

Parametric Density Estimation using Gaussian Mixture Models

An Application of the EM Algorithm

Based on tutorials by Jeff A. Bilmes and by Ludwig Schwardt

Pardis Noorzad – Amirkabir University of Technology – Bahman 12, 1389

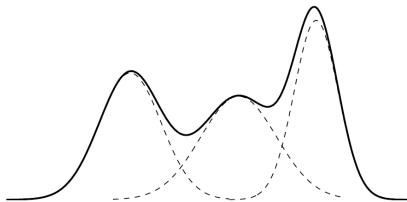
Outline

- 1 Introduction
- 2 Density Estimation
- 3 Gaussian Mixture Model
- 4 Some Results
- 5 Other Applications of EM

What does EM do?

- **iterative** maximization of the likelihood function
 - when data has missing values or
 - when maximization of the likelihood is difficult
 - **data augmentation**
- \mathbf{X} is incomplete data
- $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ is the complete data set

EM and MLE



- **maximum likelihood estimation**
 - estimates density parameters
 - \implies EM app: **parametric density estimation**
 - when density is a **mixture of Gaussians**

The density estimation problem

Definition

Our **parametric density estimation** problem:

$$\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N, \mathbf{x}_i \in \mathbb{R}^D$$

$f(\mathbf{x}_i|\Theta)$, i.i.d. assumption

$$f(\mathbf{x}_i|\Theta) = \sum_{k=1}^K \alpha_k p(\mathbf{x}_i|\theta_k)$$

K components (given)

$$\Theta = (\alpha_1, \dots, \alpha_K, \theta_1, \dots, \theta_K)$$

Constraint on mixing probabilities α_i

Sum to unity

$$\int_{\mathbb{R}^D} p(\mathbf{x}_i|\theta_k)d\mathbf{x}_i = 1$$

$$\begin{aligned} 1 &= \int_{\mathbb{R}^D} f(\mathbf{x}_i|\Theta)d\mathbf{x}_i \\ &= \int_{\mathbb{R}^D} \sum_{k=1}^K \alpha_k p(\mathbf{x}_i|\theta_k)d\mathbf{x}_i \\ &= \sum_{k=1}^K \alpha_k. \end{aligned}$$

The maximum likelihood problem

Definition

Our **MLE** problem:

$$p(\mathbf{X}|\Theta) = \prod_{i=1}^N p(\mathbf{x}_i|\Theta) = \mathcal{L}(\Theta|\mathbf{X})$$

$$\Theta^* = \arg \max_{\Theta} \mathcal{L}(\Theta|\mathbf{X})$$

the log-likelihood: $\log(\mathcal{L}(\Theta|\mathbf{X}))$

$$\log(\mathcal{L}(\Theta|\mathbf{X})) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \alpha_k p(\mathbf{x}_i, \theta_k) \right)$$

difficult to optimize b/c it contains log of sum

The EM formulation

Definition

Our **EM** problem:

$$\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$$

$$\mathcal{L}(\Theta|\mathbf{Z}) = \mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y}) = p(\mathbf{X}, \mathbf{Y}|\Theta)$$

$$y_i \in 1 \dots K$$

$y_i = k$ if the i th sample was generated by the k th component

\mathbf{Y} is a random vector

$$\log(\mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y})) = \sum_{i=1}^N \log(\alpha_{y_i} p(\mathbf{x}_i|\theta_{y_i}))$$

The EM formulation

Continued

Definition

The E- and M-steps with the auxiliary function Q :

$$\text{E-step: } Q(\Theta, \Theta^{(i-1)}) = E[\log(p(\mathbf{X}, \mathbf{Y}|\Theta))|\mathbf{X}, \Theta^{(i-1)}]$$

$$\text{M-step: } \Theta^{(i)} = \arg \max_{\Theta} Q(\Theta, \Theta^{(i-1)})$$

The EM formulation

The Q function

$$\begin{aligned} Q(\Theta, \Theta^{(t)}) &= E[\log(p(\mathbf{X}, \mathbf{Y}|\Theta))|\mathbf{X}, \Theta^{(t)}] \\ &= \sum_{k=1}^K \sum_{i=1}^N \log(\alpha_k p(\mathbf{x}_i|\theta_k)) p(k|\mathbf{x}_i, \Theta^{(t)}) \\ &= \sum_{k=1}^K \sum_{i=1}^N \log(\alpha_k) p(k|\mathbf{x}_i, \Theta^{(t)}) \\ &\quad + \sum_{k=1}^K \sum_{i=1}^N \log(p(\mathbf{x}_i|\theta_k)) p(k|\mathbf{x}_i, \Theta^{(t)}) \end{aligned}$$

The EM formulation

Solving for α_k

We use the Lagrange multiplier to enforce $\sum_k \alpha_k = 1$.

$$\frac{\partial}{\partial \alpha_k} \left[\sum_k \sum_i \log(\alpha_k) p(k|\mathbf{x}_i, \Theta^{(t)}) + \lambda (\sum_k \alpha_k - 1) \right] = 0$$

$$\sum_i \frac{1}{\alpha_k} p(k|\mathbf{x}_i, \Theta^{(t)}) + \lambda = 0$$

Summing both sides over k , we get $\lambda = -N$.

$$\alpha_k = \frac{1}{N} \sum_i p(k|\mathbf{x}_i, \Theta^{(t)}).$$

GMM

Solving for \mathbf{m}_k and Σ_k

Gaussian components

$$p(\mathbf{x}_i | \mathbf{x}_k, \Sigma_k) = (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{m}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mathbf{m}_k)\right)$$

$$\mathbf{m}_k = \frac{\sum_i \mathbf{x}_i p(k | \mathbf{x}_i, \Theta^{(t)})}{\sum_i p(k | \mathbf{x}_i, \Theta^{(t)})}$$

$$\Sigma_k = \frac{p(k | \mathbf{x}_i, \Theta^{(t)}) (\mathbf{x}_i - \mathbf{m}_k) (\mathbf{x}_i - \mathbf{m}_k)^T}{p(k | \mathbf{x}_i, \Theta^{(t)})}$$

Choice of Σ_k

Full covariance

Fits data best, but costly in high-dimensional space.

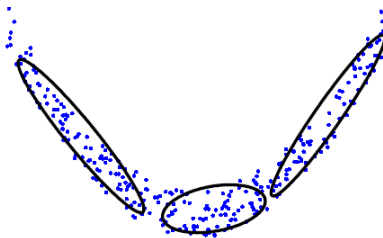


Figure: full covariance

Choice of Σ_k Diagonal covariance

A tradeoff between cost and quality.

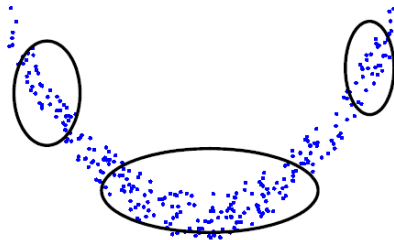


Figure: diagonal covariance

Choice of Σ_k

Spherical covariance, $\Sigma_k = \sigma_k^2 \mathbf{I}$

Needs many components to cover data.

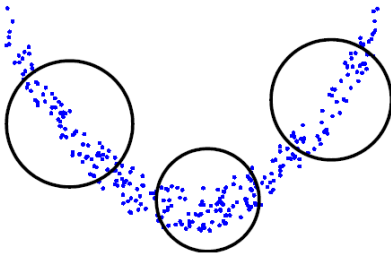


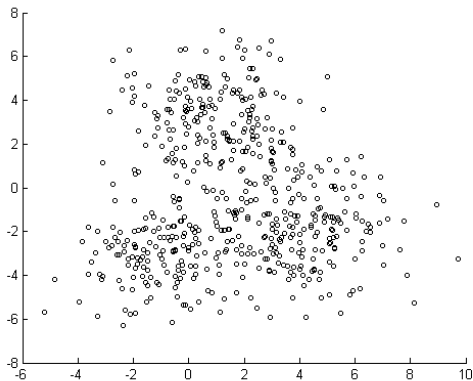
Figure: spherical covariance

Decision should be based on the size of training data set

- so that all parameters may be tuned

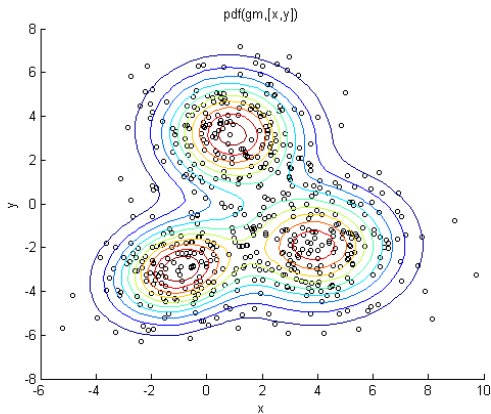
Data

data generated from a mixture of three bivariate Gaussians



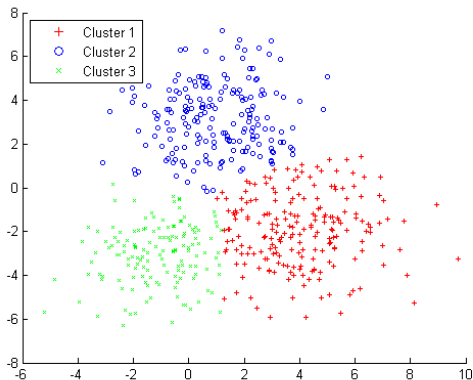
PDFs

estimated pdf contours (after 43 iterations)



Clusters

\mathbf{x}_i assigned to the cluster k corresponding to the highest $p(k|\mathbf{x}_i, \Theta)$



A problem with missing data

Life-testing experiment

- lifetime of light bulbs follows an exponential distribution with mean θ
- two separate experiments
- N bulbs tested until failed
 - failure times recorded as x_1, \dots, x_N
- M bulbs tested
 - failure times not recorded
 - only number of bulbs r , that failed at time t
 - missing data are failure times u_1, \dots, u_M

EM formulation

Life-testing experiment

- $\log(\mathcal{L}(\theta|x, u)) = -N(\log \theta + \bar{x}/\theta) - \sum_i^M (\log \theta + u_i/\theta)$
- expected value for bulb still burning: $t + \theta$
- one that burned out: $\theta - \frac{te^{-t/\theta^{(k)}}}{1-e^{-t/\theta^{(k)}}} = \theta - th^{(k)}$
- E-step:
$$Q(\theta, \theta^{(k)}) = -(N + M) \log \theta - \frac{1}{\theta} (N\bar{x} + (M - r)(t + \theta^{(k)}) + r(\theta^{(k)} - th^{(k)}))$$
- M-step: $\frac{1}{N+M} (N\bar{x} + (M - r)(t + \theta^{(k)}) + r(\theta^{(k)} - th^{(k)}))$

Thank you for your attention. Any questions?