

A particular Gaussian mixture model for clustering and its application to image retrieval

Hichem Sahbi

© Springer-Verlag 2007

Abstract We introduce a new method for data clustering based on a particular Gaussian mixture model (GMM). Each cluster of data, modeled as a GMM into an input space, is interpreted as a hyperplane in a high dimensional mapping space where the underlying coefficients are found by solving a quadratic programming (QP) problem. The main contributions of this work are (1) an original probabilistic framework for GMM estimation based on QP which only requires finding the mixture parameters, (2) this QP is interpreted as the minimization of the pairwise correlations between cluster hyperplanes in a high dimensional space and (3) it is solved easily using a new decomposition algorithm involving trivial linear programming sub-problems. The validity of the method is demonstrated for clustering 2D toy examples as well as image databases.

Keywords Gaussian mixture models · Clustering · Kernel methods and image retrieval

1 Introduction

Consider a training set $\mathcal{S}_N = \{x_1, \dots, x_N\}$ generated i.i.d. according to an unknown but a fixed probability distribution $P(X)$. Here X is a random variable standing for training data in an *input space* denoted \mathcal{X} . Our goal is to define a

partition of this training set, i.e., assign each training example to its actual cluster y_k , $k = 1, \dots, C$. Existing clustering methods can be categorized depending on the “crispness” of data memberships, i.e., whether each training example belongs to only one cluster or not. The first category includes hierarchical approaches (Jain and Dubes 1988; Posse 2001), EM (expectation maximization) (Dempster et al. 1977), C (or K) means (MacQueen 1965), self organizing maps (SOMs) (Tsao et al. 1994) and recently original approaches based on kernel methods (Ben-Hur et al. 2000; Bach and Jordan 2003; Wu and Schoelkopf 2006). In the second family the membership of a training example to one or another cluster is fuzzy and a family of algorithms dealing with fuzziness exists in the literature for instance (Yang 1993; Bezdek 1981; Frigui and Krishnapuram 1999; Dave 1991; Ichihashi et al. 2000; Sahbi and Boujemaa 2005). All these methods have been used in different domains including Gene expression (Orengo et al. 2003), image segmentation (Carson et al. 2002) and database categorization (Le-Saux and Boujemaa 2002).

One of the main issues in the existing clustering methods remains setting the appropriate number of classes for a given problem. Many methods for instance hierarchical clustering (Sneath and Sokal 1973; Fraley 1998; Posse 2001) have proven to perform well when the application allows us to know a priori the number of clusters or when the user sets it manually. Of course, the estimation of this number is application-dependent, for instance in image segmentation it can be set a priori to the number of targeted regions. Unfortunately, for some applications such as database categorization, it is not always possible to predict automatically and even manually the appropriate number of classes.

Given a training set \mathcal{S}_N in the input space \mathcal{X} , in our method each cluster in \mathcal{X} is modeled as a GMM and interpreted as a hyperplane in a high dimensional space referred to as the mapping space, denoted \mathcal{H} . *Clustering consists in*

H. Sahbi (✉)
Machine Intelligence Laboratory, Department of Engineering,
Cambridge University, Cambridge, UK
e-mail: hs385@cam.ac.uk

H. Sahbi
Certis Laboratory, Ecole Nationale des Ponts et Chaussees,
Paris, France
e-mail: sahbi@certis.enpc.fr

maximizing a likelihood function of data memberships and is equivalent to finding the parameters of the cluster hyperplanes by solving a QP problem. We will show that this QP can be tackled efficiently by solving trivial linear programming sub-problems. Notice that when solving this QP, the number of clusters (denoted C), fixed initially, might be overestimated and this leads to several overlapping clusters; therefore the actual clusters are found as constellations of highly correlated hyperplanes in the mapping space \mathcal{H} .

In the remainder of this paper, we refer to a cluster as a set of data gathered using an algorithm while a class (or a category) is the actual membership of this data according to a well defined ground truth. Among notations, i and j stand for data indices while k, c and l stand for cluster indices. Other notations will be introduced as we go along through different sections of this paper which is organized as following: in Sect. 2 we provide a short reminder on GMMs followed by our clustering minimization problem in Sect. 3. In Sect. 4 we show that this minimization problem can be solved as a succession of simple linear programming subproblems which are handled efficiently. We present in Sect. 5 experiments on simple as well as challenging problems in content based image retrieval. Finally, we conclude in Sect. 6 and we provide some directions for a future work.

2 A short reminder on GMMs

Given a training set \mathcal{S}_N of size N , we are interested in a particular Gaussian mixture model (denoted \mathcal{M}_k) Bishop (1995) where the number of its components is equal to the size of the training set. The parameters of \mathcal{M}_k are denoted $\Theta_k = \{\Theta_{i|k}, i = 1, \dots, N\}$. Here $\Theta_{i|k}$ stands for the i th component parameter of \mathcal{M}_k , i.e., the mean and the covariance matrices of a Gaussian density function. In this case, the output of the likelihood GMM function related to a cluster y_k is a weighted sum of N component densities:

$$p(x|y_k) = \sum_i^N P(\Theta_{i|k}|y_k) p(x|y_k, \Theta_{i|k}), \tag{1}$$

here $P(\Theta_{i|k}|y_k)$ (also denoted μ_{ik}) is the prior probability of the i th component parameter $\Theta_{i|k}$ given the cluster y_k and $p(x|y_k, \Theta_{i|k}) = \mathcal{N}_k(x; \Theta_{i|k})$ is the normal density function of the i th component. In order to guarantee that $p(x|y_k)$ is a density function, the mixture parameters $\{\mu_{ik}\}$ are chosen such that $\sum_i \mu_{ik} = 1$. Usually the training of GMMs can be formulated as a maximum likelihood problem where the parameters Θ_k and the mixture parameters $\{\mu_{ik}\}$ are estimated using expectation maximization (Dempster et al. 1977; Bishop 1995).

We consider in our work, $\mathcal{N}_k(x; \Theta_{i|k})$ with a fixed mean $x_i \in \mathcal{S}_N$ and covariance matrix $\sigma \Sigma_i$ ($\Sigma_i \in \mathbb{R}^{p \times p}$,

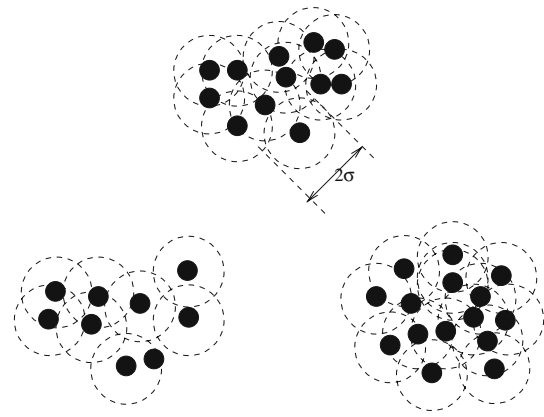


Fig. 1 This figure shows a particular GMM model where the centers and the variances are fixed. The only free parameters are the GMM mixture coefficients and the scale σ

$\sigma \in \mathbb{R}^+$). Now, each cluster y_k is modeled as a GMM where the *only* free parameters are the mixture coefficients $\{\mu_{ik}\}$; the means and the covariances are assumed constant but of course dependent on a priori knowledge of the training set (see Fig. 1 and Sect. 5).

3 Clustering

The goal of a clustering algorithm is to make clusters, containing data from different classes, as different as possible while keeping data from the same classes as close as possible to their *actual* clusters. Usually this implies the optimization of an objective function (see for instance Frigui and Krishnapuram 1997) involving a fidelity term which measures the fitness of each training sample to its model and a regularizer which reduces the number of clusters, i.e., the complexity of the model.

3.1 Our formulation

Following notations and definitions in Sect. 2, a given $x \in \mathbb{R}^p$ is assigned to a cluster y_k if:

$$y_k = \arg \max_{y_l} p(x|y_l), \tag{2}$$

here the weights $\mu = \{\mu_{il}\}$, of the GMM functions $p(x|y_l)$ $l = 1, \dots, C$, are found by solving the following constrained minimization problem (see motivation below):

$$\begin{aligned} \min_{\mu} \sum_{k,i} \mu_{ik} & \left(\sum_{l \neq k} p(x_i|y_l) \right) \\ \text{s.t.} \sum_i \mu_{ic} &= 1, \quad \mu_{ic} \in [0, 1], \quad i = 1, \dots, N, \\ & c = 1, \dots, C \end{aligned} \tag{3}$$

In the above objective function, the training data belonging to a cluster y_l are assumed drawn from a GMM with a likelihood function $p(x|y_l)$. Each mixture parameter μ_{il} stands for the degree of membership (or the contribution) of the training example x_i to the cluster y_l . The overall objective is to maximize the membership of each training example to its actual cluster while keeping the memberships to the remaining clusters relatively low. Usually, existing clustering methods (see for instance [Bezdek 1981](#)) find the parameters $\{\mu_{il}\}$ as those which maximize the membership of each training example to its actual cluster. In contrast to these methods, our formulation proceeds using a dual principle; *the purpose is to minimize the memberships of training examples to their non-actual clusters*.

Using (1), we can expand the objective function (3) as:

$$\min_{\mu} \sum_{k \neq l} \sum_{i,j} \mu_{ik} \mu_{jl} \mathcal{N}_l(x_i; \Theta_{j|l}), \tag{4}$$

here $\mathcal{N}_l(x_i; \Theta_{j|l})$ is the response of a Gaussian density function (also referred to as kernel), with fixed parameters $\Theta_{j|l} = \{x_j, \sigma, \Sigma_j\}$; x_j is the mean, Σ_j is the covariance and σ is the scale. We will denote this kernel simply as $K_{\sigma}(\|x_i - x_j\|)$. It is known that the Gaussian kernel is positive definite ([Cristianini and Shawe-Taylor 2000](#)) so this function corresponds to a scalar product in the mapping space \mathcal{H} , i.e., there is a mapping Φ_{σ} from the input space into an infinite dimensional space such that $K_{\sigma}(\|x_i - x_j\|) = \langle \Phi_{\sigma}(x_i), \Phi_{\sigma}(x_j) \rangle$ where $\langle \rangle$ stands for the inner product in \mathcal{H} . At this stage, the response of a GMM function $p(x|y_k)$ is equal to the scalar product $\langle \omega_k, \Phi_{\sigma}(x) \rangle$ where $\omega_k = \sum_i \mu_{ik} \Phi_{\sigma}(x_i)$ is the normal of a hyperplane in \mathcal{H} (see Fig. 2). Now, the objective function (4) can be rewritten:

$$\min_{\mu} \sum_{k \neq l} \langle \omega_k, \omega_l \rangle \tag{5}$$

The above objective function minimizes the sum of hyperplane *correlations* taken pairwise among all different clusters. Now, we can derive the new form of the constrained minimization problem (3):

$$\begin{aligned} \min_{\mu} & \sum_{k,i} \sum_{l \neq k,j} \mu_{ik} \mu_{jl} K_{\sigma}(\|x_i - x_j\|) \\ \text{s.t.} & \sum_i \mu_{ic} = 1, \quad \mu_{ic} \in [0, 1], \quad i = 1, \dots, N, \\ & c = 1, \dots, C \end{aligned} \tag{6}$$

This defines a constrained QP which can be solved using standard QP libraries (see for instance [Vanderbei 1999](#)). When solving this problem, training examples $\{x_i\}$ for which the mixing parameter $\{\mu_{ik}\}$ are positive will be referred to as the *GMM vectors* of the cluster y_k (see Fig. 2).

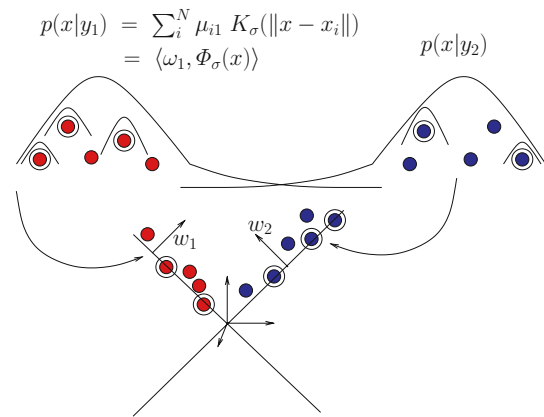


Fig. 2 This figure shows the mapping of training samples into a high dimensional space. Data in the original space characterizes Gaussian blobs while in the mapping space they correspond to hyperplanes. The GMM vectors are surrounded with circles and correspond to the centers of the Gaussian kernels for which the mixture parameters do not vanish

4 Training

The number of parameters intervening in (6) is $N \times C$, so for clustering problems of reasonable size, for instance $N = 1.000$ and $C = 20$, solving this QP, using standard packages, can quickly get out of hand. Chunking methods have been successfully used to solve QP for large scale training problems such as SVM ([Osuna et al. 1997](#)). The idea consists in solving a QP problem using an active subset of parameters, referred to as a *chunk*, which is updated iteratively. When the QP is convex, by checking that the Gram matrix is positive definite, the process is guaranteed to converge to the global optimum after a sufficient number of iterations ([Osuna et al. 1997](#)).

Using the same principle as [Osuna et al. \(1997\)](#) and [Platt \(1999\)](#), we will show in this section that for a particular choice of the active chunk, the QP (6) can be decomposed into linear programming subproblems each one can be solved trivially.

4.1 Decomposition

Let us fix one cluster index $p \in \{1, \dots, C\}$ and rewrite the objective function (6) as:

$$\begin{aligned} \min_{\mu} & 2 \sum_i \mu_{ip} c_{ip} \\ & + \sum_{i,k \neq p} \sum_{j,l \neq k,p} \mu_{ik} \mu_{jl} K_{\sigma}(\|x_i - x_j\|), \end{aligned} \tag{7}$$

here:

$$c_{ip} = \sum_{j,l \neq p} \mu_{jl} K_{\sigma}(\|x_i - x_j\|) \tag{8}$$

Consider a chunk $\{\mu_{ip}\}_{i=1}^N$ and fix $\{\mu_{jk}, k \neq p\}_{j=1}^N$, the objective function (7) is linear in terms of $\{\mu_{ip}\}$ and can be solved, with respect to this chunk, using linear programming. Only one equality constraint $\sum_i \mu_{ip} = 1$ is taken into account in the linear programming problem as the other equality constraints in (6) are independent from the chunk $\{\mu_{ip}\}$.

We can further reduce the size of the chunk $\{\mu_{ip}\}_{i=1}^N$ to only two parameters. Given $i_1, i_2 \in \{1, \dots, N\}$, we can rewrite (7) as:

$$\begin{aligned} & \min_{\mu_{i_1 p}, \mu_{i_2 p}} 2(\mu_{i_1 p} c_{i_1 p} + \mu_{i_2 p} c_{i_2 p}) + b \\ \text{s.t.} \quad & \mu_{i_1 p} + \mu_{i_2 p} = 1 - \sum_{i \neq i_1, i_2} \mu_{ip} \\ & 0 \leq \mu_{i_1 p} \leq 1 \\ & 0 \leq \mu_{i_2 p} \leq 1, \end{aligned} \tag{9}$$

here b is a constant independent from the chunk $\{\mu_{i_1 p}, \mu_{i_2 p}\}$ (see 7). When $c_{i_1 p} = c_{i_2 p}$ the above objective function and the equality constraints are linearly dependent, so we have an infinite set of optimal solutions; any *one* which satisfies the constraints in (9) can be considered. Let us assume $c_{i_1 p} \neq c_{i_2 p}$ and denote $d = 1 - \sum_{i \neq i_1, i_2} \mu_{ip}$, since $\mu_{i_2 p} = d - \mu_{i_1 p}$, the above linear programming problem can be written as:

$$\begin{aligned} & \min_{\mu_{i_1 p}} \mu_{i_1 p} (c_{i_1 p} - c_{i_2 p}) \\ \text{s.t.} \quad & \max(d - 1, 0) \leq \mu_{i_1 p} \leq \min(d, 1) \end{aligned} \tag{10}$$

Depending on the sign of $c_{i_1 p} - c_{i_2 p}$, the solution of (10) is simply taken as one of the bounds of the inequality constraint in (10). At this stage, the parameters $\{\mu_{1p}, \dots, \mu_{Np}\}$ which solve (7) are found by solving iteratively the trivial minimization problem (10) for different chunks, as shown in **Algorithm1**,

Algorithm1(p, μ)

```

Do the following steps ITERMAX1 iterations
Select randomly  $(i_1, i_2) \in \{1, \dots, N\}^2$ .
 $(c_{i_1 p}, c_{i_2 p}) \leftarrow$  (8).
if  $(c_{i_1 p} - c_{i_2 p} < 0)$   $\mu_{i_1 p} \leftarrow \min(d, 1)$  (see 10)
else  $\mu_{i_1 p} \leftarrow \max(d - 1, 0)$ 
 $\mu_{i_2 p} \leftarrow d - \mu_{i_1 p}$  (see 9)
endDo
return  $(\mu_{i_1 p}, \dots, \mu_{i_N p})$ 

```

and the whole QP (6) is solved by looping several times through different cluster indices as shown in **Algorithm2**. For a convex form of the QP (6) and for a large value of ITERMAX1 and ITERMAX2, each subproblem will be convex and the convergence of the decomposition algorithms (1,2) is guaranteed as also discussed for support vector machine training (Platt 1999). In practice ITERMAX1 and ITERMAX2 are set to C and N^2 respectively.

Algorithm2

```

Set  $\{\mu\}$  to random values.
Do the following steps ITERMAX2 iterations
Fix  $p$  in  $\{1, \dots, C\}$ , update  $\mu_{1p}, \dots, \mu_{Np}$  using:
 $(\mu_{1p}, \dots, \mu_{Np}) \leftarrow$  Algorithm1( $p, \mu$ )
endDo

```

4.2 Agglomeration

Initially, the algorithm described above considers an over-estimated number of clusters (denoted C). This will result into several overlapping cluster GMMs which might be detected using a similarity measure for instance the Kullback Leibler divergence (Bishop 1995).

Given two cluster GMMs $p(x|y_k), p(x|y_l)$, for a fixed $\epsilon > 0$ we declare $p(x|y_k), p(x|y_l)$ as similar if:

$$\int_{\mathcal{X}} \|p(x|y_k) - p(x|y_l)\|^2 dx < \epsilon \tag{11}$$

As $p(x|y_k) = \langle \omega_k, \Phi_{\sigma}(x) \rangle$ and $p(x|y_l) = \langle \omega_l, \Phi_{\sigma}(x) \rangle$, we can simply rewrite (11) as:

$$\int_{\mathcal{X}} \langle \omega_k - \omega_l, \Phi_{\sigma}(x) \rangle^2 dx < \epsilon \tag{12}$$

We can set the memberships in ω_k, ω_l such that $\|\omega_k\| = \|\omega_l\| = 1$. Now, from (12) it is clear that two overlapping GMMs correspond to two highly correlated hyperplanes in the mapping space, i.e., $\langle \omega_k, \omega_l \rangle \geq \tau$ (for a fixed threshold $\tau \in [0, 1]$).

Regularization: notice that,

$$\frac{\partial \langle \omega_k, \omega_l \rangle}{\partial \sigma} \geq 0, \tag{13}$$

so the correlation is an increasing function of the scale σ . Assuming $\|\omega_k\| = \|\omega_l\| = 1$, it results that $\forall \tau \in [0, 1], \exists \sigma$ such that $\langle \omega_k, \omega_l \rangle \geq \tau$.

Let us define the following adjacency matrix:

$$A_{k,l} = \begin{cases} 1 & \text{if } \langle \omega_k, \omega_l \rangle \geq \tau \\ 0 & \text{otherwise} \end{cases} \tag{14}$$

This matrix characterizes a graph where a node is a hyperplane and an arc corresponds to two highly correlated hyperplanes (i.e., $\langle \omega_k, \omega_l \rangle \geq \tau$). Now, the actual number of clusters is found as the number of connected components in this graph. For a fixed threshold τ , the scale parameter σ acts as a *regularizer*. When σ is overestimated, the graph $A = \{A_{k,l}\}$ will be connected resulting into only one big cluster whereas a

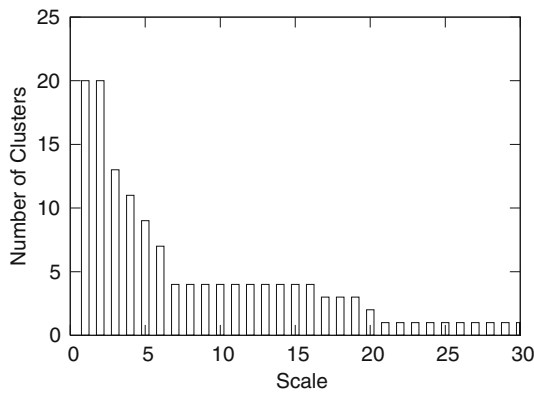


Fig. 3 This figure shows the decrease of the number of clusters with respect to the scale parameter σ ($N = 100, C = 20$). These experiments are performed on the training samples shown in Fig. 4

small σ results into lots of disconnected subgraphs, therefore the total number of clusters will be large (see Fig. 3).

4.3 Mixing different scales

We can extend the formulation presented in (3) to handle clusters with large variation in scale. Let us denote $\sigma(y_l)$ the scale of the GMM related to the cluster y_l . We can rewrite the objective function (6) as:

$$\begin{aligned} \min_{\mu} & \sum_{k,i} \sum_{l \neq k,j} \mu_{ik} \mu_{jl} K_{\sigma(y_l)}(\|x_i - x_j\|) \\ \text{s.t.} & \sum_i \mu_{ic} = 1, \quad \mu_{ic} \in [0, 1], \quad i = 1, \dots, N, \\ & c = 1, \dots, C \end{aligned} \quad (15)$$

Following the same steps (see Sect. 4.1), we can solve (15). However, the above objective function cannot be interpreted

as the minimization of pairwise correlations between hyperplanes since the GMMs, with different scale parameters, correspond to hyperplanes living in different mapping spaces. It follows that, the criteria (in 14) used to detect and merge overlapping GMMs cannot be used. Instead, other criteria, for instance the Kullback Leibler divergence, might be used.

4.4 Toy examples

These experiments are targeted to show the performance brought by this decomposition algorithm both in term of precision and speed. We consider a two dimensional training set, of size N , generated using four Gaussian densities centered at four different locations (see Fig. 4, top-left). We cluster these data using **Algorithm2**, for different values of N (40, 80, 100, 500, 1000). For $N = 100$, Table 1 shows all the pairwise correlations between the hyperplanes found (i) when running **Algorithm2** and (ii) when solving the whole QP (6) using a standard package LOQO (Vanderbei 1999). We can see that each hyperplane found using **Algorithm2** is highly correlated with *only* one hyperplane found using LOQO. Table 2 shows a comparison of the underlying runtime performances using a 1 GHz Pentium III.

5 Database categorization

Experiments have been conducted on both the Olivetti and Columbia databases in order to show the good performance of our clustering method. The Olivetti subset contains 20 persons each one represented by ten faces. The Columbia subset contains 15 categories of objects each one represented by 72 images. Each image from both the Olivetti and Columbia datasets is processed using histogram equalization and

Fig. 4 (Top-left) Data randomly generated from four Gaussian densities centered at four different locations ($N = 1,000$). When running **Algorithm2**, C is fixed to 20, so the total number of parameters in the training problem is 20×1000 . (Other figures) Different clusters are shown with *different colors* and the underlying GMM vectors are shown in *gray*. In these experiments σ is set, respectively, from *top* to *bottom-right* to 10, 4 and 50

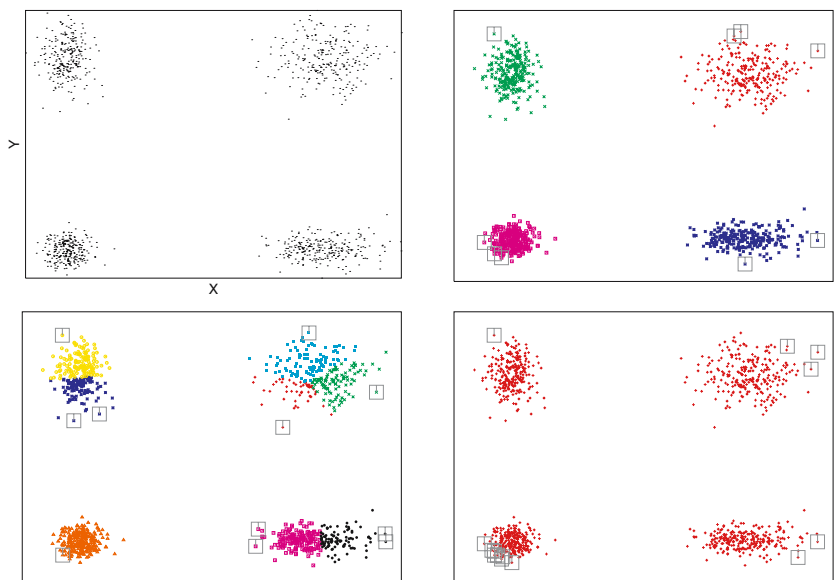


Table 1 This table shows the correlations (normalized inner products) between the hyperplanes found, when (i) running **Algorithm2** and (ii) solving the QP (6) using the LOQO package, on the 2D toy data of Fig. 4 ($\sigma = 10$)

LOQO pack vs. Algorithm2	Class 1	Class 2	Class 3	Class 4
Cluster 1	0.994	0.043	0.331	0.187
Cluster 2	0.036	0.999	0.238	0.133
Cluster 3	0.284	0.313	0.985	0.058
Cluster 4	0.152	0.147	0.057	0.998

Rows characterize the clusters found after running our algorithm while columns characterize the classes (i.e., ground truth)

Table 2 Comparison of the run-time performances on the training samples in Fig. 4 for different values of N ($C = 20, \sigma = 10$)

# $\mathcal{S}(N)$	40	80	100	500
# of parameters in the QP ($N \times C$)	280	480	700	3,500
LOQO pack	24.6 (s)	87.9 (s)	317.8 (s)	> 1 (H)
Algorithm2	0.09 (s)	0.32 (s)	0.56 (s)	24.99 (s)

encoded using 20 coefficients of projection into a subspace learned using ISOMAP (Tanenbaum et al. 2000). Notice that ISOMAP embeds a training database into a subspace such that non-linearly separable classes become more separable, homogeneous in terms of scale, and easier to cluster.

5.1 Scale

It might be easier for database categorization to predict the variance of data rather than the number of clusters, mainly for large databases living in high dimensional spaces. As the number of clusters is initially set to an overestimated value (see Sect. 4.2), our method relies mainly on one parameter, i.e., the scale σ , which makes it possible to find automatically the actual number of clusters. In our experiments, we predict σ by sampling manually few images from some categories, estimating the scales of the underlying classes, then setting σ to the expectation of the scale through these classes. While this setting is not automatic, it has at least the advantage of introducing a priori knowledge on the variance of the data at

Fig. 5 These diagrams show the variation of the scale with respect to the number of samples per class on (left) the Olivetti and (right) Columbia sets. The scale is given as the expectation of the distance between pairs of images taken from 3, 5, 7 and 9 classes

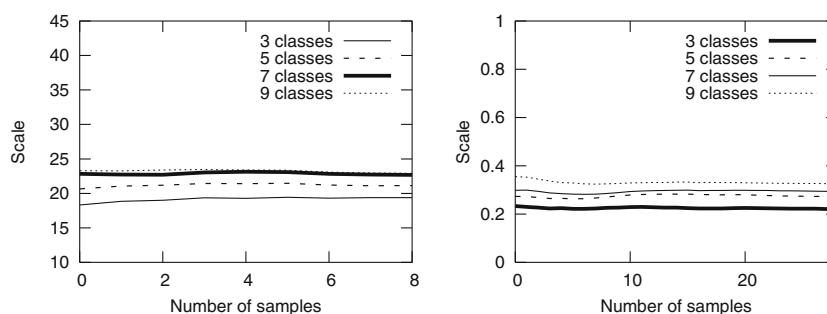


Fig. 6 (Top) Decrease of the number of clusters with respect to the scale parameter. (Bottom) Probability of error with respect to σ . These observations are given using (left) the Olivetti and (right) the Columbia sets

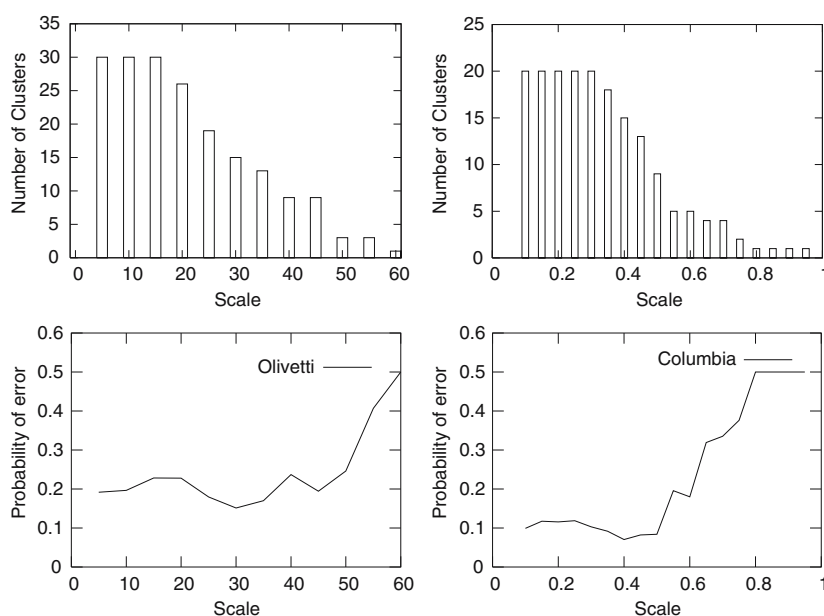


Table 3 This table shows the distribution of the categories (ground truth) through the clusters on the Olivetti set using our clustering algorithm

Categories clusters	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	Total
1	9	9
2	.	10	.	.	1	3	.	.	.	1	3	.	18
3	.	.	10	5	.	.	15
4	.	.	.	10	3	2	.	.	.	15
5	9	1	10
6	5	5
7	5	.	2	7
8	1	.	.	.	9	2	12
9	3	3
10	1	4	5	.	2	12
11	5	5
12	6	3	9
13	2	10	.	1	13
14	1	8	9
15	1	9	10
16	4	1	2	2	8	2	.	.	19
17	3	.	5	8
18	4	7	.	11
19	2	2
Total	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	

We can see that this distribution is concentrated near the diagonal. The scale parameter σ is set to 24 and $C = 30$. Errors in the cardinality of the clusters are mentioned in bold

Table 4 This table shows the distribution of the categories through the clusters on the Columbia set using our clustering algorithm

Categories clusters	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	Total
1	72	72
2	.	72	72
3	.	.	72	72
4	.	.	.	72	72
5	72	72
6	37	.	34	71
7	35	.	38	73
8	72	72
9	72	72
10	72	72
11	72	72
12	72	.	.	.	72
13	30	.	.	30
14	42	72	186
Total	72	72	72	72	72	72	72	72	72	72	72	72	72	72	72	

We can see that this distribution is concentrated near the diagonal. The scale σ is set to 0.4 and $C = 20$. Errors in the cardinality of the clusters are mentioned in bold

the expense of few interactions and reasonable effort from the user.

Figure 5 shows that this heuristic provides “good guess” of the scale on the Olivetti and Columbia sets as it is close to the

optimal scale shown in Fig. 6. Indeed, the estimated scale in Fig. 5 is approximately 24 and 0.4 for, respectively, Olivetti and Columbia sets. For these scales, the number of clusters found by our algorithm on Olivetti and Columbia (resp. 19



Fig. 7 This figure shows 18 face prototypes, from different clusters, found after the application of our method. Each prototype corresponds to a GMM vector



Fig. 8 This figure shows five clusters from the Olivetti database found after the application of our algorithm

and 15, see Fig. 6, top) is close to the actual numbers, resp. 20 and 15. Furthermore, the probability of error (17) is close to its minimum for both Olivetti and Columbia (see Fig. 6, bottom).

5.2 Precision

A Clustering method can be *objectively* evaluated when the ground truth is available, otherwise the meaning of clustering can differ from one intelligent observer to another. The validity criteria are introduced in order to measure the quality of a clustering algorithm, i.e., its capacity to assign data to their actual classes. For a survey on these methods, see for example (Halkidi et al. 2002).

In the presence of a ground truth, we consider in our work a simple validity criteria based on the probability of misclassification. The latter occurs when either two examples belonging to two different classes are assigned to the same cluster, or when two elements belonging to the same class are assigned different clusters. We denote X and Y as two random variables standing respectively for the training examples and their different possible classes $\{y_1, \dots, y_C\}$ and X', Y' respectively similar to X, Y . We denote by $f(X)$ the index of the cluster of X in $\{y_1, \dots, y_C\}$. Formally, we define the misclassification error as:

$$P(\mathbb{1}_{\{f(X) = f(X')\}} \neq \mathbb{1}_{\{Y = Y'\}}) = (*), \tag{16}$$

here:

$$(*) = P(f(X) \neq f(X') \mid Y = Y') P(Y = Y') + P(f(X) = f(X') \mid Y \neq Y') P(Y \neq Y') \tag{17}$$

Again, Fig. 6 (bottom) shows the misclassification error (17) with respect to the scale parameter σ . Tables 3 and 4 are the confusion matrices which show the distribution of 20 and 15 categories from, respectively, the Olivetti and the Columbia sets through the clusters after the application of our clustering method. We can see that this distribution is concentrated in the diagonal and this clearly shows that most of the training data are assigned to their actual categories (see Figs. 7, 8).

5.3 Comparison

Figure 9 shows a comparison of our kernel clustering methods with respect to existing state of the art approaches including fuzzy clustering (Sahbi and Boujemaa 2005), K-means (MacQueen 1965) and hierarchical clustering (Posse 2001). These results are shown for the Olivetti database. The error (17) is measured with respect to the number of clusters. Notice that the latter is fixed for hierarchical clustering and K-means while it corresponds to the resulting number of clusters after setting σ (see Sect. 4.2) for fuzzy clustering and our

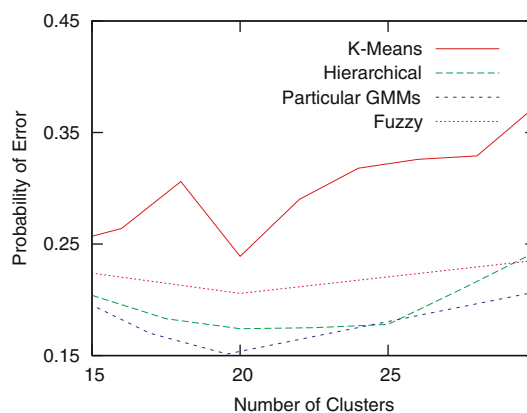


Fig. 9 This figure shows a comparison of the generalization error on the Olivetti set of our method and other existing state of the art methods including fuzzy, k-means and hierarchical clustering

method. We clearly see that the generalization error takes its smallest value when the number of clusters is close to the actual one (i.e., 20).

6 Conclusion and future work

We introduced in this work an original approach for clustering based on a particular Gaussian mixture model (GMM). The method considers an objective function which acts as a regularizer and minimizes the overlap between the cluster GMMs. The GMM parameters are found by solving a quadratic programming problem using a new decomposition algorithm which considers trivial linear programming sub-problems. The actual number of clusters is found by controlling the scale parameter of these GMMs; in practice, it turns out that predicting this parameter is easier than predicting the actual number of clusters mainly for large databases living in high dimensional spaces.

The concept presented in this paper is different from kernel regression which might be assimilated to density estimation while our approach performs this estimation for each cluster. Obviously, the proposed approach performs clustering and density estimation at the same time.

The validity of the method is demonstrated on toy data as well as database categorization problems. As a future work, we will investigate the application of the method to handle noisy data.

References

Bach FR, Jordan MI (2003) Learning spectral clustering. Neural information processing systems
 Ben-Hur A, Horn D, Siegelmann HT, Vapnik V (2000) Support vector clustering. Neural information processing systems, pp 367–373
 Bezdek JC (1981) Pattern recognition with fuzzy objective function algorithms. Plenum Press, New York

- Bishop CM (1995) *Neural networks for pattern recognition*. Clarendon Press, Oxford
- Carson C, Belongie S, Greenspan H, Malik J (2002) Blobworld: image segmentation using expectation-maximization and its application to image querying. *IEEE Trans Pattern Anal Mach Intell* 24(8):1026–1038
- Cristianini N, Shawe-Taylor J (2000) *An introduction to support vector machines*. Cambridge University Press, Cambridge
- Dave RN (1991) Characterization and detection of noise in clustering. In *Pattern Recognit* 12(11):657–664
- Dempster A, Laird N, Rubin D (1977) Maximum likelihood from incomplete data via the em algorithm. *J R Stat Soc B* 39(1):1–38
- Fraley C (1998) Algorithms for model-based gaussian hierarchical clustering. *SIAM J Sci Comput* 20(1):270–281
- Frigui H, Krishnapuram R (1997) Clustering by competitive agglomeration. *Pattern recognition*, vol 30, no. 7
- Frigui H, Krishnapuram R (1999) A robust competitive clustering algorithm with applications in computer vision. *IEEE Trans Pattern Anal Mach Intell* 21(5):450–465
- Halkidi M, Batistakis Y, Vazirgiannis M (2002) Cluster validity methods: Part i and ii. *SIGMOD Record*
- Ichihashai H, Honda K, Tani N (2000) Gaussian mixture pdf approximation and fuzzy c-means clustering with entropy regularization. In: *Proceedings of the 4th Asian fuzzy system symposium*, pp 217–221
- Jain AK, Dubes RC (1988) *Algorithms for clustering data*. Prentice Hall, Englewood Cliffs
- Le-Saux B, Boujemaa N (2002) Unsupervised robust clustering for image database categorization. *IEEE-IAPR International Conference on Pattern Recognition*, pp 259–262
- MacQueen J (1965) Some methods for classification and analysis of multivariate observations. In: *Proceedings of the 5th Berkeley symposium on mathematical statistics and probability*, vol 1
- Orengo CA, Jones DT, Thornton JM (2003) *Bioinformatics—genes, protein and computers*. bios. ISBN: 1-85996-054-5
- Osuna E, Freund R, Girosi F (1997) Training support vector machines: an application to face detection. In: *Proceedings of the international conference on computer vision and pattern recognition*, pp 130–136
- Platt J (1999) Fast training of support vector machines using sequential minimal optimization. In: Schölkopf B, Burges C, Smola AJ (eds) *Advances in kernel methods—support vector learning*. MIT Press, Cambridge, pp 185–208
- Posse C (2001) Hierarchical model-based clustering for large datasets. *J Comput Graph Stat* 10(3):464–486
- Sahbi H, Boujemaa N (2005) Validity of fuzzy clustering using entropy regularization. In: *proceedings of the IEEE conference on fuzzy systems*
- Sneath PH, Sokal RR (1973) *Numerical taxonomy—the principles and practice of numerical classification*. W. H. Freeman, San Francisco
- Tanenbaum J, de Silva V, Langford J (2000) A global geometric framework for non-linear dimensionality reduction. *Science* 290(5500):2319–2323
- Tsao EK, Bezdek JC, Pal NR (1994) Fuzzy kohonen clustering networks. *PR* 27(5):757–764
- Vanderbei RJ (1999) LOQO: an interior point code for quadratic programming. *Optim Methods Softw* 11:451–484
- Wu M, Schoelkopf B (2006) A local learning approach for clustering. *Advances neural information processing systems*
- Yang MS (1993) A survey of fuzzy clustering. *MathCompMod* 18:1–16