

# The EM Algorithm

Héctor Corrada Bravo

University of Maryland, College Park, USA

CMSC 644: 2019-03-13

# Soft K-means Clustering

Instead of the combinatorial approach of the  $K$ -means algorithm, take a more direct probabilistic approach to modeling distribution  $Pr(X)$ .

Assume each of the  $K$  clusters corresponds to a multivariate distribution  $Pr_k(X)$ ,

$Pr(X)$  is given by *mixture* of these distributions as  $Pr(X) = \sum_{k=1}^K \pi_k Pr_k(X)$ .

# Soft K-means Clustering

Specifically, take  $Pr_k(X)$  as a multivariate normal distribution  $f_k(X) = N(\mu_k, \sigma_k^2 I)$

and mixture density  $f(X) = \sum_{k=1}^K \pi_k f_k(X)$ .

# Soft K-means Clustering

Use Maximum Likelihood to estimate parameters

$$\theta = (\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2, \pi_1, \dots, \pi_K)$$

based on their log-likelihood

$$\ell(\theta; X) = \sum_{i=1}^N \log \left[ \sum_{k=1}^K \pi_k f_k(x_i; \theta) \right]$$

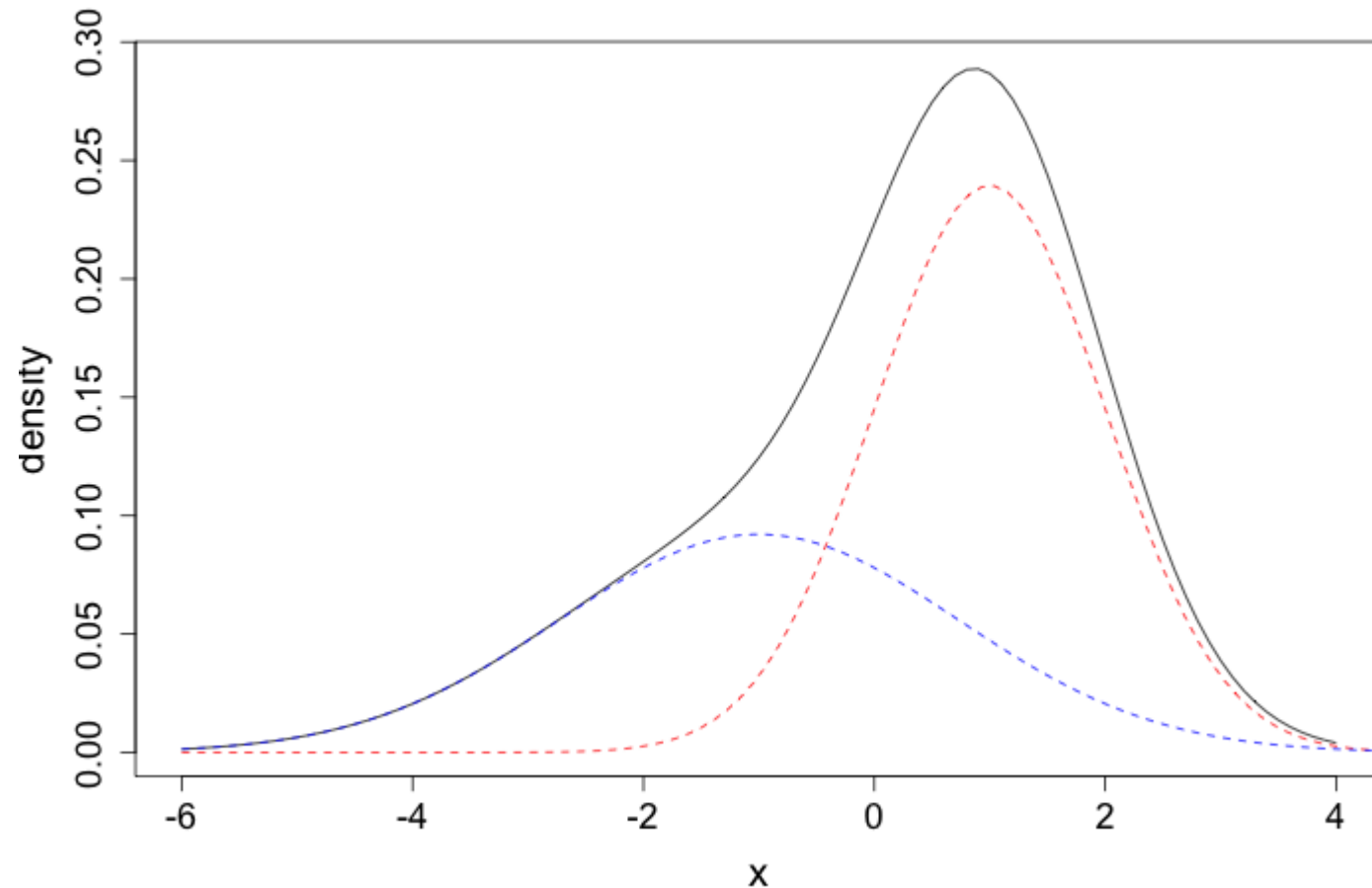
# Soft K-means Clustering

$$\ell(\theta; X) = \sum_{i=1}^N \log \left[ \sum_{k=1}^K \pi_k f_k(x_i; \theta) \right]$$

Maximizing this likelihood directly is computationally difficult

Use Expectation Maximization algorithm (EM) instead.

# Example: Mixture of Two Univariate Gaussians



# Soft K-means Clustering

Consider unobserved latent variables  $\Delta_{ik}$  taking values 0 or 1,

$\Delta_{ij} = 1$  specifies observation  $x_i$  was generated by component  $k$  of the mixture distribution.

# Soft K-means Clustering

Now set  $Pr(\Delta_{ik} = 1) = \pi_k$ , and assume we *observed* values for latent variables  $\Delta_{ik}$ .

We can write the log-likelihood in this case as

$$\ell_0(\theta; X, \Delta) = \sum_{i=1}^N \sum_{k=1}^K \Delta_{ik} \log f_k(x_i; \theta) + \sum_{i=1}^N \sum_{k=1}^K \Delta_{ik} \log \pi_k$$



# Soft K-means Clustering

We have closed-form solutions for maximum likelihood estimates:

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \Delta_{ik} x_i}{\sum_{i=1}^N \Delta_{ik}}$$

$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^N \Delta_{ik} (x_i - \hat{\mu}_k)^2}{\sum_{i=1}^N \Delta_{ik}}$$

$$\hat{\pi}_k = \frac{\sum_{i=1}^K \Delta_{ik}}{N}.$$

# Soft K-means Clustering

Of course, this result depends on observing values for  $\Delta_{ik}$  which *we don't observe*. Use an iterative approach as well:

- given current estimate of parameters  $\theta$ ,
- Substitute  $E[\Delta_{ik}|X_i, \theta]$  for  $\Delta_{ik}$ .

# Soft K-means Clustering

Of course, this result depends on observing values for  $\Delta_{ik}$  which *we don't observe*. Use an iterative approach as well:

- given current estimate of parameters  $\theta$ ,
- Substitute  $E[\Delta_{ik}|X_i, \theta]$  for  $\Delta_{ik}$ .

We will prove that this maximizes the likelihood we need  $\ell(\theta; X)$ .

# Soft K-means Clustering

In the mixture case, what does this look like?

Define

$$\gamma_{ik}(\theta) = E(\Delta_{ik}|X_i, \theta) = Pr(\Delta_{ik} = 1|X_i, \theta)$$

# Soft K-means Clustering

Use Bayes' Rule to write this in terms of the multivariate normal densities with respect to current estimates  $\theta$ :

$$\begin{aligned}\gamma_{ik} &= \frac{Pr(X_i | \Delta_{ik} = 1) Pr(\Delta_{ik} = 1)}{Pr(X_i)} \\ &= \frac{f_k(x_i; \mu_k, \sigma_k^2) \pi_k}{\sum_{l=1}^K f_l(x_i; \mu_l, \sigma_l^2) \pi_l}\end{aligned}$$

# Soft K-means Clustering

Quantity  $\gamma_{ik}(\theta)$  is referred to as the *responsibility* of cluster  $k$  for observation  $i$ , according to current parameter estimate  $\theta$ .

# Soft K-means Clustering

We can now give a complete specification of the EM algorithm for mixture model clustering.

1. Take initial guesses for parameters  $\theta$
2. *Expectation Step*: Compute responsibilities  $\gamma_{ik}(\theta)$
3. *Maximization Step*: Estimate new parameters based on responsibilities as below.
4. Iterate steps 2 and 3 until convergence

# Soft K-means Algorithm

Estimates in the Maximization step are given by

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \gamma_{ik}(\theta) x_i}{\sum_{i=1}^N \gamma_{ik}}$$

$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^N \gamma_{ik}(\theta) (x_i - \mu_k)^2}{\sum_{i=1}^N \gamma_{ik}(\theta)}$$

and

$$\hat{\pi}_k = \frac{\sum_{i=1}^N \gamma_{ik}(\theta)}{N}$$



# Soft K-means Algorithm

The name "soft" K-means refers to the fact that parameter estimates for each cluster are obtained by weighted averages across all observations.

# The EM Algorithm in General

So, why does that work?

Why does plugging in  $\gamma_{ik}(\theta)$  for the latent variables  $\Delta_{ik}$  work?

Why does that maximize log-likelihood  $\ell(\theta; X)$ ?

# The EM Algorithm in General

Think of it as follows:

$z$ : observed data

$z^m$ : missing *latent* data  $T = (Z, Z^m)$ : complete data (observed and missing)

# The EM Algorithm in General

Think of it as follows:

$z$ : observed data

$z^m$ : missing *latent* data  $T = (Z, Z^m)$ : complete data (observed and missing)

$\ell(\theta'; Z)$ : log-likelihood w.r.t. *observed* data

$\ell_0(\theta'; T)$ : log-likelihood w.r.t. *complete* data

# The EM Algorithm in General

Next, notice that

$$Pr(Z|\theta') = \frac{Pr(T|\theta')}{Pr(Z^m|Z, \theta')}$$

# The EM Algorithm in General

Next, notice that

$$Pr(Z|\theta') = \frac{Pr(T|\theta')}{Pr(Z^m|Z, \theta')}$$

As likelihood:

$$\ell(\theta'; Z) = \ell_0(\theta'; T) - \ell_1(\theta'; Z^m|Z)$$

# The EM Algorithm in General

Iterative approach: given parameters  $\theta$  take expectation of log-likelihoods

$$\begin{aligned}\ell(\theta'; Z) &= E[\ell_0(\theta'; T) | Z, \theta] - E[\ell_1(\theta'; Z^m | Z) | Z, \theta] \\ &\equiv Q(\theta', \theta) - R(\theta', \theta)\end{aligned}$$

# The EM Algorithm in General

Iterative approach: given parameters  $\theta$  take expectation of log-likelihoods

$$\begin{aligned}\ell(\theta'; Z) &= E[\ell_0(\theta'; T) | Z, \theta] - E[\ell_1(\theta'; Z^m | Z) | Z, \theta] \\ &\equiv Q(\theta', \theta) - R(\theta', \theta)\end{aligned}$$

In soft k-means,  $Q(\theta', \theta)$  is the log likelihood of complete data with  $\Delta_{ik}$  replaced by  $\gamma_{ik}(\theta)$



# The EM Algorithm in General

The general EM algorithm

1. Initialize parameters  $\theta^{(0)}$
2. Construct *function*  $Q(\theta', \theta^{(j)})$
3. Find next set of parameters  $\theta^{(j+1)} = \arg \max_{\theta'} Q(\theta', \theta^{(j)})$
4. Iterate steps 2 and 3 until convergence

# The EM Algorithm in General

So, why does that work?

$$\begin{aligned}\ell(\theta^{(j+1)}; Z) - \ell(\theta^{(j)}; Z) &= [Q(\theta^{(j+1)}, \theta^{(j)}) - Q(\theta^{(j)}, \theta^{(j)})] \\ &\quad - [R(\theta^{(j+1)}, \theta^{(j)}) - R(\theta^{(j)}, \theta^{(j)})] \\ &\geq 0\end{aligned}$$

# The EM Algorithm in General

So, why does that work?

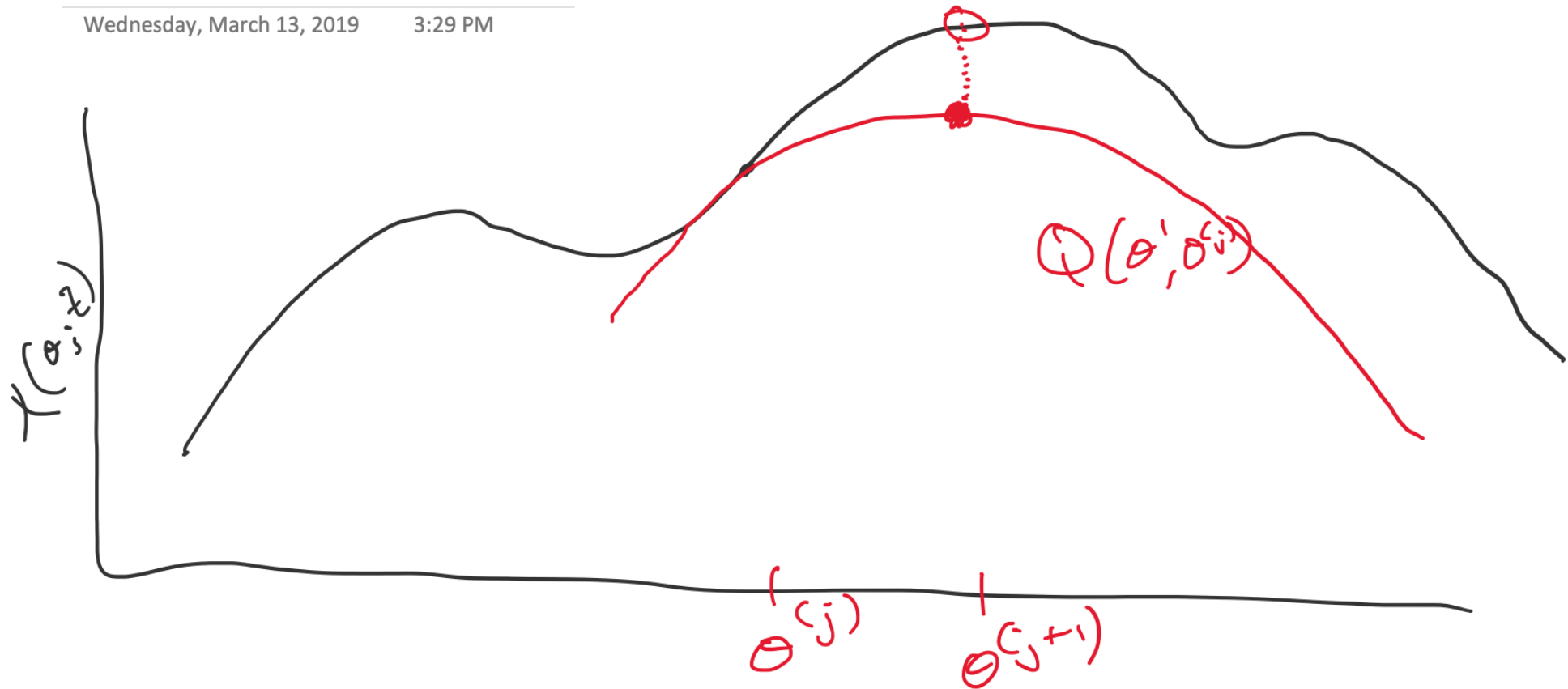
$$\begin{aligned}\ell(\theta^{(j+1)}; Z) - \ell(\theta^{(j)}; Z) &= [Q(\theta^{(j+1)}, \theta^{(j)}) - Q(\theta^{(j)}, \theta^{(j)})] \\ &\quad - [R(\theta^{(j+1)}, \theta^{(j)}) - R(\theta^{(j)}, \theta^{(j)})] \\ &\geq 0\end{aligned}$$

I.E., every step makes log-likelihood larger

# The EM Algorithm in General

Why else does it work?  $Q(\theta', \theta)$  *minorizes*  $\ell(\theta'; Z)$

Wednesday, March 13, 2019 3:29 PM



# The EM Algorithm in General

General algorithmic concept:

Iterative approach:

- Initialize parameters
- Construct bound based on current parameters
- Optimize bound

# Imputing missing data

$z$ : observed data

$z^m$ : missing observations

Requires a likelihood model...

# Latent semantic analysis

Documents as *mixtures* of topics (Hoffman 1998)

