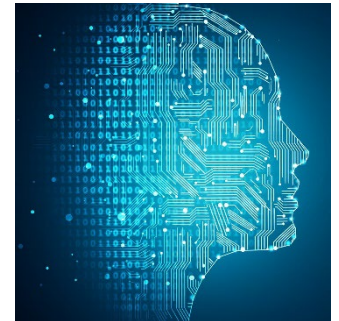


Machine Learning

GMMS and the EM Algorithm



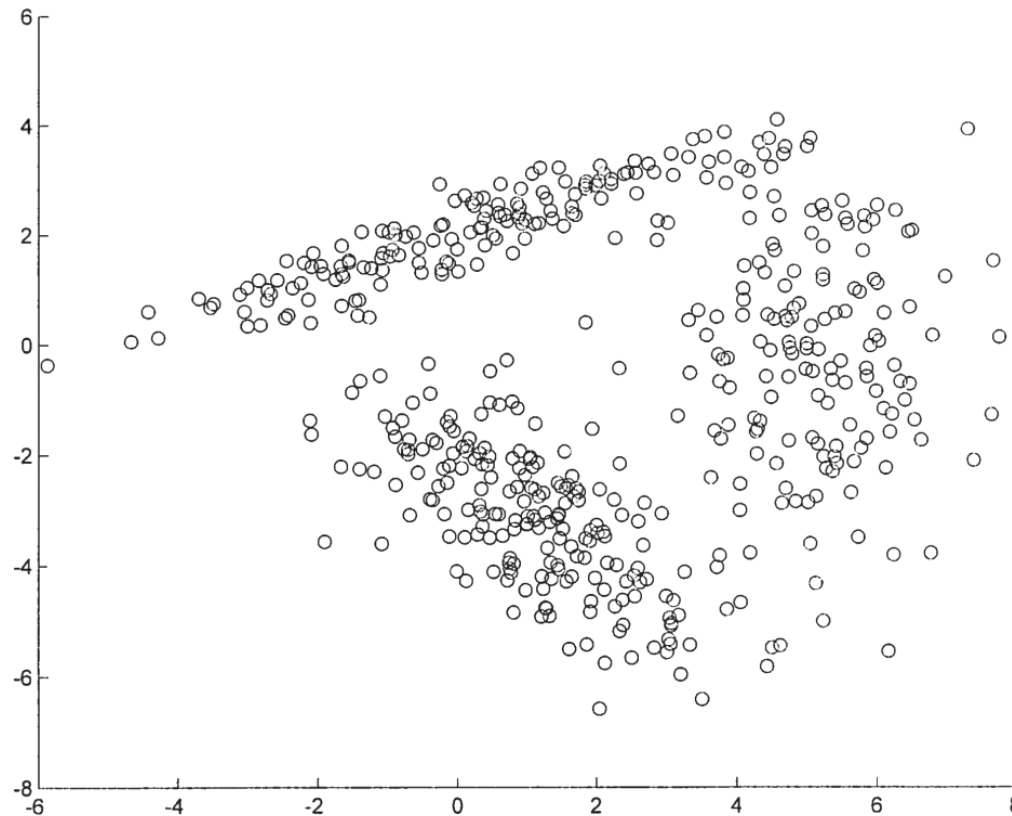
Kevin Moon (kevin.moon@usu.edu)
STAT/CS 5810/6655



Motivating Example



- Suppose we want to cluster the following dataset

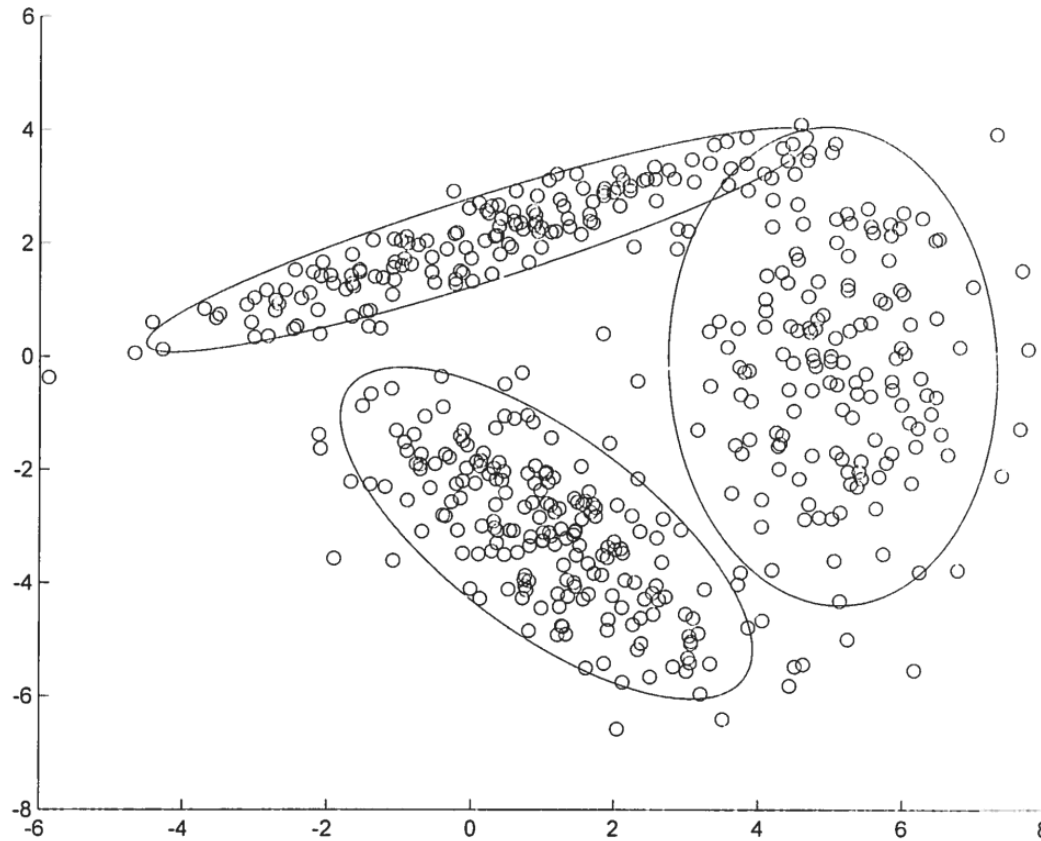


- The data naturally cluster into 3 groups which are each described by bivariate Gaussian densities

Motivating Example



- Below ellipses represent 90% contours of a Gaussian mixture model (GMM) learned using MLE



- Today, we'll learn about GMMs and the EM algorithm, an iterative algorithm for MLE

Outline



1. Gaussian Mixture Models
2. MLE of GMMs
3. The EM algorithm for GMMs
4. The EM algorithm in general

Multivariate Gaussian RV



- Probability density:

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

- $\mathbf{x} \in \mathbb{R}^d$
- $\boldsymbol{\mu} \in \mathbb{R}^d$
- $\Sigma \in \mathbb{R}^{d \times d}, \Sigma \succ 0$ (PD)
- A random variable \mathbf{X} follows a *Gaussian Mixture Model* (GMM) if its probability density function f has the form

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

- $w_k \geq 0, \sum_k w_k = 1$
- $\boldsymbol{\mu}_k \in \mathbb{R}^d, \Sigma_k \in \mathbb{R}^{d \times d}, \Sigma_k \succ 0$

GMM vs. Sum of Gaussians

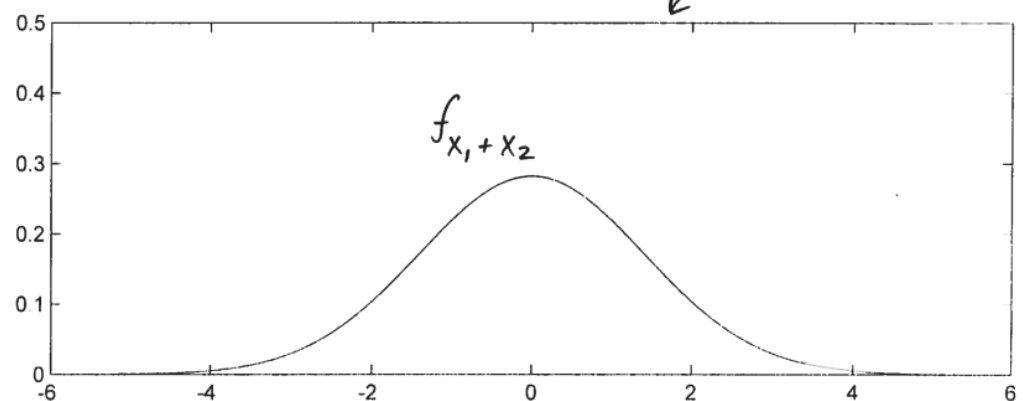
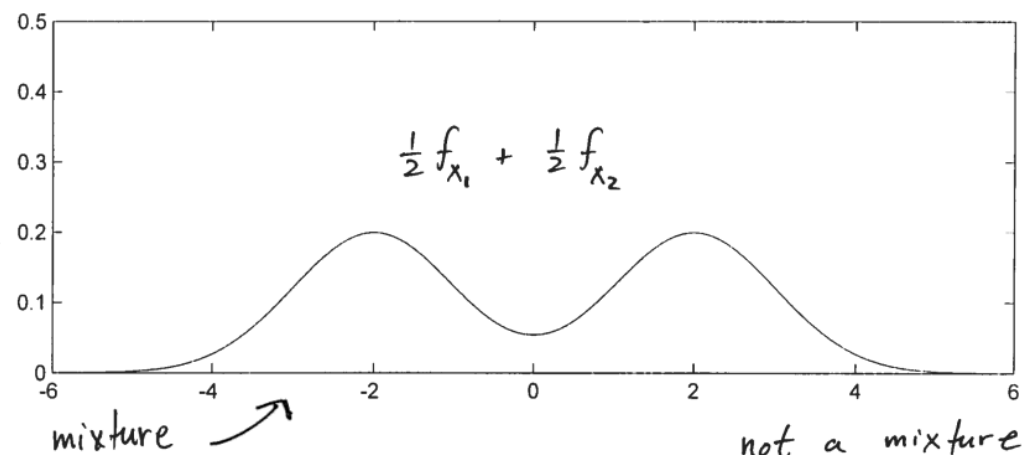
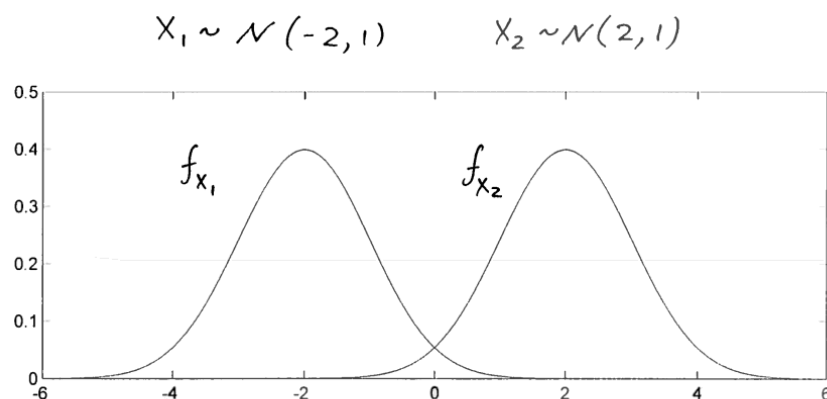


- It can be easy to confuse a *mixture* of Gaussians with a *sum* of Gaussians
- A sum of Gaussian RVs is another Gaussian
 - Unimodal
- A mixture of Gaussians can be multimodal and therefore not Gaussian

GMM vs Sum of Gaussians



Example



$X_1 + X_2 \sim \mathcal{N}(0, 2)$ if X_1 and X_2
are independent

Simulating a GMM



- Suppose the following are known:

$$\boldsymbol{\theta} = (w_1, \dots, w_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \Sigma_1, \dots, \Sigma_K)$$

- How can we simulate a realization of the GMM?
 1. Select a “component” k at random weighted according to the weights w_1, \dots, w_K
 2. Draw a realization $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$
- Why does this work?

Simulating a GMM



- Let $S \in \{1, \dots, K\}$ be a discrete RV such that
$$\Pr(S = k) = w_k$$
- The pdf $f(\mathbf{x})$ of \mathbf{X} is such that for any event A ,
$$\Pr(\mathbf{X} \in A) = \int_A f(\mathbf{x}) d\mathbf{x}$$
- By the law of total expectation:

$$\begin{aligned}\Pr(\mathbf{X} \in A) &= \sum_{k=1}^K \Pr(\mathbf{X} \in A | S = k) \cdot \Pr(S = k) \\ &= \sum_{k=1}^K \left(\int_A \phi(\mathbf{x}; \boldsymbol{\mu}_k, \Sigma_k) d\mathbf{x} \right) w_k \\ &= \int_A \left(\sum_{k=1}^K w_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \Sigma_k) \right) d\mathbf{x}\end{aligned}$$

- Hence $f(\mathbf{x}) = \sum_{k=1}^K w_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \Sigma_k)$ as claimed before

Latent variables



- The variable S is an example of a *state* variable and is said to be *hidden* or *latent* because it is usually unobserved
- We will assume that every realization of a GMM is associated with a hidden state variable

MLE for GMMs



- From observations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, we want to infer the parameters $\boldsymbol{\theta} = (w_1, \dots, w_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \Sigma_1, \dots, \Sigma_K)$
 - I.e. to cluster the data
- We will use maximum likelihood estimation viewing K as fixed
- When $K = 1$, the MLE has a closed-form solution:
 - $\hat{\boldsymbol{\mu}}_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$
 - $\hat{\Sigma}_1 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_1)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_1)^T$
- However, when $K > 1$, there is no closed-form solution

MLE for GMMs



- Denote $\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ for brevity
- Likelihood function:

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

- Log-likelihood:

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

MLE for GMMs



- Since we can't find an analytical solution that maximizes ℓ wrt θ , we'll use an iterative strategy
- This strategy depends critically on the state variables associated with the observations:

$$\underline{s} = (s_1, \dots, s_n)$$

- A natural idea is an alternating algorithm like in k -means:
 - Given θ , update the estimate of \underline{s}
 - Given \underline{s} , update the estimate of θ
- Each step can be performed efficiently
- Learning these parameters can be thought of as a variant of k -means where the cluster assignments are *soft*

The Expectation Maximization (EM) Algorithm



- The *complete data* (observations with labels) is

$$\underline{\mathbf{z}} = (\underline{\mathbf{x}}, \underline{\mathbf{s}})$$

- Define the indicator variable

$$\Delta_{i,k} = \begin{cases} 1 & \text{if } s_i = k \\ 0 & \text{if } s_i \neq k \end{cases}$$

The EM Algorithm



- The *complete-data log-likelihood* is

$$\begin{aligned}\ell(\boldsymbol{\theta}; \underline{\mathbf{x}}, \underline{\mathbf{s}}) &= \log \left(\prod_{i=1}^n \Pr(S_i = s_i; \theta) f(x_i | s_i; \theta) \right) \\ &= \log \left(\prod_{i=1}^n w_{s_i} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_{s_i}, \Sigma_{s_i}) \right) \\ &= \sum_{i=1}^n \log \left(w_{s_i} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_{s_i}, \Sigma_{s_i}) \right) \\ &= \sum_{i=1}^n \log \left(\sum_{k=1}^K \Delta_{i,k} w_k \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k) \right) \\ &= \sum_{i=1}^n \sum_{k=1}^K \Delta_{i,k} [\log w_k + \log \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k)]\end{aligned}$$

Only term that depends on \mathbf{s}

The EM Algorithm



- Ideally, if we knew $S_i/\Delta_{i,k}$, we could maximize the complete data log-likelihood
- Since we don't know these state variables, we can replace $\Delta_{i,k}$ with its expected value and then maximize wrt θ
- However, the expected value of $\Delta_{i,k}$ depends on θ
 - A “chicken and egg” problem
 - Suggests an iterative approach
- The EM algorithm is such an approach
 - Produces a sequence of estimates $\theta^{(1)}, \theta^{(2)}, \dots$
 - Alternates between two steps: the E step and the M step

The E-step



- Calculate the expected complete-data log-likelihood:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)}) = \mathbb{E}_{\underline{\mathbf{S}}|\underline{\mathbf{x}}}[\ell(\boldsymbol{\theta}; \underline{\mathbf{x}}, \underline{\mathbf{s}})|\underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}]$$

Conditional expectation wrt $\underline{\mathbf{S}}|\underline{\mathbf{x}}$. \mathbf{S} is capitalized here since it is viewed as random.

The pmf of $\underline{\mathbf{S}}|\underline{\mathbf{x}}$ depends on the GMM parameters; use the current estimate

- This becomes

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)}) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{i,k}^{(j)} [\log w_k + \log \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k)]$$

Let's find an expression for this

The E-Step



$$\begin{aligned}\gamma_{i,k}^{(j)} &= \mathbb{E}[\Delta_{i,k} | \underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}] \\ &= \Pr(\Delta_{i,k} = 1 | \underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}) \\ &= \Pr(S_i = k | \underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}) \\ &= \frac{\Pr(S_i = k; \boldsymbol{\theta}^{(j)}) f(\mathbf{x}_i | S_i = k; \boldsymbol{\theta}^{(j)})}{f(\mathbf{x}_i | \boldsymbol{\theta}^{(j)})} \\ &= \frac{w_k^{(j)} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_k^{(j)}, \Sigma_k^{(j)})}{\sum_{\ell=1}^K w_{\ell}^{(j)} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_{\ell}^{(j)}, \Sigma_{\ell}^{(j)})}\end{aligned}$$

- This is the fraction of the density value at \mathbf{x}_i explained by the k th component
 - Sometimes called the *responsibility* of cluster k for \mathbf{x}_i
 - A soft measure of cluster membership

The M-Step



- Compute

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

- Recall

$$\begin{aligned} & Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)}) \\ &= \sum_{i=1}^n \sum_{k=1}^K \gamma_{i,k}^{(j)} \left[\log w_k - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_k| \right. \\ & \quad \left. - \frac{1}{2} (x_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (x_i - \boldsymbol{\mu}_k) \right] \end{aligned}$$

The M-Step



- It can be shown without too much trouble that the solution is

$$\boldsymbol{\mu}_k^{(j+1)} = \frac{\sum_{i=1}^n \gamma_{i,k}^{(j)} \mathbf{x}_i}{\sum_{i=1}^n \gamma_{i,k}^{(j)}} \quad \leftarrow \text{Weighted sample mean and covariance}$$

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

$$w_k^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}$$

Fraction of all data explained by k th component

Terminating the EM Algorithm



- The algorithm can be terminated when the increase in the likelihood is small:

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

- ϵ is a small number
- We'll see that the likelihood function is increasing with j

Initialization



- The EM algorithm is sensitive to initialization (like k -means)
 - I.e. it is not convex
- One possibility:
 - $\mu_k^{(0)}$ = random x_i (distinct for each k)
 - $\Sigma_k^{(0)}$ = sample covariance of all data (same for all k)
 - $w_k^{(0)} = \frac{1}{K}$
 - As with k -means, it may be beneficial to run the algorithm many times and pick the result with the best likelihood
- Another possibility: initialize with k -means result

Connection to k -means



- k -means is a special case of the EM algorithm with GMM

- Consider the GMM $f(\mathbf{x}) = \sum_{k=1}^K w_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \sigma^2 I)$
 - σ^2 is fixed

Common, isotropic
covariance

- The EM algorithm iterates over:

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

$$w_k^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}$$

$$\gamma_{i,k}^{(j+1)} = \frac{w_k^{(j)} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_k^{(j)}, \sigma^2 I)}{\sum_{\ell=1}^K w_{\ell}^{(j)} \cdot \phi(\mathbf{x}_i; \boldsymbol{\mu}_{\ell}^{(j)}, \sigma^2 I)}$$

Connection to k -means



- Now as $\sigma^2 \rightarrow 0$, we have

$$\gamma_{i,k} \rightarrow \begin{cases} 1 & \text{if } k = \arg \min_{\ell} \|\mathbf{x}_i - \boldsymbol{\mu}_{\ell}\| \\ 0 & \text{otherwise} \end{cases}$$

- This is equivalent to the k -means algorithm

The EM Algorithm in General



- The EM algorithm applies to many other MLE problems where unobserved variables would make computation easier
- As before, denote
 - $\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
 - $\underline{\mathbf{s}} = (s_1, \dots, s_n)$
 - s_i is some variable that explains how \mathbf{x}_i was generated
- Let $\ell(\boldsymbol{\theta}; \underline{\mathbf{x}})$ denote the log-likelihood and $\ell(\boldsymbol{\theta}; \underline{\mathbf{x}}, \underline{\mathbf{s}})$ the complete-data log-likelihood

The EM Algorithm in General



The general EM algorithm:

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

- Repeat

E-Step: Form

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(j)}) = \mathbb{E}_{\underline{\mathbf{s}}|\underline{\mathbf{x}}}[\ell(\boldsymbol{\theta}; \underline{\mathbf{x}}, \underline{\mathbf{s}})|\underline{\mathbf{x}}; \boldsymbol{\theta}^{(j)}]$$

M-Step: Compute

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

- Until termination criterion satisfied

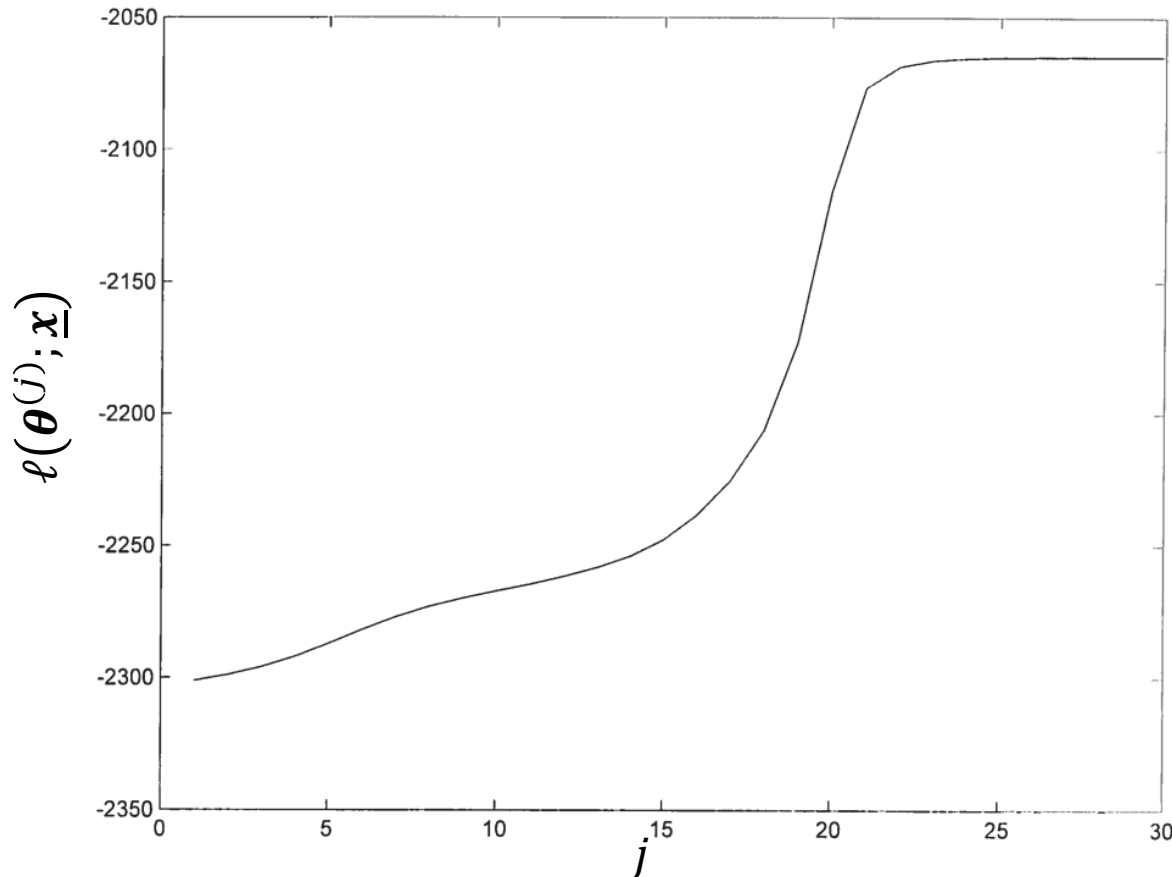
Ascent property of the EM Algorithm



- **Theorem:** For each $j = 0, 1, 2, \dots$

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

- 3 Gaussians example from beginning:



Final Thoughts/Summary



- The EM algorithm is a general algorithm for MLE whenever there may be unobserved (latent) or missing variables
- Initialization matters (generally not a convex problem)
- Convergence can sometimes be slow

Further reading



- A tutorial:
<https://ieeexplore.ieee.org/abstract/document/543975>
- Proof of the final theorem:
http://web.eecs.umich.edu/~cscott/past_courses/eecs545f16/20_em_gmm.pdf
- ESL Section 8.5