate priors and appropriate factorized form for the testing or approximation q and finally give the details of the parameters updating algorithm. To be able to propose conjugate priors for all the parameters, we change slightly the model by replacing ν_k in the Gamma expression $\mathcal{G}(z_{n,k} | \frac{\nu_k}{2}, \frac{\nu_k}{2})$ of the Student-t expression by two parameters $\mathcal{G}(z_{n,k} | \alpha_k, \beta_k)$:

$$p(\mathbf{x}_n, c_n = k, z_{nk} | a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \beta_k, K) = a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, z_{nk}^{-1} \boldsymbol{\Sigma}_k) \mathcal{G}(z_{nk} | \alpha_k, \beta_k). \quad (33)$$

The final hierarchical model that we propose is shown in the Figure 1.

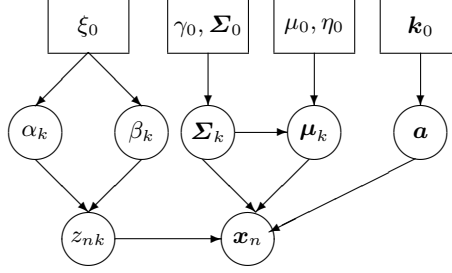


Fig. 1. Graphical representation of the model.

5.1 Conjugate priors

In the following, noting by $\boldsymbol{\Theta} = \{(a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \beta_k), k = 1, \dots, K\}$, we propose to use the factorized prior laws:

$$p(\boldsymbol{\Theta}) = p(\mathbf{a}) \sum_k [p(\alpha_k) p(\beta_k) p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) p(\boldsymbol{\Sigma}_k)] \quad (34)$$

with the following components:

$$\begin{cases} p(\mathbf{a}) = \mathcal{D}(\mathbf{a} | \mathbf{k}_0), & \mathbf{k}_0 = [k_0, \dots, k_0] = k_0 \mathbf{1} \\ p(\alpha_k) = \mathcal{E}(\alpha_k | \zeta_0) = \mathcal{G}(\alpha_k | \mathbf{1}, \zeta_0) \\ p(\beta_k) = \mathcal{E}(\beta_k | \zeta_0) = \mathcal{G}(\beta_k | \mathbf{1}, \zeta_0) \\ p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \boldsymbol{\mu}_0 \mathbf{1}, \eta_0^{-1} \boldsymbol{\Sigma}_k) \\ p(\boldsymbol{\Sigma}_k) = \mathcal{IW}(\boldsymbol{\Sigma}_k | \gamma_0, \gamma_0 \boldsymbol{\Sigma}_0) \end{cases} \quad (35)$$

where

$$\mathcal{D}(\mathbf{a} | \mathbf{k}) = \frac{\Gamma(\sum_l k_l)}{\prod_l \Gamma(k_l)} \prod_l a_l^{k_l - 1} \quad (36)$$

is the Dirichlet pdf,

$$\mathcal{E}(t | \zeta_0) = \zeta_0 \exp[-\zeta_0 t] \quad (37)$$

is the Exponential pdf,

$$\mathcal{G}(t | a, b) = \frac{b^a}{\Gamma(a)} t^{a-1} \exp[-bt] \quad (38)$$

is the Gamma pdf and

$$\mathcal{IW}(\boldsymbol{\Sigma}|\gamma, \gamma \boldsymbol{\Delta}) = \frac{|\frac{1}{2}\boldsymbol{\Delta}|^{\gamma/2} \exp[-\frac{1}{2}\text{Tr}\{\boldsymbol{\Delta}\boldsymbol{\Sigma}^{-1}\}]}{\Gamma_D(\gamma/2)|\boldsymbol{\Sigma}|^{\frac{\gamma+D+1}{2}}}. \quad (39)$$

is the inverse Wishart pdf.

With these prior laws and the likelihood: $p(\mathbf{x}_n|c(n), \mathbf{z}_k(n), \boldsymbol{\Theta}, k)$ we can obtain the joint posterior law:

$$p_k(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})}{p(\mathbf{X})}. \quad (40)$$

Now, we have to choose a factored form for q in such a way that we can transform the optimization of the $KL(q : p)$ or the free energy $\mathcal{F}(q)$ to the updating of the parameters of the different components of q . We propose to use the following decomposition:

$$\begin{aligned} q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}) &= q(\mathbf{c}, \mathbf{Z}) q(\boldsymbol{\Theta}) = \\ &= \prod_n \prod_k [q(c_n = k|z_{nk}) q(z_{nk})] \\ &= \prod_k [q(\alpha_k) q(\beta_k) q(\boldsymbol{\mu}_k|\boldsymbol{\Sigma}_k) q(\boldsymbol{\Sigma}_k)] q(\mathbf{a}). \end{aligned} \quad (41)$$

with:

$$\begin{cases} q(\mathbf{a}) = \mathcal{D}(\mathbf{a}|\tilde{\mathbf{k}}), & \tilde{\mathbf{k}} = [\tilde{k}_1, \dots, \tilde{k}_K] \\ q(\alpha_k) = \mathcal{G}(\alpha_k|\tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\beta_k) = \mathcal{G}(\beta_k|\tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\boldsymbol{\mu}_k|\boldsymbol{\Sigma}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\eta}}^{-1}\boldsymbol{\Sigma}_k) \\ q(\boldsymbol{\Sigma}_k) = \mathcal{IW}(\boldsymbol{\Sigma}_k|\tilde{\gamma}, \tilde{\gamma}\tilde{\boldsymbol{\Sigma}}) \end{cases} \quad (42)$$

With these choices, we have

$$\begin{aligned} \mathcal{F}(q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})) &= \langle \ln p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}|K) \rangle_{q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})} \\ &= \prod_k \prod_n \mathcal{F}_{1_{kn}} + \prod_k \mathcal{F}_{2_k} \end{aligned} \quad (43)$$

with

$$\begin{aligned} \mathcal{F}_{1_{kn}} &= \langle \ln p(\mathbf{x}_n, c_n, z_{nk}, \boldsymbol{\theta}_k) \rangle_{q(c_n=k|z_{nk})q(z_{nk})} \\ \mathcal{F}_{2_k} &= \langle \ln p(\mathbf{x}_n, c_n, z_{nk}, \boldsymbol{\theta}_k) \rangle_{q(\boldsymbol{\theta}_k)} \end{aligned} \quad (44)$$

Now, to obtain the expressions of the updating expressions of the tilded parameters, we need to go to the following three steps:

- E step: Optimizing \mathcal{F} with respect to $q(\mathbf{c}, \mathbf{Z})$ when keeping $q(\boldsymbol{\Theta})$ fixed, we obtain the expression of $q(c_n = k|z_{nk}) = \tilde{a}_k$, $q(z_{nk}) = \mathcal{G}(z_{nk}|\tilde{\alpha}_k, \tilde{\beta}_k)$.
- M step: Optimizing \mathcal{F} with respect to $q(\boldsymbol{\Theta})$ when keeping $q(\mathbf{c}, \mathbf{Z})$ fixed, we obtain the expression of $q(\mathbf{a}) = \mathcal{D}(\mathbf{a}|\tilde{\mathbf{k}})$, $\tilde{\mathbf{k}} = [\tilde{k}_1, \dots, \tilde{k}_K]$, $q(\alpha_k) = \mathcal{G}(\alpha_k|\tilde{\zeta}_k, \tilde{\eta}_k)$, $q(\beta_k) = \mathcal{G}(\beta_k|\tilde{\zeta}_k, \tilde{\eta}_k)$, $q(\boldsymbol{\mu}_k|\boldsymbol{\Sigma}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\eta}}^{-1}\boldsymbol{\Sigma}_k)$, and $q(\boldsymbol{\Sigma}_k) = \mathcal{IW}(\boldsymbol{\Sigma}_k|\tilde{\gamma}, \tilde{\gamma}\tilde{\boldsymbol{\Sigma}})$, which gives the updating algorithm for the corresponding tilded parameters.

- \mathcal{F} evaluation After each E step and M step, we can also evaluate the expression of $\mathcal{F}(q)$ which can be used for stopping rule of the iterative algorithm. Also, final value of this expression for each value of K , noted \mathcal{F}_k , can be used as a criterion for the model selection, i.e.; the determination of the number of clusters.

The expressions of all the tilded parameters update as well as the expression of \mathcal{F}_K are easily obtained thanks to the properties of the conjugate priors. However, these expressions are cumbersome and will be given in the appendix.

6 Conclusion

Clustering and classification of a set of data are between the most important tasks in statistical researches for many applications such as data mining in biology. Mixture models and in particular Mixture of Gaussians are classical models for these tasks. In this paper, we proposed to use a mixture of Student-t distribution model for the data via a hierarchical graphical model. Then, we proposed a Bayesian framework to do these tasks. The main advantages of this model is that the model accounts for the uncertainties of variances and covariances and we can use the Variational Bayesian Approximation (VBA) methods. To obtain fast algorithms and be able to handle large data sets, we used conjugate priors everywhere it was possible. The proposed algorithm has been used for clustering, classification and discriminant analysis of some biological data, but in this paper, we only presented the main algorithm.

Acknowledgment

This work has been supported partially by the C5SYS (<https://www.erasysbio.net/index.php?index=272>) project of ERASYSBIO.

References

1. Forgy, E.W.: Cluster analysis of multivariate data: efficiency vs interpretability of classifications. *Biometrics* (1965)
2. MacKay, D.J.C.: A practical Bayesian framework for backpropagation networks. *Neural Computation* **4** (1992) 448–472
3. Redner, R., Walker, H.: Mixture densities, maximum likelihood and the em algorithm. *SIAM Review* **26** (1984)
4. Husmeier, D., Penny, W., Roberts, S.: Empirical evaluation of Bayesian sampling for neural classifiers. In L.Niklason, M., T.Ziemke, eds.: *ICANN 98: Proceedings of the 8th International Conference on Artificial Neural Networks*. (1998)
5. Lee, T., Lewicki, M., Sejnowski, T.: Independent component analysis using an extended infomax algorithm for mixed sub-gaussian and super-gaussian sources. *Neural Computation* **11** (1999) 409–433
6. A., H., E., O.: Independent component analysis: Algorithms and applications. *Neural Networks* **13** (2000) 411–430

7. Ma, J., Xu, L., Jordan, M.I.: Asymptotic convergence rate of the EM algorithm for Gaussian mixtures. *Neural Computation* **12** (2001) 2881–2907
8. Nielsen, F., Nock, R.: Clustering multivariate normal distributions. *Neural Computation* **Springer-Verlag, Berlin, Heidelberg** (2009) 164–174
9. Quandt, R., Ramsey, J.: Estimating mixtures of normal distributions and switching regressions. *J. American Statistical Association* **73** (1978) 2
10. Hathaway, R.: A constrained formulation of maximum-likelihood estimation for normal mixture distributions. *Annals of Statistics* **13** (1985) 1
11. Hastie, T., Tibshirani, R.: Discriminant analysis by Gaussian mixture. *Journal of Royal Statistical Society* **58** (1996) 155–176
12. Dempster, A., Laird, N., Rubin, D.: Maximum likelihood from incomplete data via the em algorithm. *J. R. Statist. Soc. B* **39** (1977) 3
13. Neal, R., Hinton, G.: A view of the EM algorithm that justifies incremental, sparse, and other variants. *Learning in graphical models* **89** (1998) 355–368
14. Tipping, M.E., Bishop, C.M.: Mixtures of probabilistic principal components analysis. *Neural Computation* **11** (1999) 443–482
15. Jordan, M., Ghahramani, Z., Jaakkola, T., , Saul, L.: An introduction to variational methods for graphical models. *Machine Learning* **37** (2006) 183–233
16. Jaakkola, T.S., Jordan, M.I.: Bayesian parameter estimation via variational methods. *Statistics and Computing* **10** (2000) 25–37
17. Friston, K., Penny, W.: Bayesian inference and posterior probability maps. In: *Proceedings of the 9th International Conference on Neural Information Processing (ICONIP'02)*. (2002) 413–417
18. Winn, J., Bishop, C.M., Jaakkola, T.: Variational message passing. *Journal of Machine Learning Research* **6** (2005) 661–694
19. David, M.B., Michael, I.J.: Variational inference for the dirichlet process mixtures. **1** (2006) 121–144
20. Beal, M.: Variational Algorithms for Approximate Bayesian Inference. PhD thesis, Gatsby Computational Neuroscience Unit, University College London (2003)
21. Beal, M., Ghahramani, Z.: Variational Bayesian learning of directed graphical models with hidden variables. *Bayesian Statistics* **1** (2006) 793–832
22. Kim, H., Ghahramani, Z.: Bayesian gaussian process classification with the em-ep algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **28** (2006) 1948–1959
23. Nasios, N., Bors, A.: Variational learning for gaussian mixture models. *IEEE Transactions on Systems, Man and Cybernetics, Part B* **36** (2006) 849–862
24. Ghahramani, Z., Griffiths, T., Sollich, P.: Bayesian nonparametric latent feature models. *Bayesian Statistics* **8** (2007)
25. McGrory, C., Titterton, D.: Variational approximations in Bayesian model selection for finite mixture distributions. *Computational Statistics and Data Analysis* **51** (2007) 5352–5367
26. Qiao, Z., Zhou, L., Huang, J.Z.: Sparse linear discriminant analysis with applications to high dimensional low sample size data. *International Journal of Applied Mathematics* **39** (2008) 48–60