

UNIVERSITY OF CAEN NORMANDY

INTERNSHIP REPORT

Hierarchical Clustering of a Mixture Model



UNIVERSITÉ
CAEN
NORMANDIE

student : Abdelmounim
EL-KORCHI

Professor : M. Faicel
CHAMROUKHI
Professor : M. Thien Pham

1^{er} 2021 - May 25

Résumé

[1] this chapter we propose an efficient algorithm for reducing a large mixture of Gaussians into a smaller mixture while still preserving the component structure of the original model; this is achieved by clustering (grouping) the components. The method minimizes a new, easily computed distance measure between two Gaussian mixtures that can be motivated from a suitable stochastic model and the iterations of the algorithm use only the model parameters, avoiding the need for explicit resampling of datapoints. We demonstrate the method by performing hierarchical clustering of handwritten digits.

I introduction

The Gaussian mixture model (MoG) is a flexible and powerful parametric framework for unsupervised data grouping. Mixture models, however, are often involved in other learning processes whose goals extend beyond simple density estimation to hierarchical clustering, grouping of discrete categories or model simplification. In many such situations we need to group the Gaussians components and re-represent each group by a new single Gaussian density. This grouping results in a compact representation of the original mixture of many Gaussians that respects the original component structure in the sense that no original component is split in the reduced representation. We can view the problem of Gaussian component clustering as general data-point clustering with side information that points belonging to the same original Gaussian component should end up in the same final cluster. Several algorithms that perform clustering of data points given such constraints were recently proposed [11, 5, 12]. In this study we extend these approaches to model-based rather than datapoint based settings. Of course, one could always generate data by sampling from the model, enforcing the constraint that any two samples generated by the same mixture component must end up in the same final cluster. We show that if we already have a parametric representation of the constraint via the MoG density, there is no need for an explicit sampling phase to generate representative datapoints and their associated constraints. In other situations we want to collapse a MoG into a mixture of fewer components in order to reduce computation complexity. One example is statistical inference in switching dynamic linear models, where performing exact inference with a MoG prior causes the number of Gaussian components representing the current belief to grow exponentially in time. One common solution to this problem is grouping the Gaussians according to their common history in recent timesteps and collapsing Gaussians grouped together into a single Gaussian [1]. Such a reduction, however, is not based on the parameters of the Gaussians. Other instances in which collapsing MoGs is relevant are variants of particle filtering [10], non-parametric belief propagation [7] and fault detection in dynamical systems [3]. A straight-forward solution for these situations is first to produce samples from the original MoG and then to apply the EM algorithm to learn a reduced model; however this is computationally inefficient and does not preserve the component structure of the original mixture.

II The Clustering Algorithm

We assume that we are given a mixture density f composed of k d -dimensional Gaussian components :

$$f(y) = \sum_{i=1}^k \alpha_i N(y, \mu_i \Sigma_i) = \sum_{i=1}^k \alpha_i f_i(y) \quad (1)$$

where :

$$f_i(y) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{\text{Det}(\Sigma_i)}} \exp\left(-\frac{1}{2} {}^t(y - \mu_i) \Sigma_i^{-1} (y - \mu_i)\right), \quad y \in R^d. \quad (2)$$

$$y = {}^t(y_1, \dots, y_d) \in R^d \quad \mu_i = {}^t(\mu_i^1, \dots, \mu_i^d) \in R^d \quad \Sigma_i = (\sigma_{j,j'}^i, j, j' = 1, \dots, d) \in M(d, d) \quad (3)$$

We want to cluster the components of f into a reduced mixture of $m < k$ components. we denote the set of all (d -dimensional) Gaussian mixture models with at most m components by $\text{MoG}(m)$, one way to formalize the goal of clustering is to say that we wish to find the element g of $\text{MoG}(m)$ “closest” to f under some distance measure

$$\hat{g} = \operatorname{argmin}_g KL(f||g) = \operatorname{argmax}_g \int f \log(g) \quad (4)$$

where KL is kullback-Leibler divergence . This criterion leads to an intractable optimization problem; there is not even a closed-form expression for the KL-divergence between two MoGs let alone an analytic minimizer of its second argument. Furthermore, minimizing a KL-based criterion does not pre- serving the original component structure of f . Instead, we introduce the following new distance measure between but to find g in $\text{MoG}(m)$ which verified (4). This leads to an intractable optimization problem, Instead, we introduce the following new distance measure between

$$f(y) = \sum_{i=1}^k \alpha_i f_i(y) \quad \text{and} \quad g(y) = \sum_{i=1}^k \beta_i g_i(y) \quad (5)$$

$$d(f, g) = \sum_{i=1}^k \alpha_i \min_{j=1}^m KL(f_i || g_j) \quad (6)$$

this distance can be analytically computed. In particular, for example between two Gaussian distributions $N(y, \mu_1 \Sigma_1)$ and $N(y, \mu_2 \Sigma_2)$ which is given by :

$$d(N(y, \mu_1 \Sigma_1), N(y, \mu_2 \Sigma_2)) = \frac{1}{2} \left(\log \left(\frac{|\sigma_2|}{|\sigma_1|} \right) + \operatorname{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1) \Sigma_2^{-1} (\mu_2 - \mu_1) - d \right) \quad (7)$$

Under this distance, the optimal reduced MOG representation \hat{g} is the solution to the minimization of (6) over $\hat{g} = \operatorname{argmin}_g d(f, g)$

we prove that the optimal density \hat{g} is a MoG obtained from grouping the components of f into clusters and collapsing all Gaussians within a cluster into a single Gaussian, to obtain a locally optimal solution.

Denote the set of all the m^k mappings from $\{1, \dots, k\}$ to $\{1, \dots, m\}$ by S . For each $\pi \in S$ et $g \in \text{MoG}(m)$ define :

$$d(f, g, \pi) = \sum_{i=1}^k \alpha_i KL(f_i || g_{\pi(i)}) \quad (8)$$

for $g \in \text{MoG}(m)$, we associate a matching function $\pi^g \in S$:

$$\pi^g(j) = \operatorname{argmin}_{i=1}^m (KL(f_j || g_i)) \quad j = 1, \dots, k \quad (9)$$

One can easily verify that :

$$d(f, g) = d(f, g, \pi^g) = \min_{\pi \in S} d(f, g, \pi) \quad (10)$$

$$\pi^g(j) = \operatorname{argmin}_g \min_{\pi \in S} (d(f, g, \pi)) \quad (11)$$

Unfortunately, we can not solve (11) analytically, Instead we can use alternating minimization to obtain a local minimum.

we define the following functions :

$$f_j^\pi = \frac{\sum_{i \in \pi^{-1}(j)} \alpha_i f_i}{\sum_{i \in \pi^{-1}(j)} \alpha_i} \quad (12)$$

$$\mu'_j = 1/\beta_j \sum_{i \in \pi^{-1}(j)} \alpha_i \mu_i \quad \Sigma'_j = 1/\beta_j \sum_{i \in \pi^{-1}(j)} \alpha_i (\Sigma_i + \|\mu_i - \mu'_j\|^2) \quad (13)$$

where

$$\beta_j = \sum_{i \in \pi^{-1}(j)} \alpha_i \quad (14)$$

let

$$g_j^\pi = N(\mu'_j, \Sigma'_j) = \operatorname{argmin}_g (KL(f_i^\pi \| g)) = \operatorname{argmin}_g (d(f_i^\pi, g)) \quad (15)$$

finally the function we are looking for is :

$$g^\pi = \sum_{j=1}^m \beta_j g_j^\pi \quad (16)$$

$\pi^g = \arg \min_{\pi} d(f, g, \pi) \quad \textbf{(REGROUP)}$
$g^\pi = \arg \min_g d(f, g, \pi) \quad \textbf{(REFIT)}$

III Hierarchical Clustering algorithm of a Mixture Model

...

Algorithm 1: Hierarchical Clustering algorithm of a Mixture Model

Input : k (the number of clusters of F function)

m (the number of clusters of G function)

X (data)

initialization : g^π
 π^g

while $\sum_{i=1}^k \alpha_i KL(f_i \| g_{\pi(i)})$ *no change* **do**

Update : REGROUP

$$\pi^g(i) = \operatorname{argmin}_{j=1}^m (KL(f_i \| g_j)) \quad i = 1, \dots, k$$

Update : REFIT

$$\beta_j = \sum_{i \in \pi^{-1}(j)} \alpha_i \quad j = 1, \dots, m$$

$$\mu'_j = 1/\beta_j \sum_{i \in \pi^{-1}(j)} \alpha_i \mu_i \quad j = 1, \dots, m$$

$$\Sigma'_j = 1/\beta_j \sum_{i \in \pi^{-1}(j)} \alpha_i (\Sigma_i + \|\mu_i - \mu'_j\|^2) \quad j = 1, \dots, m$$

$$g_j^\pi(y) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{\operatorname{Det}(\Sigma'_j)}} \exp \left(-\frac{1}{2} {}^t(y - \mu'_j) \Sigma'^{-1}_j (y - \mu'_j) \right), \quad y \in \mathbb{R}^d, \quad j = 1, \dots, m$$

$$g^\pi = \sum_{j=1}^m \beta_j g_j^\pi$$

end

Result: a set of m cluster

IV performing hierarchical clustering of data from GMM and handwritten digits

In this section, in the first hand we will test the performance of the HCMM algorithm on a dataset that follows the Gaussian mixture model law, we see that the algorithm merge together the Gaussian clusters in a hierarchical way which is the case, in fact the algorithm works in the following way it merge all the Gaussians clusters compared to the distance KL, the precision of the algorithm on a data set follows a Gaussine law and better than on any data set.

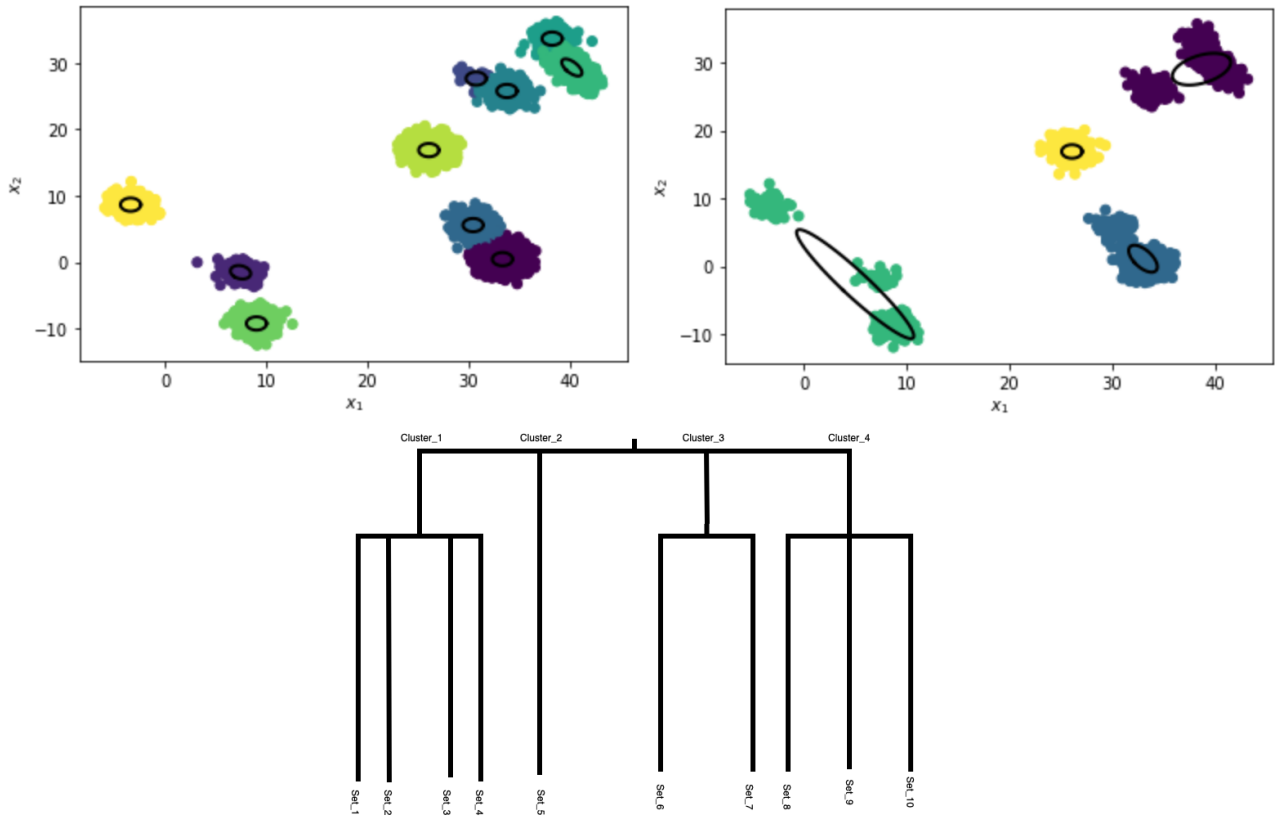


FIGURE 1 – Dendrogram of HCMM for $K = 10$ $m = 4$

IV.I experimental results on data from Gaussian mixture model.

We test the algorithm on different datasets that all follows the Gaussian mixture model law here are the results

TABLE 1 – HCMM on GMM-data

Data	Gaussian mixture models	Hierarchical Clustering of a Mixture Model
The number of mixture components :	$K = 4$	$k = 4$ and $M = 2$
The number of mixture components :	$K = 10$	$K = 10$ and $M = 4$
The number of mixture components :	$K = 150$	$K = 150$ and $M = 3$
The number of mixture components :	$K = 40$	$K = 40$ and $M = 10$

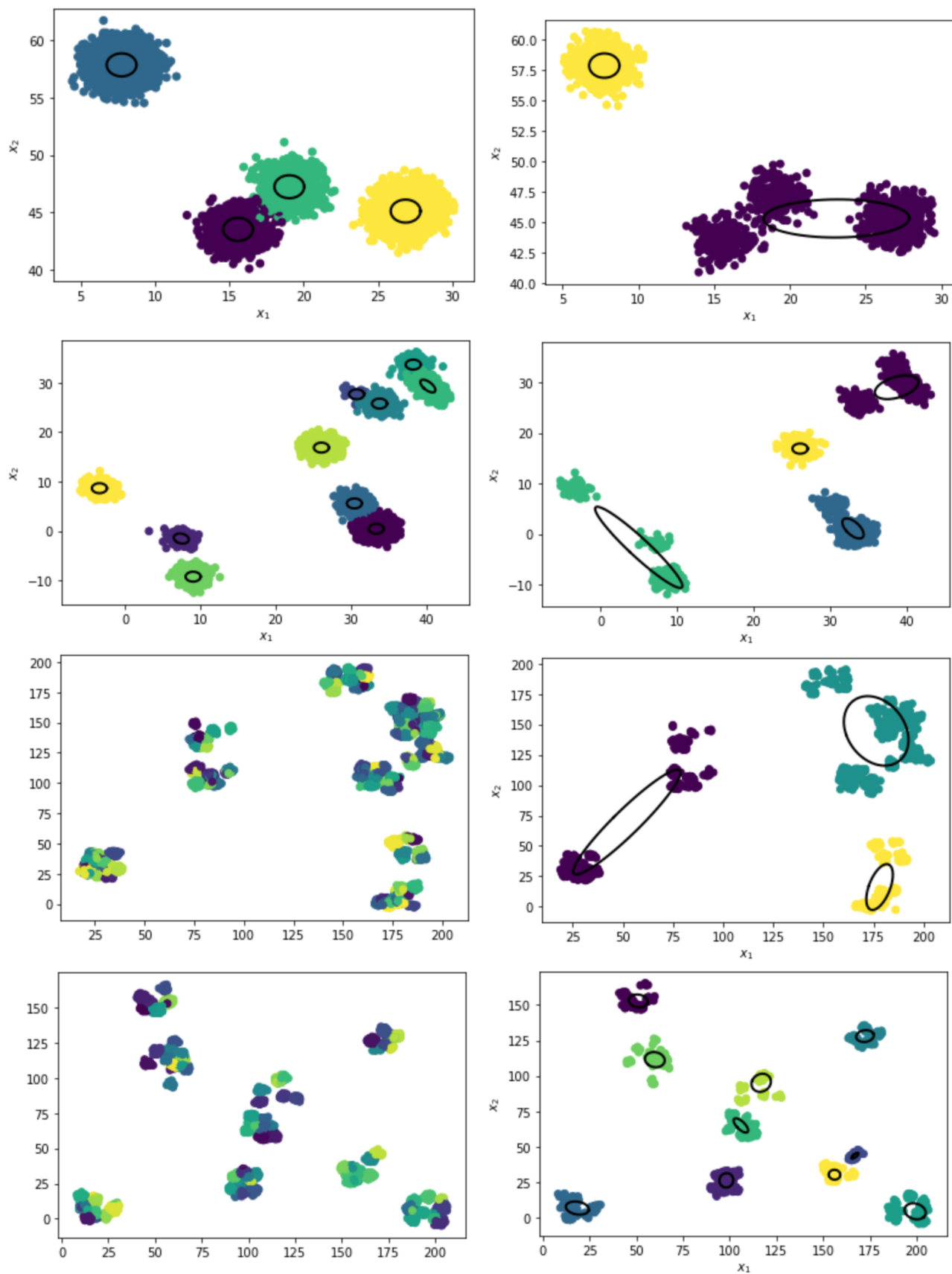


FIGURE 2 – HCMM for GMM data

IV.II experimental results on data handwritten digits.

in the second paragraph we evaluate the performance of our semi-supervised clustering algorithm and compare it to EM clustering approach that does not respect the original component structure. We have applied both methods to clustering handwritten digits images.

In each case, a set of objects is organized in predefined categories. For each category c we learn from a labeled training set a Gaussian distribution $f(x|c)$. A prior distribution over the categories $p(c)$ can be also extracted from the labeled training set. The goal is to cluster the objects into a small number of clusters (fewer than the number of class labels). The EM algorithm approach is to apply an unsupervised clustering to entire collection of original objects, ignoring their class labels. Alternatively we can utilize the given categorization as side-information in order to obtain an improved reduced clustering which also respects the original labels, thus inducing a hierarchical structure. the following image represents an extract of the results of HCMM we see that it separates the figures under two sets :

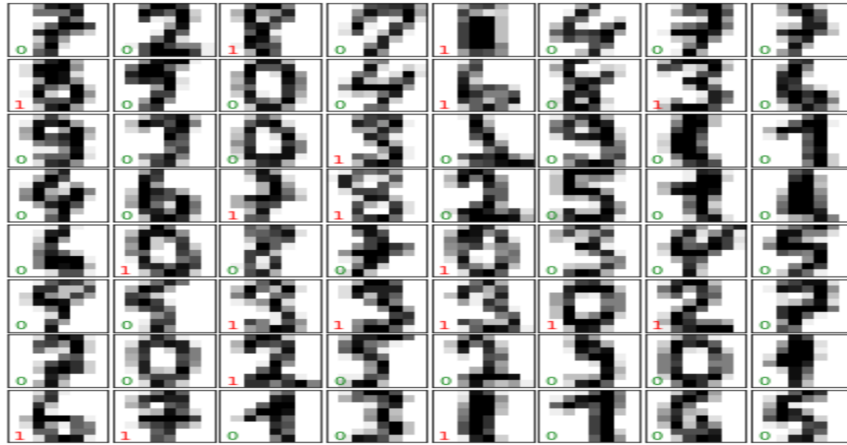


FIGURE 3 – results of HCMM for 64 first images

Our first experiment used a database of handwritten digits. Each example is represented by a 8×8 grayscale pixel. image we use 1437 images for training and 360 images for testing . In the next step we want to divide the digits into two natural clusters, while taking into account their original 10-way structure. We applied our semi-supervised algorithm to reduce the mixture of 10 Gaussians into a mixture of two Gaussians. The minimal distance (2) is obtained when the ten digits are divided into the two groups $[1\ 2\ 3\ 5\ 7\ 8\ 9]$ and $[0, 4, 6]$.

***** Hierarchical_Clustering_Mixture_Model *****		
	cluster = 0	
	unique_elements of cluster	[0 1 3 4 5 7 8 9]
	counts_elements of cluster	[3 62 10 64 39 63 46 10]
	cluster = 1	
	unique_elements of cluster	[0 1 2 3 5 6 8 9]
	counts_elements of cluster	[59 2 62 54 25 64 13 53]
***** Gaussian_Mixture_Model *****		
	cluster = 0	
	unique_elements of cluster	[0 1 3 4 5 7 8 9]
	counts_elements of cluster	[5 55 11 64 42 63 37 11]
	cluster = 1	
	unique_elements of cluster	[0 1 2 3 5 6 8 9]
	counts_elements of cluster	[57 9 62 53 22 64 22 52]
***** the performance of HCMM & GMM *****		
	the purity of assignments provided by HCMM	{0: 3, 1: 2, 3: 10, 5: 25, 8: 13, 9: 10}
	the purity of assignments provided by GM	{0: 5, 1: 9, 3: 11, 5: 22, 8: 22, 9: 11}

FIGURE 4 – results of HCMM for 64 first images

Table 1 : Clustering results showing the purity of a 2-cluster reduced model learned from a training set of handwritten digits in 10 original classes. For each true label, the percentage of cases (from an unseen test set) falling into each of the two reduced classes is shown. The top two lines show the purity of assignments provided by our clustering algorithm; the bottom two lines show assignments from a flat unsupervised fitting of a two component mixture. Table 1 (top) presents, for each digit, the percentage of images that were affiliated with each of the two clusters. Alternatively we can apply a standard EM algorithm to learn by maximum likelihood a flat mixture of 2 Gaussians directly from the 7000 training examples, without utilizing their class labels. Table 1 (bottom) shows the results of such an unsupervised clustering, evaluated on the same test set. Although the likelihood of the unsupervised mixture model was significantly better than the semi-supervised model, both on train and test data-sets it is obvious that the purity of the clusters it learns is much worse since it is not preserving the hierarchical class structure. Comparing the top and bottom of Table 1, we can see that using the side information we obtain a clustering of the digit data-base which is much more correlated with categorization of the set into ten digits than the unsupervised procedure.