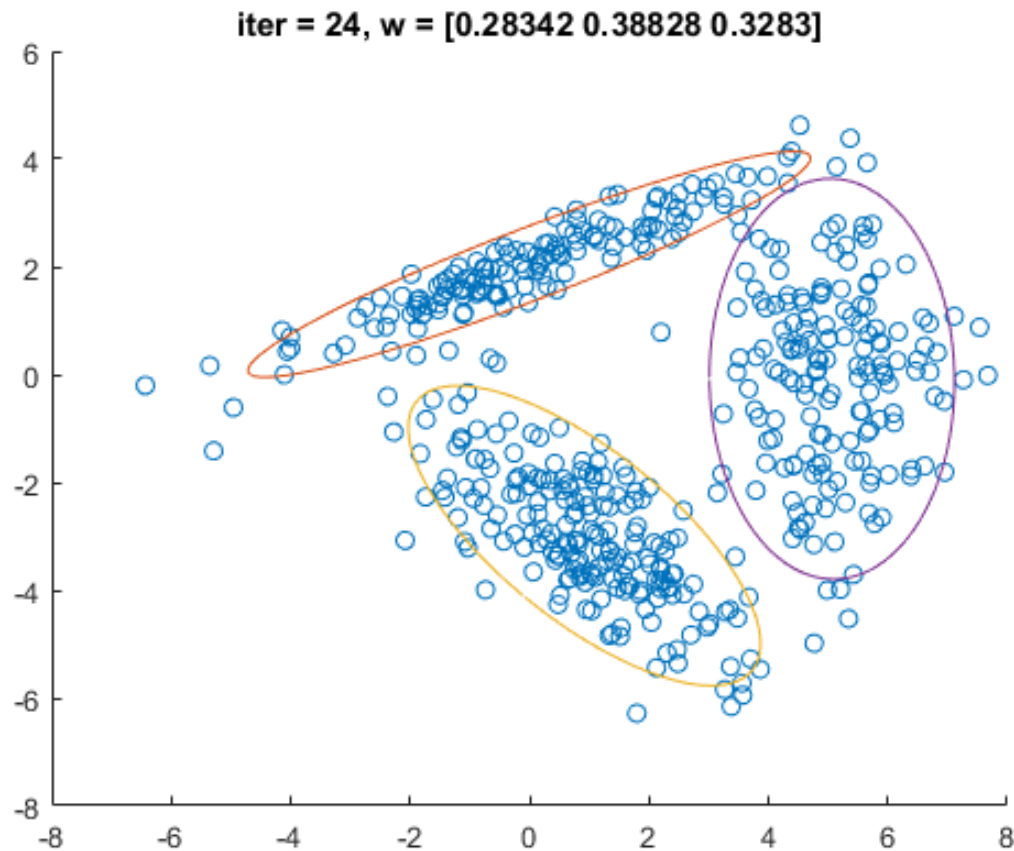


Latent Variable Models; The Expectation-Maximization Algorithm

Outline

- Review of Gaussian Mixture Models
- Latent Variable Models
- The Expectation-Maximization (EM) algorithm

Gaussian Mixture Models



Gaussian Mixture Models

- Let $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} > 0$.
- Recall the multivariate Gaussian density

$$\phi(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

- A random variable \mathbf{X} follows a *Gaussian mixture model* (GMM) with K components if its probability density function f has the form

$$f(\mathbf{x}) = \sum_{k=1}^K w_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where $w_k \geq 0$, $\sum_k w_k = 1$, and for all k , $\boldsymbol{\mu}_k \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_k \in \mathbb{R}^{d \times d}$, $\boldsymbol{\Sigma}_k > 0$.

K known

Maximum Likelihood Estimation

- Observed data: $\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
- The likelihood is

$$\begin{aligned} L(\boldsymbol{\theta}; \underline{\mathbf{x}}) &:= \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\theta}) \\ &= \prod_{i=1}^n \left(\sum_k w_k \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \end{aligned}$$

and the log-likelihood is

$$\ell(\boldsymbol{\theta}; \underline{\mathbf{x}}) := \sum_{i=1}^n \log \left(\sum_{k=1}^K w_k \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right).$$

- When $K > 1$, there is no closed form solution.

Algorithm for Learning GMMs

Initialize $\boldsymbol{\theta}^{(0)}$, $j = 0$

Repeat

E-step:

$$\gamma_{i,k}^{(j)} = \frac{w_k^{(j)} \phi(\mathbf{x}_i; \boldsymbol{\mu}_k^{(j)}, \boldsymbol{\Sigma}_k^{(j)})}{\sum_{\ell=1}^k w_{\ell}^{(j)} \phi(\mathbf{x}_i; \boldsymbol{\mu}_{\ell}^{(j)}, \boldsymbol{\Sigma}_{\ell}^{(j)})}$$

M-step:

$$\begin{aligned}\boldsymbol{\mu}_k^{(j+1)} &= \frac{\sum_i \gamma_{i,k}^{(j)} \mathbf{x}_i}{\sum_i \gamma_{i,k}^{(j)}} \\ \boldsymbol{\Sigma}_k^{(j+1)} &= \frac{\sum_i \gamma_{i,k}^{(j)} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)})^T}{\sum_i \gamma_{i,k}^{(j)}} \\ w_k^{(j+1)} &= \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}.\end{aligned}$$

$$j = j + 1$$

Until convergence criterion satisfied

Today's Main Message

- The previous algorithm is an instance of a more general algorithm called the Expectation-Maximization (EM) algorithm
- The EM algorithm is an iterative algorithm for maximum likelihood estimation for latent variable models
- A latent variable model is a probabilistic model for data, where each observation is explained by one or more latent, or unobserved, variables
- More concretely, in a LVM for an unsupervised learning problem, each observation \mathbf{x}_i is explained by a latent variable \mathbf{z}_i

GMMs are LVMs

- A key to understanding GMMs is to know how to simulate a realization from a known GMM.

- Suppose

$$\boldsymbol{\theta} = (w_1, \dots, w_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$$

is known.

- The following two-step procedure generates a realization of the GMM with parameter vector $\boldsymbol{\theta}$.
 - First, select $k \in \{1, \dots, K\}$ at random, according to the pmf w_1, \dots, w_K .
 - Then draw a realization of $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

$z_i = \text{component from which } x_i \text{ is drawn}$
 $\in \{1, \dots, K\}$

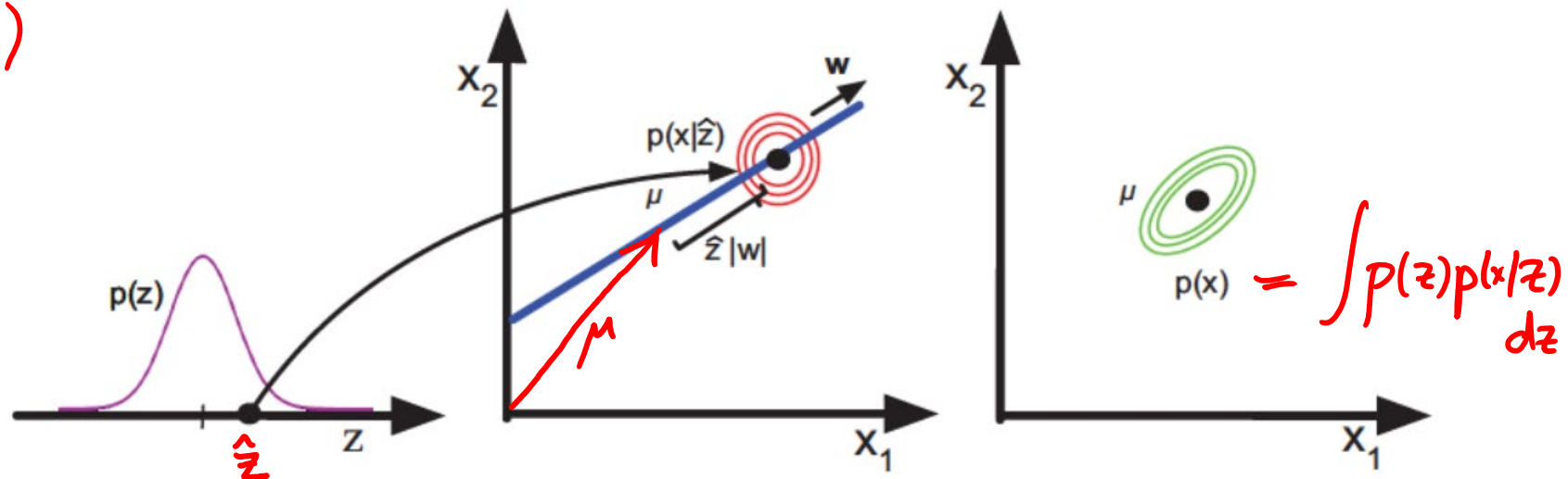
Probabilistic PCA

- Generative model for $\mathbf{x}_1, \dots, \mathbf{x}_n$
- Assume each \mathbf{x}_i is generated as follows:
 - $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), 1 \leq i \leq n$
 - $\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}), 1 \leq i \leq n$

where $\mathbf{z}_i \in \mathbb{R}^k$, $\mathbf{x}_i \in \mathbb{R}^d$, $\mathbf{W} \in \mathbb{R}^{d \times k}$, $\boldsymbol{\mu} \in \mathbb{R}^d$, $\sigma^2 > 0$

$$\mathbf{W} = [\mathbf{w}] \in \mathbb{R}^{2 \times 1}$$

$p(\mathbf{z}|\mathbf{x})$



$$p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) d\mathbf{z}$$

Murphy

$k=1$

$d=2$

Figure 12.1 Illustration of the PPCA generative process, where we have $k=1$ latent dimension generating $d=2$ observed dimensions. Based on Figure 12.9 of (Bishop 2006b).

Latent Variable Models

- We will assume that every observation \mathbf{X}_i is associated to an unobserved (or *hidden* or *latent*) variable \mathbf{Z}_i
- Let $\underline{\mathbf{Z}} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)$ denote all the latent variables.
- Note that the random variables \mathbf{X}_i and \mathbf{Z}_i are jointly distributed

Goal: $\max_{\theta} \underbrace{L(\theta; x_1, \dots, x_n)}_{f(x_1, \dots, x_n; \theta)}$

$$\underline{z} = (z_1, \dots, z_n)$$

Complete Data

- We refer to $(\underline{x}, \underline{z})$ as the *complete data*
- The *complete-data likelihood* is

$$L(\theta; \underline{x}, \underline{z}) = f(\underline{x}, \underline{z}; \theta) = f(\underline{z}; \theta) \cdot f(\underline{x} | \underline{z}; \theta)$$

and the *complete-data log-likelihood* is

$$l(\theta; \underline{x}, \underline{z}) = \log L(\theta; \underline{x}, \underline{z})$$

- The basic idea behind the EM algorithm is to replace $\ell(\theta; \underline{x})$ with

$$\mathbb{E}[l(\theta; \underline{x}, \underline{z}) | \underline{x} = \underline{x}; \theta]$$

↑ randomness due to $\underline{z} | \underline{x} = \underline{x}$

- Since θ is unknown, we use an estimate $\theta^{(j)}$ to compute the expectation, and proceed *iteratively*

want to maximize

EM Algorithm

Initialize $\boldsymbol{\theta}^{(0)}$

$j \leftarrow 0$

Repeat

E-step: Compute

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(j)}) = \mathbb{E}[\ell(\boldsymbol{\theta}; \underline{\mathbf{X}}, \underline{\mathbf{Z}}) \mid \underline{\mathbf{X}} = \underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}]$$

M-step: Solve

$$\boldsymbol{\theta}^{(j+1)} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(j)})$$

$j \leftarrow j + 1$

Until convergence criterion satisfied

Poll

True or False: A reasonable termination criterion for the EM algorithm is to stop iterating when

$$|\ell(\boldsymbol{\theta}^{(j+1)}; \underline{\mathbf{x}}, \underline{\mathbf{z}}) - \ell(\boldsymbol{\theta}^{(j)}; \underline{\mathbf{x}}, \underline{\mathbf{z}})| \leq \epsilon$$

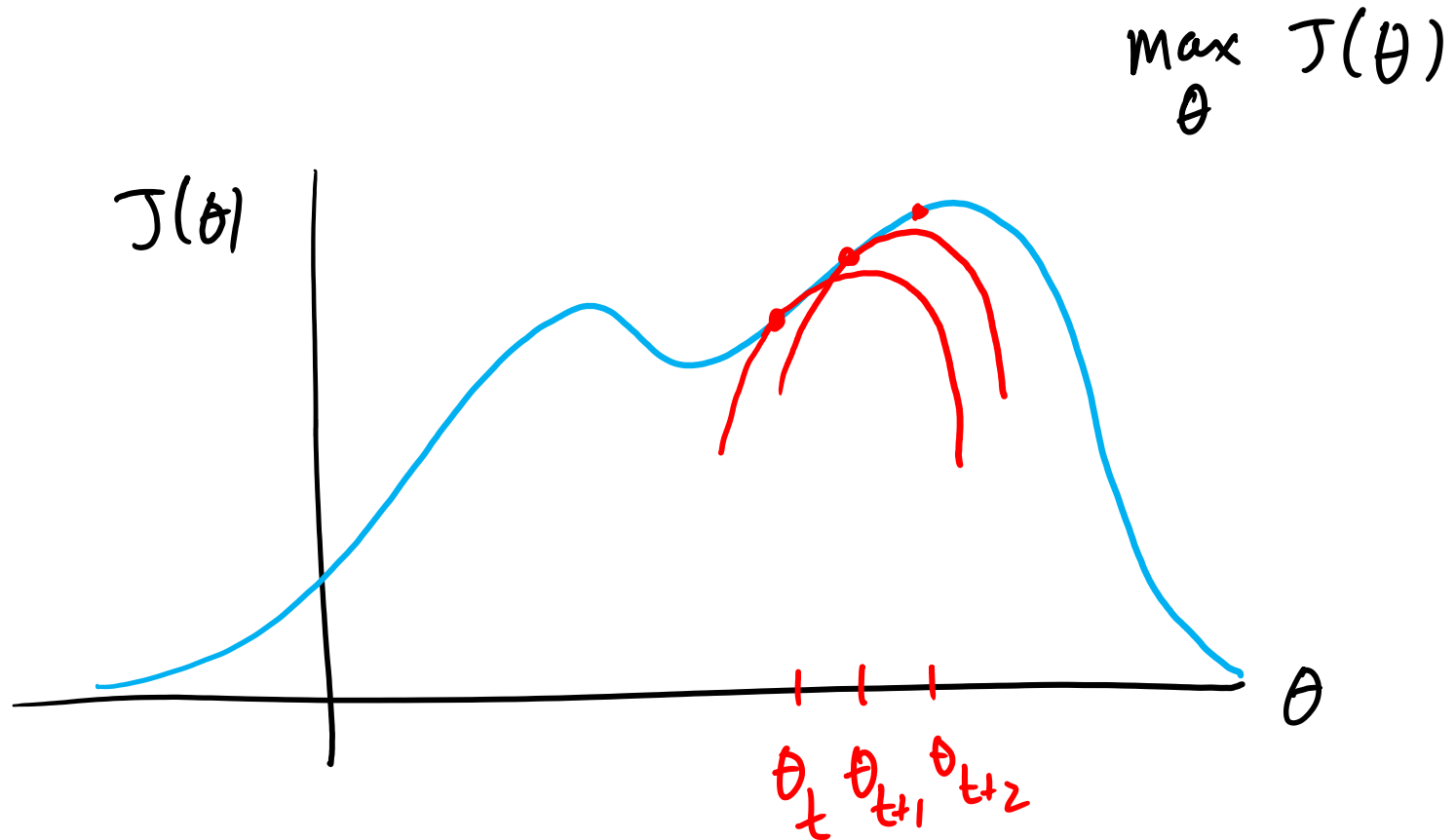
for some small ϵ .

(A) True

(B) False

$$|\ell(\theta^{(j+1)}; \underline{x}) - \ell(\theta^{(j)}; \underline{x})| \leq \epsilon$$

Minorize-Maximize Algorithms



Minorize-Maximize Algorithms

$$\theta_t \leftrightarrow \theta^{(j)}$$

Suppose we wish to maximize the objective function $J(\theta)$

Initialize θ_0

$t \leftarrow 0$

Repeat

Minorize: Find a function $J_t(\theta)$ such that

$$J(\theta_t) = J_t(\theta_t)$$

$$J(\theta) \geq J_t(\theta) \quad \forall \theta$$

Maximize: Solve

$$\theta_{t+1} \leftarrow \operatorname{argmax}_{\theta} J_t(\theta)$$

$t \leftarrow t + 1$

Until convergence

Ascent Property of MM

- MM algorithms never decrease the objective: for all $t \geq 1$

$$J(\boldsymbol{\theta}_{t+1}) \geq J(\boldsymbol{\theta}_t)$$

$$J(\boldsymbol{\theta}_t) = J_t(\boldsymbol{\theta}_t) \leq J_t(\boldsymbol{\theta}_{t+1}) \leq J(\boldsymbol{\theta}_{t+1})$$

def of
minorizing function

$\boldsymbol{\theta}_{t+1}$ maximizes J_t

EM as Minorize-Maximize

- It can be shown that

$$\mathcal{J}_t(\theta) = Q(\theta, \theta_t) + \ell(\theta_t; \underline{x}) - Q(\theta_t, \theta_t)$$

minorizes the log-likelihood $\ell(\theta; \underline{x})$. $= \mathcal{J}(\theta)$

- Hence EM is an MM algorithm

~~• Ascent property: $\ell(\theta_0; \underline{x}) \leq \ell(\theta_1; \underline{x}) \leq \ell(\theta_2; \underline{x}) \leq \dots$~~

$$\begin{aligned} \mathcal{J}_t(\theta_t) &= Q(\cancel{\theta_t}, \theta_t) + \ell(\theta_t; \underline{x}) - Q(\cancel{\theta_t}, \theta_t) \\ &= \ell(\theta_t; \underline{x}) = \mathcal{J}(\theta_t) \end{aligned}$$

EM as Minorize-Maximize

Let p + q be two probability distributions

$$D_{KL}(p \parallel q) = \mathbb{E}_{Y \sim p} \left[\log \left(\frac{p(Y)}{q(Y)} \right) \right] \geq 0$$

↑
Jensen's inequality

$$\begin{aligned} l(\theta; \underline{x}) - l(\theta_t; \underline{x}) &= Q(\theta, \theta_t) - Q(\theta_t, \theta_t) \\ &\quad + D_{KL}(p(z|\underline{x}; \theta_t) \parallel p(z|\underline{x}; \theta)) \\ &\geq Q(\theta, \theta_t) - Q(\theta_t, \theta_t) \end{aligned}$$

Poll

- True or False: The EM algorithm produces a sequence $\boldsymbol{\theta}^{(0)}, \boldsymbol{\theta}^{(1)}, \dots$ satisfying

$$\ell(\boldsymbol{\theta}^{(j)}; \underline{\boldsymbol{x}}) \leq \ell(\boldsymbol{\theta}^{(j+1)}; \underline{\boldsymbol{x}})$$

for all $j \geq 0$.

(A) True

(B) False

EM for GMMs

- Let $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} > 0$.
- Recall the multivariate Gaussian density

$$\phi(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

- A random variable \mathbf{X} follows a *Gaussian mixture model* (GMM) with K components if its probability density function f has the form

$$f(\mathbf{x}) = \sum_{k=1}^K w_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where $w_k \geq 0$, $\sum_k w_k = 1$, and for all k , $\boldsymbol{\mu}_k \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_k \in \mathbb{R}^{d \times d}$, $\boldsymbol{\Sigma}_k > 0$.

K known

$$\theta = (w_1, \dots, w_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K)$$

Complete-Data Log-Likelihood

$$l(\theta; \underline{x}, \underline{z}) = \log (f(\underline{z}; \theta) \cdot f(\underline{x} | \underline{z}; \theta)) \quad \Delta_{i,k} = \begin{cases} 1 & \text{if } z_i = k \\ 0 & \text{if } z_i \neq k \end{cases}$$

$$= \log \left(\prod_{i=1}^n f(z_i; \theta) f(x_i | z_i; \theta) \right)$$

$$= \log \left(\prod_{i=1}^n \Pr(Z_i = z_i; \theta) f(x_i | z_i; \theta) \right)$$

$$= \log \left(\prod_{i=1}^n w_{z_i} \phi(x_i; \mu_{z_i}, \Sigma_{z_i}) \right)$$

$$= \sum_{i=1}^n \log \left(\sum_{k=1}^K \Delta_{i,k} w_k \phi(x_i; \mu_k, \Sigma_k) \right)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \Delta_{i,k} \log (w_k \phi(x_i, \mu_k, \Sigma_k))$$

E-Step for GMMs

- Denote the current iterate $\theta^{(j)} = (w_1^{(j)}, \dots, w_K^{(j)}, \mu_1^{(j)}, \dots, \mu_K^{(j)}, \Sigma_1^{(j)}, \dots, \Sigma_K^{(j)})$
- The E-step amounts to calculating, for all i, k ,

$$\begin{aligned}
 \gamma_{ik}^{(j)} &= \mathbb{E}[\Delta_{i,k} \mid X_i = x_i ; \theta^{(j)}] \\
 &= \Pr\{\Delta_{i,k} = 1 \mid X_i = x_i ; \theta^{(j)}\} \\
 &= \Pr\{Z_i = k \mid X_i = x_i ; \theta^{(j)}\} \\
 &= \frac{\Pr\{Z_i = k ; \theta^{(j)}\} \cdot f(x_i \mid Z_i = k ; \theta^{(j)})}{f(x_i ; \theta^{(j)})} \\
 &= \frac{w_k^{(j)} \phi(x_i ; \mu_k^{(j)}, \Sigma_k^{(j)})}{\sum_{l=1}^K w_l^{(j)} \phi(x_i ; \mu_l^{(j)}, \Sigma_l^{(j)})}
 \end{aligned}$$

M-Step for GMMs

- We need to compute

$$\boldsymbol{\theta}^{(j+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)})$$

where

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)}) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{i,k}^{(j)} \left[\log w_k - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right]$$

Handwritten red note: $E[\Delta_{ik} | \mathbf{x}_i = \mathbf{x}_i; \boldsymbol{\theta}^{(j)}]$

- The solution is

$$w_k^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}$$
$$\boldsymbol{\mu}_k^{(j+1)} = \frac{\sum_i \gamma_{i,k}^{(j)} \mathbf{x}_i}{\sum_i \gamma_{i,k}^{(j)}}$$
$$\boldsymbol{\Sigma}_k^{(j+1)} = \frac{\sum_i \gamma_{i,k}^{(j)} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)})^T}{\sum_i \gamma_{i,k}^{(j)}}$$

Summary: EM Algorithm for GMMs

Initialize $\boldsymbol{\theta}^{(0)}$, $j = 0$

Repeat

E-step:

$$\gamma_{i,k}^{(j)} = \frac{w_k^{(j)} \phi(\mathbf{x}_i; \boldsymbol{\mu}_k^{(j)}, \boldsymbol{\Sigma}_k^{(j)})}{\sum_{\ell=1}^k w_{\ell}^{(j)} \phi(\mathbf{x}_i; \boldsymbol{\mu}_{\ell}^{(j)}, \boldsymbol{\Sigma}_{\ell}^{(j)})}$$

M-step:

$$\begin{aligned}\boldsymbol{\mu}_k^{(j+1)} &= \frac{\sum_i \gamma_{i,k}^{(j)} \mathbf{x}_i}{\sum_i \gamma_{i,k}^{(j)}} \\ \boldsymbol{\Sigma}_k^{(j+1)} &= \frac{\sum_i \gamma_{i,k}^{(j)} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{(j+1)})^T}{\sum_i \gamma_{i,k}^{(j)}} \\ w_k^{(j+1)} &= \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}.\end{aligned}$$

$j = j + 1$

Until convergence criterion satisfied