

CIS606 – Lecture 6

Woon Wei Lee, Jacob Crandall
Spring 2014, 10:00am-11:15am,
Mondays and Thursdays

For today:

- Admin stuff: project
- EM algorithm wrap-up
- Tutorial

Gaussian Mixture Models (GMM)

- **Most well known application of EM**
 - Simply another class of mixture model!
 - Recall: MNB model was:

$$\begin{aligned} P(x; \theta) &= \sum_{z=1}^k p(x|z; \theta) p(z; \theta) \\ &= \sum_{z=1}^k \left[\prod_{j=1}^m p(x_j|z) p(z; \theta) \right] \end{aligned}$$

- GMM definition very similar:

$$\begin{aligned} P(x; \theta) &= \sum_{i=1}^k p(x|z; \theta) p(z; \theta) \\ &= \sum_{z=1}^k N(\mu_z, \sigma_z^2) p(z; \theta) \\ &= \sum_{z=1}^k \left[p(z; \theta) \frac{1}{\sqrt{2\pi}^d |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x-\mu_z)^T \Sigma_z^{-1} (x-\mu_z)} \right] \end{aligned}$$

Cont'd

- **E-Step**

$$Q(z(i)|x(i); \theta) = \frac{p(x(i)|z; \theta) p(z(i); \theta)}{p(x(i); \theta)}$$

$$= \frac{N(\mu_z, \sigma_z^2) p(z(i); \theta)}{\sum_{z(i)=1}^k N(\mu_z, \sigma_z^2) p(z(i); \theta)}$$

$$= \frac{p(z(i); \theta) \frac{1}{\sqrt{2\pi}^d |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x(i)-\mu_z)^T \Sigma_z^{-1} (x(i)-\mu_z)}}{\sum_{z=1}^k p(z(i); \theta) \frac{1}{\sqrt{2\pi}^d |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x(i)-\mu_z)^T \Sigma_z^{-1} (x(i)-\mu_z)}}$$

Cont'd

- **M-Step**

- Expected complete data log-likelihood is:

$$\begin{aligned} R(\theta; \hat{\theta}) &= \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \log \left[p(z(i); \theta) \frac{1}{\sqrt{2\pi}^d |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x(i)-\mu_z)^T \Sigma_z^{-1} (x(i)-\mu_z)} \right] \\ &= \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\underbrace{\log p(z(i); \theta)}_{(a)} + \underbrace{\log \frac{1}{\sqrt{2\pi}^d |\Sigma_z|^{1/2}}}_{(b)} - \underbrace{\frac{1}{2}(x(i)-\mu_z)^T \Sigma_z^{-1} (x(i)-\mu_z)}_{(c)} \right] \end{aligned}$$

- Maximizing w.r.t. μ_z (only element (c) involved):

$$\frac{\partial R(\theta; \hat{\theta})}{\partial \mu_z} = \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) \frac{1}{2} \Sigma_z^{-1} (x(i) - \mu_z)$$

$$\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (x(i) - \mu_z) = [0, 0, \dots, 0]^T$$

$$\mu_z^* = \frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) x(i)}{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}$$

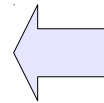
Cont'd

- Maximizing w.r.t. $p(z(i); \theta)$ (only element (a) involved):

$$\frac{\partial R(\theta; \hat{\theta})}{\partial p(z(i); \theta)} = \frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}{p(z(i)|\theta)} - \lambda = 0$$

$$p(z(i); \theta) = \frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}{\lambda}$$

$$p(z(i); \theta) = \frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}{\sum_z \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}$$



(Constrained optimization)

$$\sum_z \left[\frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})}{\lambda} \right] = 1$$

$$\rightarrow \lambda = \sum_z \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})$$

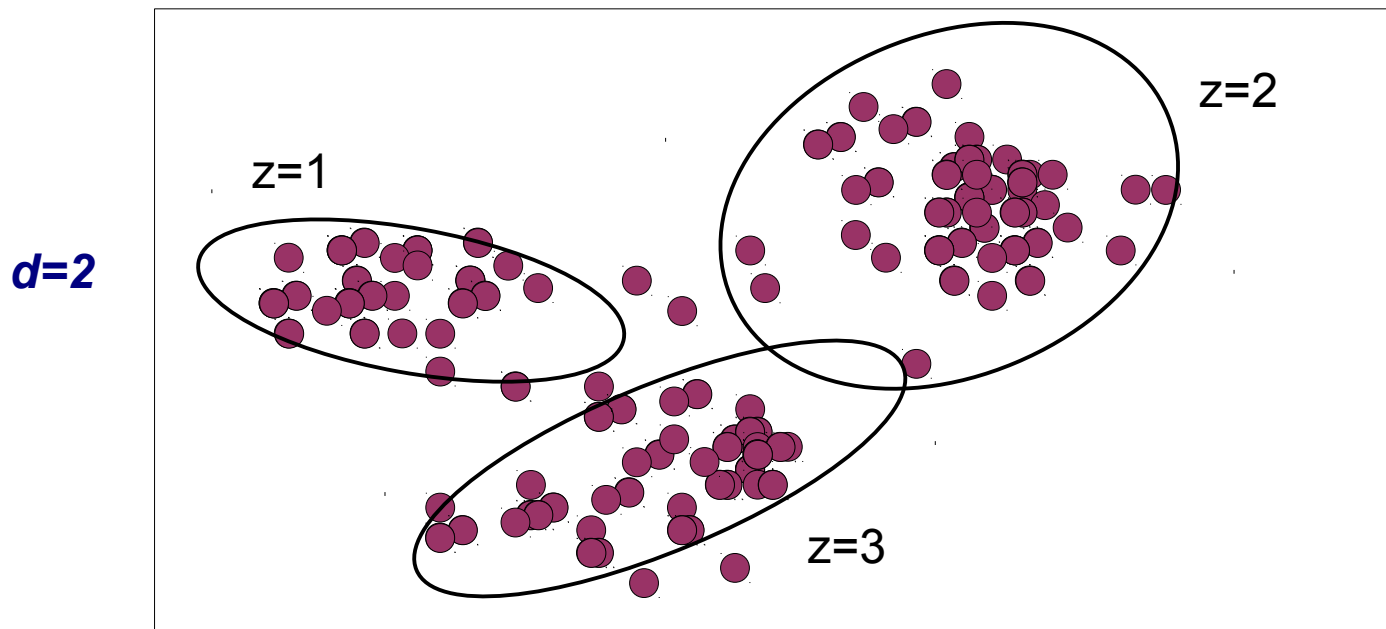
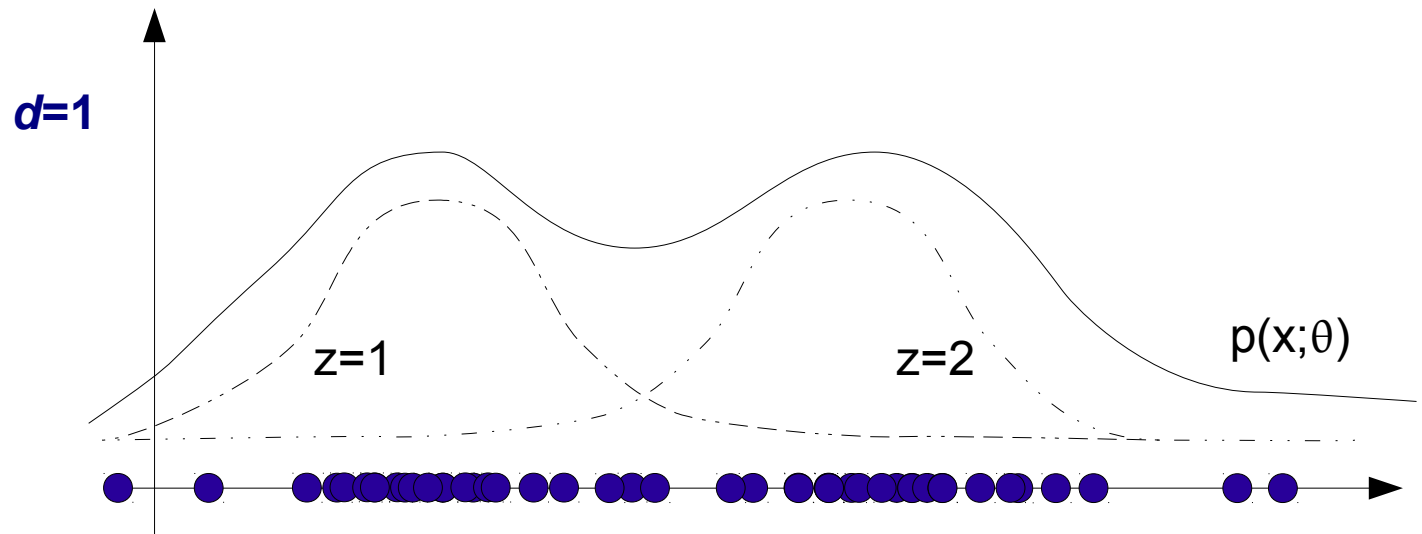
M-Step equations

$$\mu_z^* = \frac{\sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta}) x(i)}{\sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta})}$$

$$p(z(i); \theta)^* = \frac{\sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta})}{\sum_z \sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta})}$$

$$\Sigma_z^* = \frac{\sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta}) [(x(i) - \mu_z)(x(i) - \mu_z)^T]}{\sum_{i=1}^N \mathcal{Q}(z(i)|x(i); \hat{\theta})}$$

GMM as a form of clustering



Potential Problems

- **Unknown number of components!**
 - Similar problem with k-means
 - Measures of clustering quality can be used as before to evaluate candidate “ k ”s
- **Singularities**
 - A more problem occurs when the diversity of one of the Gaussians collapses.
→ e.g.: only a single data point – degenerate!
 - Results in a rank-deficient covariance matrix
 - Determinant goes to zero → not possible to invert

(Appendix – additional stuff)

- Optimize with respect to Σ_z**

- Differentiate $R(\theta; \hat{\theta})$ w.r.t. $\Sigma_z^{-1} \rightarrow$ only elements (b) and (c) involved:

$$\frac{\partial R(\theta; \hat{\theta})}{\partial \Sigma_z^{-1}} = \underbrace{\frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\log \frac{1}{|\Sigma_z|^{1/2}} \right]}{\partial \Sigma_z^{-1}}}_{(b)} - \underbrace{\frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\frac{1}{2} (x(i) - \mu_z)^T \Sigma_z^{-1} (x(i) - \mu_z) \right]}{\partial \Sigma_z^{-1}}}_{(c)}$$

$$= \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \log |\Sigma_z^{-1}|}{\partial \Sigma_z^{-1}} - \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\text{Trace} \left((x(i) - \mu_z)^T \Sigma_z^{-1} (x(i) - \mu_z) \right) \right]}{\partial \Sigma_z^{-1}}$$

***Identity**
 $\text{Trace}(AB) = \text{Trace}(BA)$

$$= \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \log |\Sigma_z^{-1}|}{\partial \Sigma_z^{-1}} - \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\text{Trace} \left(\Sigma_z^{-1} (x(i) - \mu_z) (x(i) - \mu_z)^T \right) \right]}{\partial \Sigma_z^{-1}}$$

Cont'd

From the “Matrix Cookbook”

2.8.2 Symmetric

If \mathbf{A} is symmetric, then $\mathbf{S}^{ij} = \mathbf{J}^{ij} + \mathbf{J}^{ji} - \mathbf{J}^{ij}\mathbf{J}^{ij}$ and therefore

$$\frac{df}{d\mathbf{A}} = \left[\frac{\partial f}{\partial \mathbf{A}} \right] + \left[\frac{\partial f}{\partial \mathbf{A}} \right]^T - \text{diag} \left[\frac{\partial f}{\partial \mathbf{A}} \right] \quad (127)$$

That is, e.g., ([5]):

$$\frac{\partial \text{Tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A} + \mathbf{A}^T - (\mathbf{A} \circ \mathbf{I}), \text{ see (131)} \quad (128)$$

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X})(2\mathbf{X}^{-1} - (\mathbf{X}^{-1} \circ \mathbf{I})) \quad (129)$$

$$\frac{\partial \ln \det(\mathbf{X})}{\partial \mathbf{X}} = 2\mathbf{X}^{-1} - (\mathbf{X}^{-1} \circ \mathbf{I}) \quad (130)$$

• Apply to (c)

• Apply to (b)

$$\begin{aligned} & \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \log |\Sigma_z^{-1}|}{\partial \Sigma_z^{-1}} - \frac{\partial \sum_{i=1}^N \sum_z Q(z(i)|x(i); \hat{\theta}) \left[\text{Trace} \left(\Sigma_z^{-1} (x(i) - \mu_z)(x(i) - \mu_z)^T \right) \right]}{\partial \Sigma_z^{-1}} \\ &= \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (2\Sigma_z - I \circ \Sigma_z) - \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (2S - I \circ S) = 0 \quad \left[\text{where } S = (x(i) - \mu_z)(x(i) - \mu_z)^T \right] \\ & 2 \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (\Sigma_z - S) - I \circ \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (\Sigma_z - S) = 0 \\ & \sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) (\Sigma_z - (x(i) - \mu_z)(x(i) - \mu_z)^T) = 0 \quad \left[2\mathbf{A} - I \circ \mathbf{A} = 0 \rightarrow \mathbf{A} = 0 \right] \\ & \Sigma_z = \frac{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta}) \left[(x(i) - \mu_z)(x(i) - \mu_z)^T \right]}{\sum_{i=1}^N Q(z(i)|x(i); \hat{\theta})} \end{aligned}$$

Recommended Reads

- “A gentle tutorial of the EM algorithm”, Jeff Bilmes
- Andrew Ng's EM algorithm lecture notes
- “Latent Semantic Models for Collaborative Filtering”
Thomas Hoffman
- “The Expectation Maximization Algorithm” - Frank Dellaert
- “The Matrix Cookbook”, Petersen and Pedersen
- “Elements of Statistical Learning”, pg. 272 onwards
- Various Wikis, Google Scholar, etc.