

## **Day 19 - 15 November - Mixtures of Gaussian and EM Algorithms**

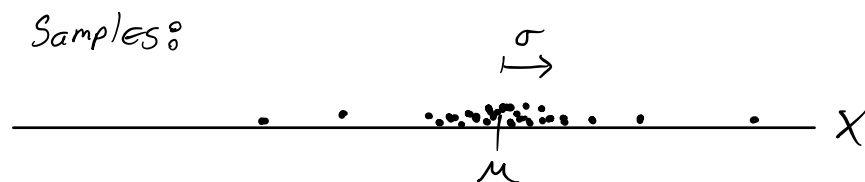
Agenda:

- Multivariate Gaussians
- Maximum Likelihood with Multivariate Gaussians
- Mixtures of Gaussians
- Expectation Maximization (EM) Algorithms

## Multivariate Gaussians

A Gaussian in  $\mathbb{R}$  follows the pdf

$$f(x | \mu, \sigma^2) = \frac{1}{(2\pi)^{1/2}} \frac{1}{(\sigma^2)^{1/2}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$$



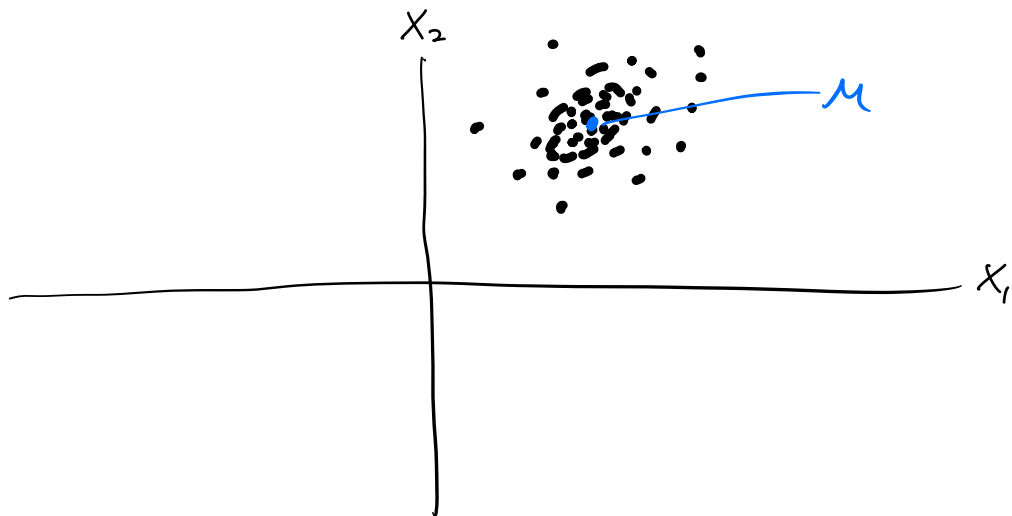
Here  $X \sim \mathcal{N}(\mu, \sigma^2)$

$$\mathbb{E}(X) = \mu$$

$$\mathbb{E}((X-\mu)^2) = \sigma^2$$

multivariate  
A Gaussian in  $\mathbb{R}^d$  follows the pdf

$$f(x | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^t \Sigma^{-1} (x-\mu)\right)$$



$$X \sim \mathcal{N}(\mu, \Sigma)$$

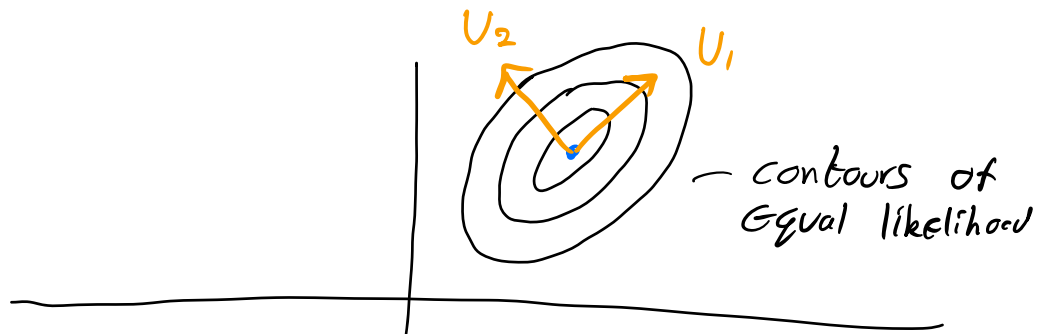
$$\mathbb{E}(X) = \mu \quad - \text{mean}$$

$$\mathbb{E}((X-\mu)(X-\mu)^t) = \Sigma \quad - \text{covariance matrix}$$

Note:  $\Sigma$  is positive semidefinite

$$\begin{aligned} \text{why?} \quad z^t \Sigma z &= z^t \mathbb{E}((X-\mu)(X-\mu)^t) z \\ &= \mathbb{E}[z^t (X-\mu)(X-\mu)^t z] \\ &= \mathbb{E}[\left((X-\mu)^t z\right)^2] \geq 0. \end{aligned}$$

Eigenvectors of  $\Sigma$  with large eigenvalues provide directions w/ most variability



$U_1$  &  $U_2$  are eigenvectors of  $\Sigma$   
w/  $\lambda_1 > \lambda_2$

## Maximum Likelihood Estimation for Gaussians

Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \Sigma)$

Estimate  $\mu$  &  $\Sigma$ .

How? Maximum Likelihood estimation

Likelihood of data  $\underline{X}$ :

$$P(\underline{X} | \mu, \Sigma) = \prod_{i=1}^n P(X_i | \mu, \Sigma)$$

$$\log P(\underline{X} | \mu, \Sigma) = \sum_{i=1}^n \log P(X_i | \mu, \Sigma)$$

$$= -\frac{1}{2} \sum_{i=1}^n (X_i - \mu)^t \Sigma^{-1} (X_i - \mu) + \text{constant}$$

MLE estimate of  $\mu$ :  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$

$$\nabla_{\mu} \log P(\underline{X} | \mu, \Sigma) = -\sum_{i=1}^n \Sigma^{-1} (X_i - \mu) = 0$$

$$\Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n X_i$$

MLE estimate of  $\Sigma$ :

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^n (X_i - \hat{\mu})(X_i - \hat{\mu})^t$$

Issue:  $E[\hat{\Sigma}] = \frac{N-1}{N} \Sigma$  (estimator is biased)

Resolution:

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{i=1}^n (X_i - \hat{\mu})(X_i - \hat{\mu})^t$$

# Mixtures of Gaussians

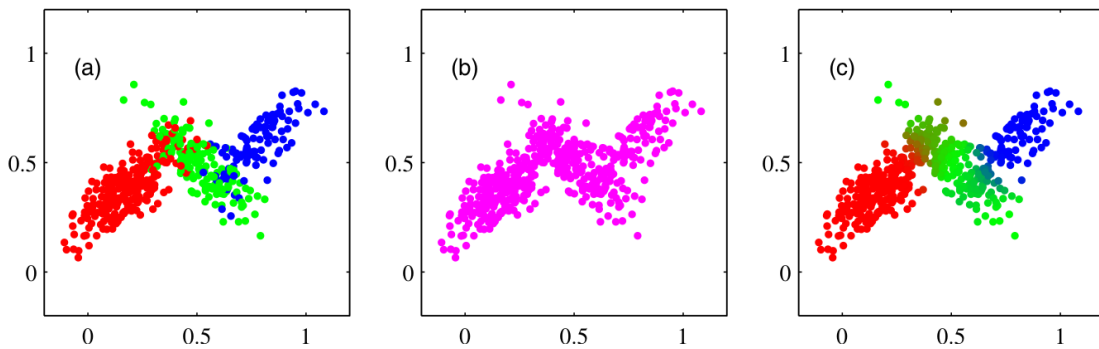
$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

where  $\pi_k \geq 0$  &  $\sum_{k=1}^K \pi_k = 1$ .

To generate a sample:

$$Z \sim \pi \quad (P(Z=k) = \pi_k)$$

$$X \sim \mathcal{N}(x | \mu_z, \Sigma_z)$$



**Figure 9.5** Example of 500 points drawn from the mixture of 3 Gaussians shown in Figure 2.23. (a) Samples from the joint distribution  $p(z)p(x|z)$  in which the three states of  $z$ , corresponding to the three components of the mixture, are depicted in red, green, and blue, and (b) the corresponding samples from the marginal distribution  $p(x)$ , which is obtained by simply ignoring the values of  $z$  and just plotting the  $x$  values. The data set in (a) is said to be *complete*, whereas that in (b) is *incomplete*. (c) The same samples in which the colours represent the value of the responsibilities  $\gamma(z_{n,k})$  associated with data point  $x_n$ , obtained by plotting the corresponding point using proportions of red, blue, and green ink given by  $\gamma(z_{n,k})$  for  $k = 1, 2, 3$ , respectively

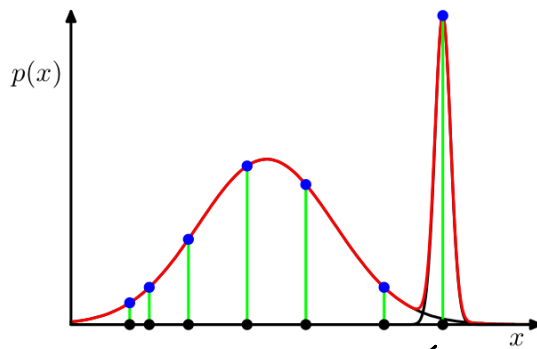
Goal: Use Maximum likelihood estimation to estimate the Gaussian mixture underlying a dataset  $\{X_i\}_{i=1 \dots n}$   $\pi, \mu, \Sigma$

Likelihood of data

$$\log p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \mathcal{N}(X_i | \mu_k, \Sigma_k)$$

Issue: a singularity could arise

**Figure 9.7** Illustration of how singularities in the likelihood function arise with mixtures of Gaussians. This should be compared with the case of a single Gaussian shown in Figure 1.14 for which no singularities arise.



Can put arbitrarily narrow Gaussian around a single data point.

## Expectation - Maximization for Gaussian Mixtures

What conditions should be satisfied at an optimum

$$\log p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \mathcal{N}(X_i | \mu_k, \Sigma_k)$$

Set  $\nabla_{\mu_k} \cdot = 0$ ,



$$0 = \nabla_{\mu_k} \log P(X | \Pi, \mu, \Sigma) = \sum_{i=1}^n \frac{\pi_k \mathcal{N}(X_i | \mu_k, \Sigma_k) \Sigma_k^{-1} (X_i - \mu_k)}{\sum_{j=1}^k \pi_j \mathcal{N}(X_i | \mu_j, \Sigma_j)}$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^n \gamma(Z_{ik}) X_i}{\sum_{i=1}^n \gamma(Z_{ik})} \quad \text{w/} \quad \gamma(Z_{ik}) = \sum_{j=1}^k \pi_j \mathcal{N}(X_i | \mu_j, \Sigma_j)$$

$$\text{Similarly, } \nabla_{\Sigma_k} = 0 \Rightarrow \Sigma_k = \frac{\sum_{i=1}^n \gamma(Z_{ik}) (X_i - \mu_k)(X_i - \mu_k)^t}{\sum_{i=1}^n \gamma(Z_{ik})}$$

$$\text{Finally, } \nabla_{\pi_k} = 0 \Rightarrow \pi_k = \frac{\sum_{i=1}^n \gamma(Z_{ik})}{n}$$

Gives rise to EM algorithm:

1) Initialize  $\mu_k, \Sigma_k, \pi_k$

2) E step

$$\text{update } \gamma(Z_{ik}) = \frac{\pi_k \mathcal{N}(X_i | \mu_k, \Sigma_k)}{\sum_{j=1}^k \pi_j \mathcal{N}(X_i | \mu_j, \Sigma_j)}$$

3) M step

$$\text{update } \mu_k = \frac{1}{N_k} \sum_{i=1}^n \gamma(z_{ik}) X_i$$

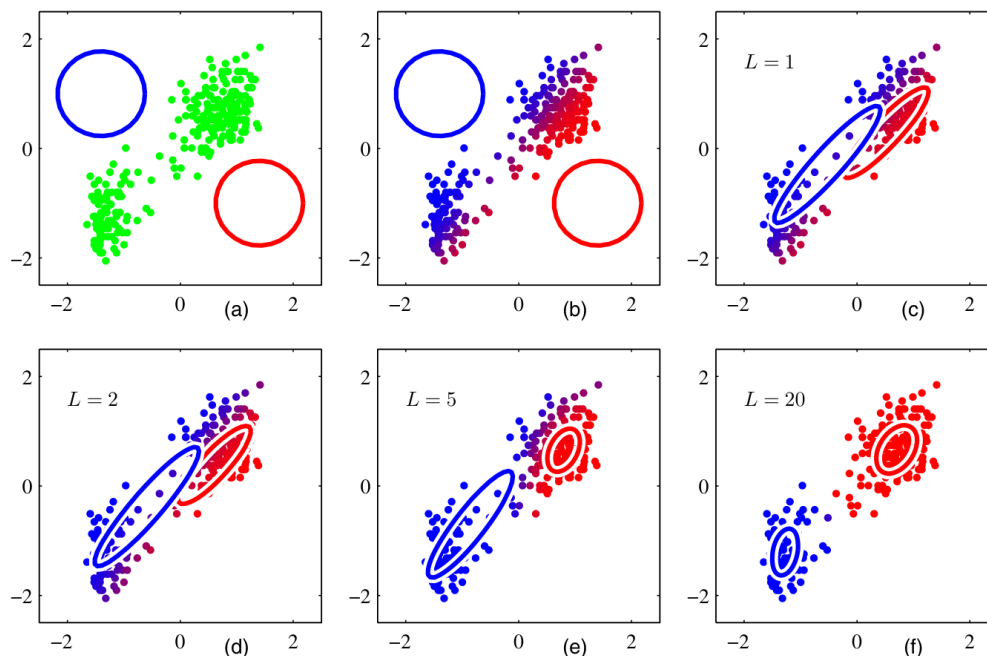
$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^n \gamma(z_{ik}) (X_i - \mu_k)(X_i - \mu_k)^t$$

$$\pi_k = N_k/n$$

$$w/ \quad N_k = \sum_{i=1}^n \gamma(z_{ik})$$

4) Repeat 2&3 until stopping condition

Visualization



**Figure 9.8** Illustration of the EM algorithm using the Old Faithful set as used for the illustration of the  $K$ -means algorithm in Figure 9.1. See the text for details.

## EM and Gaussian Mixtures more abstractly

---

Data  $\{X_i\}_{i=1 \dots n}$

Model w/ 2 Gaussians

$$Z_i = \begin{cases} 1 & \text{w/ prob } \pi_1 \\ 2 & \text{w/ prob } \pi_2 = 1 - \pi_1 \end{cases}$$

$$X_i \sim \mathcal{N}(X_i | \mu_{Z_i}, \Sigma_{Z_i})$$

Given  $\{X_i\}$  estimate  $\{\pi, \mu_1, \mu_2, \Sigma_1, \Sigma_2\} = \theta$

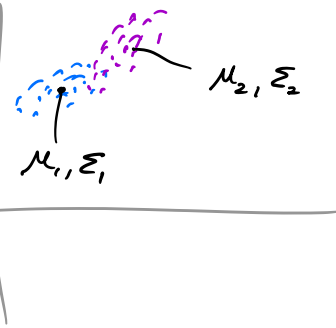
Likelihood of data

$$L(\theta; X, Z) = P(X, Z | \theta)$$

$$= \prod_{i=1}^n \prod_{j=1}^2 \left( \mathcal{N}(X_i | \mu_j, \Sigma_j) \pi_j \right)^{\mathbb{1}_{Z_i=j}}$$

$\underbrace{\hspace{10em}}_{P(X|Z, \theta)} \quad \underbrace{\hspace{10em}}_{P(Z|\theta)}$

Issue: don't know  $Z$ , so keep distribution over all values it could take and update that distribution



(E)

Estimate dist over  $Z$  given  $\theta = \hat{\theta}$

Compute  $P(Z_i = k | X_i, \hat{\theta})$

$$\begin{aligned} P(Z_i = k | X_i, \hat{\theta}) &= \frac{P(X_i | Z_i = k, \hat{\theta}) P(Z_i = k | \hat{\theta})}{P(X_i | \hat{\theta})} \\ &= \frac{\mathcal{N}(X_i | \hat{\mu}_k, \hat{\Sigma}_k) \hat{\pi}_k}{\sum_{j=1}^K \mathcal{N}(X_i | \hat{\mu}_j, \hat{\Sigma}_j) \hat{\pi}_j} \end{aligned}$$

We can rewrite the likelihood function

$$\begin{aligned} Q(\theta, \hat{\theta}) &= \mathbb{E}_{Z | X, \hat{\theta}} \log L(\theta | X, Z) \\ &= \sum_{i=1}^n \mathbb{E}_{Z_i | X_i, \hat{\theta}} \log L(\theta | x_i, z_i) \\ &= \sum_{i=1}^n \sum_{k=1}^K P(Z_i = k | X_i, \hat{\theta}) \log L(\theta | x_i, z_i) \end{aligned}$$

(M)

$$\hat{\theta} \leftarrow \underset{\theta}{\operatorname{argmax}} Q(\theta, \hat{\theta})$$

## The General EM Algorithm

Given a joint distribution  $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$  over observed variables  $\mathbf{X}$  and latent variables  $\mathbf{Z}$ , governed by parameters  $\boldsymbol{\theta}$ , the goal is to maximize the likelihood function  $p(\mathbf{X}|\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ .

1. Choose an initial setting for the parameters  $\boldsymbol{\theta}^{\text{old}}$ .

2. **E step** Evaluate  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ .

3. **M step** Evaluate  $\boldsymbol{\theta}^{\text{new}}$  given by

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (9.32)$$

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.33)$$

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}} \quad (9.34)$$

and return to step 2.