

CS 340 Lec. 20: Mixture Models and EM Algorithm

AD

March 2011

Limitations of Clustering using K-Means

- No uncertainty about cluster labels $\{z_i\}_{i=1}^N$.
- Selection of the cost function optimized quite arbitrary.
- What about if the number of clusters K has to be estimated?

- We follow a probabilistic approach where the pdf $p(\mathbf{x})$ of individual data $\{\mathbf{x}_i\}_{i=1}^N$ is modelled explicitly.
- A mixture model states that the pdf of data \mathbf{x}_i is

$$p(\mathbf{x}_i) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i)$$

where $K \geq 2$, $0 \leq \pi_k \leq 1$, $\sum_{k=1}^K \pi_k = 1$ and $\{p_k(\mathbf{x}_i)\}_{k=1}^K$ are pdf.

- You can think of $p_k(\mathbf{x}_i)$ as the pdf of cluster k .

Latent Cluster Labels

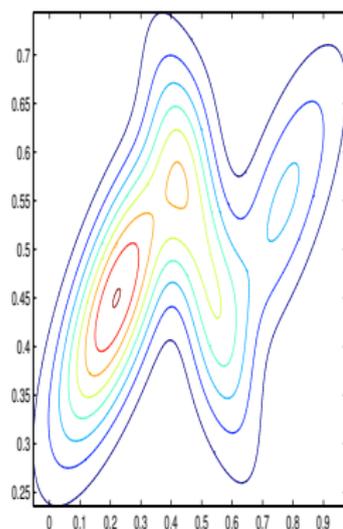
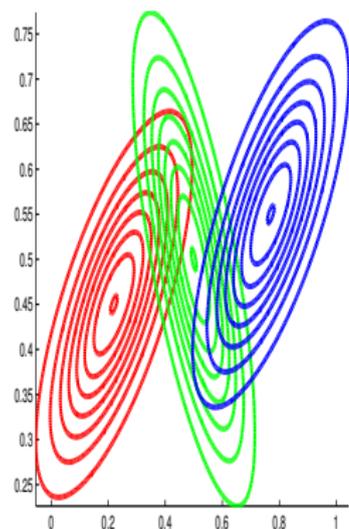
- We associate to each \mathbf{x}_i a cluster label $z_i \in \{1, 2, \dots, K\}$ as in K-means.
- If we set $p(z_i = k) = \pi_k$ then we can rewrite

$$p(\mathbf{x}_i) = \sum_{k=1}^K p(z_i = k) p_k(\mathbf{x}_i)$$

- Alternatively and equivalently, this means that we have now a joint distribution

$$\begin{aligned} p(\mathbf{x}_i, z_i = k) &= p(z_i = k) p(\mathbf{x}_i | z_i = k) \\ &= p(z_i = k) p_k(\mathbf{x}_i) \end{aligned}$$

Example: Mixture of Three 2D-Gaussians



(left) 3 Gaussians in 2D, we display contours of constant proba for each component (center) contours of constant proba of the mixture density (right) Surface plot of the pdf.

- Given \mathbf{x}_i , we can determine

$$\begin{aligned} p(z_i = k | \mathbf{x}_i) &= \frac{p(\mathbf{x}_i, z_i = k)}{\sum_{l=1}^K p(\mathbf{x}_i, z_i = l)} \\ &= \frac{\pi_k p_k(\mathbf{x}_i)}{\sum_{l=1}^K \pi_l p_l(\mathbf{x}_i)}, \end{aligned}$$

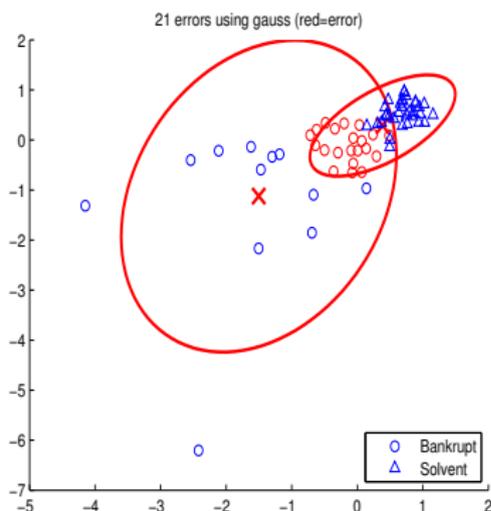
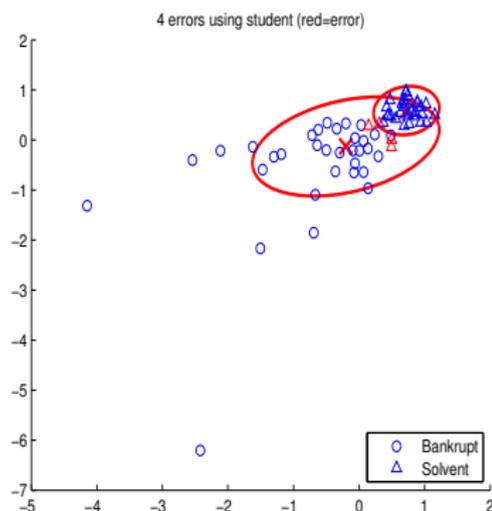
this is sometimes known as soft clustering.

- Assume we can assign data \mathbf{x}_i to a single cluster, then we could set

$$\hat{z}_i = \arg \max_{k \in \{1, 2, \dots, K\}} p(z_i = k | \mathbf{x}_i),$$

this is known as hard clustering.

Example: Mixture of Two 2D-Gaussians and 2D-Students



Mixture models trained on bankruptcy dataset modelled using a mixture of Gaussians (left) and Students (right). Estimated posterior proba is computed. If correct, blue. If incorrect, red.

- Mixture of Gaussians

$$p(\mathbf{x}_i) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Mixture of multivariate Bernoullis: $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,D}) \in \{0, 1\}^D$

$$p(\mathbf{x}_i) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i)$$

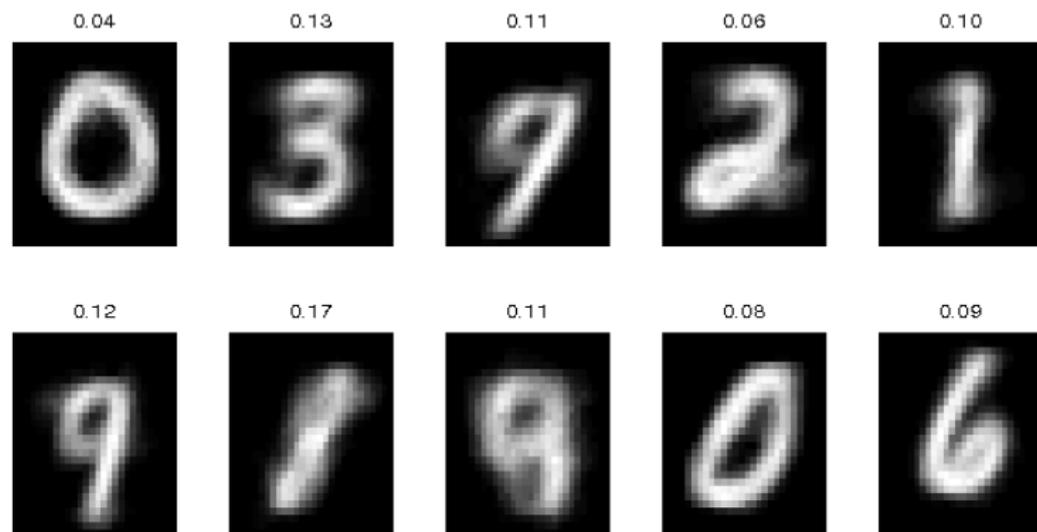
where

$$p_k(\mathbf{x}_i) = \prod_{j=1}^D (\mu_{k,j})^{x_{i,j}} (1 - \mu_{k,j})^{1-x_{i,j}}$$

Mixture of Bernoullis for MNIST Data

- Binary images of digits; $D = 784$.
- We consider applying a mixture of Bernoullis to unlabeled data.
- We set $K = 10$.
- Parameters are learned using Maximum Likelihood (more later!).

Mixture of Bernoullis for MNIST Data



A mixture of 10 multivariate Bernoulli fitted to binarized MNIST data. We display the MLE of cluster means.

- Better models of class conditional distributions for generative classifiers
- Mixture of regressions / Mixture of Experts.
- Applications: astronomy (autoclass), econometrics (mixture of Garch models, SV), genetics, marketing, speech processing.

Maximum Likelihood Parameter Estimation for Mixture Models

- In practice, we typically have

$$p(\mathbf{x}|\theta) = \sum_{k=1}^K \pi_k f(\mathbf{x}; \phi_k)$$

and we need to estimate the parameters $\theta = \{\pi_k, \phi_k\}_{k=1}^K$ given ∞ .

- The ML parameter estimates is given by

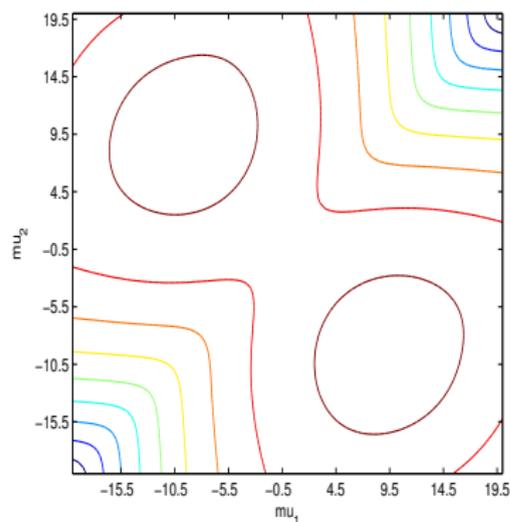
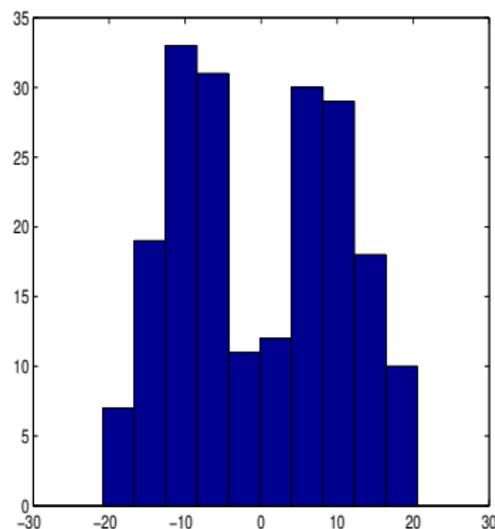
$$\hat{\theta}_{ML} = \arg \max I(\theta)$$

where

$$I(\theta) = \sum_{i=1}^N \log p(\mathbf{x}_i|\theta)$$

- No analytic solution to this problem! Gradient methods could be used but are painful to implement.

Likelihood Surface for a Simple Example



(left) $N = 200$ data points from a mixture of two 2D Gaussians with $\pi_1 = \pi_2 = 0.5$, $\sigma_1 = \sigma_2 = 5$ and $\mu_1 = -\mu_2 = 10$. (right) Log-Likelihood surface $l(\mu_1, \mu_2)$, all the other parameters being assumed known.

- EM is a very popular approach to maximize $l(\theta)$ in this context.
- The key idea is to introduce explicitly the cluster labels.
- If the cluster labels were known then we would estimate θ by maximizing the so-called complete likelihood

$$\begin{aligned}l_c(\theta) &= \sum_{i=1}^N \log p(\mathbf{x}_i, z_i | \theta) \\ &= \sum_{i=1}^N \log \pi_{z_i} f(\mathbf{x}_i; \phi_{z_i})\end{aligned}$$

- We have

$$\begin{aligned}l_c(\theta) &= \sum_{k=1}^K \left(\sum_{i=1:z_i=k}^N \log \pi_{z_i} f(\mathbf{x}_i; \phi_{z_i}) \right) \\ &= \sum_{k=1}^K N_k \log(\pi_k) + \sum_{i=1:z_i=k}^N \log f(\mathbf{x}_i; \phi_k)\end{aligned}$$

where $N_k = \sum_{i=1:z_i=k}^N 1$.

- We would obtain the MLE for the complete likelihood

$$\hat{\pi}_k = \frac{N_k}{N}, \quad \hat{\phi}_k = \arg \max_{\phi_k} \sum_{i=1:z_i=k}^N \log f(\mathbf{x}_i; \phi_k)$$

- Problem: We don't have access to the cluster labels!

Example: Finite mixture of scalar Gaussians

- In this case, $\phi = (\mu, \sigma^2)$

$$f(x; \phi) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

and $\theta = \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K$.

- In this case, the MLE estimate of the complete likelihood is

$$\hat{\pi}_k = \frac{N_k}{N}, \quad \hat{\mu}_k = \frac{1}{N_k} \sum_{i=1:z_i=k}^N x_i,$$

$$\hat{\sigma}_k^2 = \frac{1}{N_k} \sum_{i=1:z_i=k}^N (x_i - \hat{\mu}_k)^2$$

Expectation-Maximization

- EM is an iterative algorithm which generates a sequence of estimates $\{\theta^{(t)}\}$ such that

$$l(\theta^{(t)}) \geq l(\theta^{(t-1)}).$$

- At iteration t , we compute

$$\begin{aligned} Q(\theta, \theta^{(t-1)}) &= \mathbb{E} \left(l_c(\theta) \mid \mathbf{x}_{1:N}, \theta^{(t-1)} \right) \\ &= \sum_{z_{1:N} \in \{1, 2, \dots, K\}^N} \left(\sum_{i=1}^N \log p(\mathbf{x}_i, z_i \mid \theta) \right) p(z_{1:N} \mid \mathbf{x}_{1:N}, \theta^{(t-1)}) \\ &= \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, z_i = k \mid \theta) p(z_i = k \mid \mathbf{x}_i, \theta^{(t-1)}) \end{aligned}$$

and set

$$\theta^{(t)} = \arg \max_{\theta} Q(\theta, \theta^{(t-1)})$$

- We have

$$\begin{aligned} Q(\theta, \theta^{(t-1)}) &= \sum_{i=1}^N \sum_{k=1}^K \log p(\mathbf{x}_i, z_i = k | \theta) p(z_i = k | \mathbf{x}_i, \theta^{(t-1)}) \\ &= \sum_{i=1}^N \sum_{k=1}^K \{\log \pi_k + \log f(\mathbf{x}_i; \phi_k)\} p(z_i = k | \mathbf{x}_i, \theta^{(t-1)}) \\ &= \sum_{k=1}^K \left(\sum_{i=1}^N p(z_i = k | \mathbf{x}_i, \theta^{(t-1)}) \right) \log \pi_k \\ &\quad + \sum_{k=1}^K \left(\sum_{i=1}^N p(z_i = k | \mathbf{x}_i, \theta^{(t-1)}) \log f(\mathbf{x}_i; \phi_k) \right) \end{aligned}$$

- We obtain

$$\hat{\pi}_k^{(t)} = \frac{\sum_{i=1}^N p(z_i = k | \mathbf{x}_i, \theta^{(t-1)})}{N},$$

$$\phi_k^{(t)} = \arg \max_{\phi_k} \sum_{i=1}^N p(z_i = k | \mathbf{x}_i, \theta^{(t-1)}) \log f(\mathbf{x}_i; \phi_k)$$

Example: Finite mixture of scalar Gaussians

- In this case, the EM algorithm iterate

$$\hat{\pi}_k^{(t)} = \frac{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)})}{N}$$

and

$$\hat{\mu}_k^{(t)} = \frac{\sum_{i=1}^N x_i p(z_i = k | x_i, \theta^{(t-1)})}{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)})},$$
$$\hat{\sigma}_k^{2(t)} = \frac{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)}) (x_i - \hat{\mu}_k^{(t)})^2}{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)})}.$$

- We typically iterate the algorithm until $\|\theta^{(t)} - \theta^{(t-1)}\| < \varepsilon$.

Example: Finite mixture of Bernoulli

- Consider now the case where

$$p_k(\mathbf{x}) = \prod_{j=1}^D (\mu_{k,j})^{x_j} (1 - \mu_{k,j})^{1-x_j}$$

so $\theta = \{\pi_k, \mu_{k,1}, \dots, \mu_{k,D}\}_{k=1}^K$.

- In this case, the EM algorithm yields

$$\hat{\pi}_k^{(t)} = \frac{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)})}{N}$$

and

$$\hat{\mu}_{k,j}^{(t)} = \frac{\sum_{i=1}^N x_{i,j} p(z_i = k | x_i, \theta^{(t-1)})}{\sum_{i=1}^N p(z_i = k | x_i, \theta^{(t-1)})}.$$

Proof of Convergence for EM Algorithm

- We want to show that $l(\theta^{(t+1)}) \geq l(\theta^{(t)})$ for $\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$.
- *Proof.* We have

$$p(z_{1:N} | \theta, \mathbf{x}_{1:N}) = \frac{p(\mathbf{x}_{1:N}, z_{1:N} | \theta)}{p(\mathbf{x}_{1:N} | \theta)} \Leftrightarrow p(\mathbf{x}_{1:N} | \theta) = \frac{p(\mathbf{x}_{1:N}, z_{1:N} | \theta)}{p(z_{1:N} | \theta, \mathbf{x}_{1:N})}$$

thus

$$l(\theta) = \log p(\mathbf{x}_{1:N} | \theta) = \log p(\mathbf{x}_{1:N}, z_{1:N} | \theta) - \log p(z_{1:N} | \theta, \mathbf{x}_{1:N})$$

and for any value $\theta^{(t)}$

$$\begin{aligned} l(\theta) &= \underbrace{\sum_{z_{1:N}} \log p(\mathbf{x}_{1:N}, z_{1:N} | \theta) \cdot p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N})}_{=Q(\theta, \theta^{(t)})} \\ &\quad - \sum_{z_{1:N}} \log p(z_{1:N} | \theta, \mathbf{x}_{1:N}) \cdot p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}). \end{aligned}$$

Proof of Convergence for EM Algorithm

- We want to show that $l(\theta^{(t+1)}) \geq l(\theta^{(t)})$ for the EM, so we need to prove that

$$\begin{aligned} & \sum_{z_{1:N}} \log p(z_{1:N} | \theta^{(t+1)}, \mathbf{x}_{1:N}) \cdot p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}) \\ & \leq \sum_{z_{1:N}} \log p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}) \cdot p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}) \end{aligned}$$

- We have

$$\begin{aligned} & \sum_{z_{1:N}} \log \frac{p(z_{1:N} | \theta^{(t+1)}, \mathbf{x}_{1:N})}{p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N})} \cdot p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}) \\ & \leq \log \sum_{z_{1:N}} \frac{p(z_{1:N} | \theta^{(t+1)}, \mathbf{x}_{1:N})}{p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N})} p(z_{1:N} | \theta^{(t)}, \mathbf{x}_{1:N}) \quad (\text{Jensen}) \\ & = \log 1 = 0. \end{aligned}$$

About the EM Algorithm

- Some good things about EM
 - no learning rate (step-size) parameter
 - automatically enforces parameter constraints
 - very fast for low dimensions
 - each iteration guaranteed to improve likelihood
- Some bad things about EM
 - can get stuck in local minima
 - can be slower than conjugate gradient (especially near convergence)
 - requires expensive inference step
 - is a maximum likelihood/MAP method