

10-701 Introduction to Machine Learning

The EM Algorithm

Spring 2019

Ameet Talwalkar
(slide credit: Virginia Smith)

Outline

1. Gaussian mixture models
2. GMMs and Incomplete Data
3. EM Algorithm

EM Algorithm

EM algorithm: motivation and setup

- EM is a general procedure to estimate parameters for probabilistic models with hidden/latent variables
- Suppose the model is given by a joint distribution

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})$$

$$p(x, z; \theta)$$

$$\mathcal{D} = \{x_n\} : x : \text{incomplete data}$$

model: $p(x; \theta) = \sum_z p(x, z; \theta)$



$$\begin{aligned}
 l(\theta) &= \log \prod_n p(x_n; \theta) \\
 &= \sum_n \log p(x_n; \theta) \\
 &= \sum_n \log \left(\sum_{z_n} p(x_n, z_n; \theta) \right) \\
 &= \sum_n \log \sum_{z_n} q_n(z_n) \underbrace{\frac{p(x_n, z_n; \theta)}{q_n(z_n)}}_{} \\
 &= \sum_n \log \underbrace{E_q \left(\frac{p(x_n, z_n; \theta)}{q_n(z_n)} \right)}_{} \\
 &\geq \sum_n E_q \left(\log \left(\frac{p(x_n, z_n; \theta)}{q_n(z_n)} \right) \right)
 \end{aligned}$$

Jensen's Inequality

f : strictly concave

i) $f(E(x)) \geq E(f(x))$

ii) equality will hold if x is constant r.v.

$$h(\theta) = E_x \{ g(x, \theta) \}$$

Equality holds if: $\frac{p(x_n, z_n; \theta)}{q_n(z_n)} = \text{constant values of } z_n$

$$\Rightarrow q_n(z_n) \propto p(x_n, z_n; \theta)$$

$$\Rightarrow q_n(z_n) = c p(x_n, z_n; \theta) \quad \text{values of } z_n$$

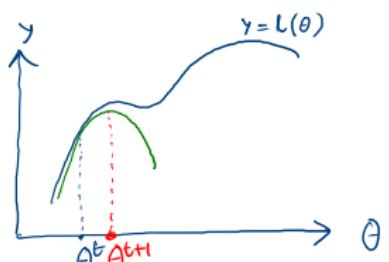
$$\Rightarrow \sum_{z_n} q_n(z_n) = 1 = c \sum_{z_n} p(x_n; z_n; \theta)$$

$$\Rightarrow c = \frac{1}{p(x_n; \theta)} = c p(x_n; \theta)$$

$$q_h(z_n) = \frac{p(z_n, z_n; \theta^t)}{p(z_n; \theta^t)} = p(z_n | z_n; \theta^t)$$

$$l(\theta^t) = \sum_n E_q \left(\log \left(\frac{p(z_n, z_n; \theta^t)}{p(z_n | z_n; \theta^t)} \right) \right)$$

$$l(\theta) \geq \underbrace{\sum_n \sum_{z_n} p(z_n | z_n; \theta^t) \log \frac{p(z_n, z_n; \theta)}{p(z_n | z_n; \theta^t)}}$$



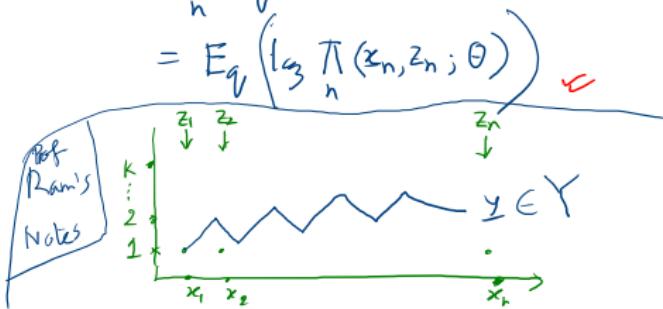
E-Step:

- (1) Compute the posteriors $E_{\theta} \log p(z_n | z_n; \theta^t)$
- (2) $Q(\theta, \theta^t) = \sum_n p(z_n | z_n; \theta^t) \log p(z_n, z_n; \theta) = \sum_n E_{\theta} \log p(z_n, z_n; \theta)$

$$= E_{\theta} \left(\log \prod_n p(z_n, z_n; \theta) \right)$$

M-Step:

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^t)$$



EM algorithm: motivation and setup

- EM is a general procedure to estimate parameters for probabilistic models with hidden/latent variables
- Suppose the model is given by a joint distribution

$$p(\mathbf{x}|\theta) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$$

- Given **incomplete data** $\mathcal{D} = \{\mathbf{x}_n\}$ our goal is to compute MLE of θ :

$$\begin{aligned}\theta &= \arg \max \ell(\theta) = \arg \max \log \mathcal{D} = \arg \max \sum_n \log p(\mathbf{x}_n|\theta) \\ &= \arg \max \sum_n \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\theta)\end{aligned}$$

$\ell(\theta)$

The objective function $\ell(\theta)$ is called **incomplete log-likelihood**

A lower bound

- log-sum form of incomplete log-likelihood is difficult to work with
- EM: construct lower bound on $\ell(\theta)$ (E-step) and optimize it (M-step)

A lower bound

- log-sum form of incomplete log-likelihood is difficult to work with
- EM: construct lower bound on $\ell(\theta)$ (E-step) and optimize it
(M-step)
- If we define $q(z)$ as a distribution over z , then

$$\ell(\theta) = \sum_n \log \sum_{z_n} p(x_n, z_n | \theta)$$

A lower bound

- log-sum form of incomplete log-likelihood is difficult to work with
- EM: construct lower bound on $\ell(\theta)$ (E-step) and optimize it (M-step)
- If we define $q(z)$ as a distribution over z , then

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) \\ &= \sum_n \log \sum_{z_n} q_n(z_n) \underbrace{\frac{p(x_n, z_n | \theta)}{q_n(z_n)}}_{\text{green}}\end{aligned}$$

A lower bound

- log-sum form of incomplete log-likelihood is difficult to work with

- EM: construct lower bound on $\ell(\theta)$ (E-step) and optimize it (M-step)

- If we define $\tilde{q}(z)$ as a distribution over z , then

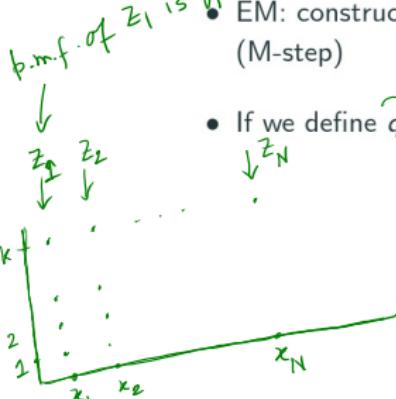
$$\ell(\theta) = \sum_n \log \sum_{z_n} p(x_n, z_n | \theta)$$

$$= \sum_n \log \left[\sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \right]$$

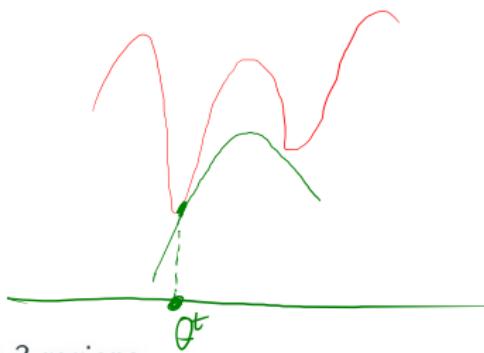
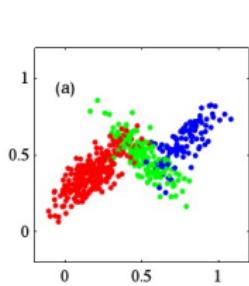
$$\geq \sum_n \left| \sum_{z_n=1}^k q_n(z_n) \log \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \right|$$

$$\begin{aligned} &= E \left\{ \frac{p(x_n, z_n; \theta)}{q_n(z_n)} \right\} \\ &\stackrel{\text{Jensen's}}{\geq} E \left\{ \log \left\{ E \left\{ \frac{p(x_n, z_n; \theta)}{q_n(z_n)} \right\} \right\} \right\} \\ &\stackrel{\text{Jensen's}}{\geq} E \left\{ \log \left\{ \sum_{z_n=1}^k q_n(z_n) \frac{p(x_n, z_n; \theta)}{q_n(z_n)} \right\} \right\} \\ &= \sum_{z_n=1}^k q_n(z_n) \log \frac{p(x_n, z_n; \theta)}{q_n(z_n)} \end{aligned}$$

- Last step follows from Jensen's inequality, i.e., $f(\mathbb{E}X) \geq \mathbb{E}f(X)$ for concave function f



GMM Example



- Consider the previous model where x could be from 3 regions
- We can choose $q(z)$ as any valid distribution
- e.g., $q(z = k) = 1/3$ for any of 3 colors
- e.g., $q(z = k) = 1/2$ for red and blue, 0 for green

Which $q(z)$ should we choose?

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \frac{p(x_n, z_n | \theta)}{q_n(z_n)}\end{aligned}$$

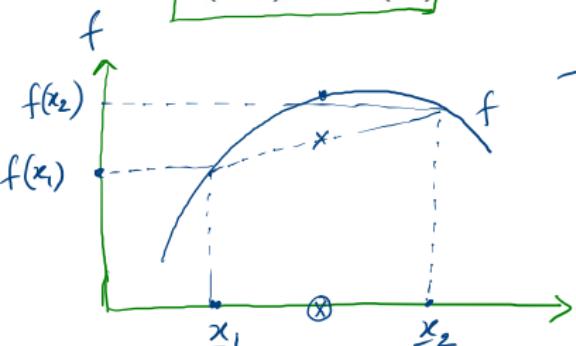
- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a tight lower bound

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \frac{p(x_n, z_n | \theta)}{q_n(z_n)}\end{aligned}$$

- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a *tight* lower bound, and given some current estimate θ^t , we will pick $q_n(\cdot)$ such that our lower bound holds *with equality* at θ^t
- $f(\mathbb{E}X) = \mathbb{E}f(X)?$



Concavity:

$$f(tz_1 + (1-t)z_2) \geq t f(z_1) + (1-t) f(z_2)$$

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \frac{p(x_n, z_n | \theta)}{q_n(z_n)}\end{aligned}$$

- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a *tight* lower bound, and given some current estimate θ^t , we will pick $q_n(\cdot)$ such that our lower bound holds *with equality* at θ^t
- $f(\mathbb{E}X) = \mathbb{E}f(X)$? It is sufficient for X to be a constant random variable!

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \boxed{\frac{p(x_n, z_n | \theta)}{q_n(z_n)}}\end{aligned}$$

- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a *tight* lower bound, and given some current estimate θ^t , we will pick $q_n(\cdot)$ such that our lower bound holds *with equality* at θ^t
- $f(\mathbb{E}X) = \mathbb{E}f(X)$? It is sufficient for X to be a constant random variable!
- Choose $\boxed{q_n(z_n) \propto p(x_n, z_n | \theta^t)}$

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \left[\frac{p(x_n, z_n | \theta)}{q_n(z_n)} \right] = C \quad \Rightarrow \quad q_{\hat{n}}(z_n) = \frac{1}{C} p(x_n, z_n | \theta) \\ &\Rightarrow \sum_{z_n} q_{\hat{n}}(z_n) = 1 = \frac{1}{C} \sum_{z_n} p(x_n, z_n | \theta)\end{aligned}$$

- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a *tight* lower bound, and given some current estimate θ^t , we will pick $q_n(\cdot)$ such that our lower bound holds *with equality* at θ^t
- $f(\mathbb{E}X) = \mathbb{E}f(X)$? It is sufficient for X to be a constant random variable!

- Choose $q_n(z_n) \propto p(x_n, z_n | \theta^t)$! Since $q_n(\cdot)$ is a distribution, we have

$$q_n(z_n) = \frac{p(x_n, z_n | \theta^t)}{\sqrt{\sum_k p(x_n, z_n = k | \theta^t)}} = \frac{p(x_n, z_n | \theta^t)}{p(x_n | \theta^t)} = p(z_n | x_n; \theta^t)$$

Which $q(z)$ to choose?

Recall:

$$\begin{aligned}\ell(\theta) &= \sum_n \log \sum_{z_n} p(x_n, z_n | \theta) = \sum_n \log \sum_{z_n} q_n(z_n) \frac{p(x_n, z_n | \theta)}{q_n(z_n)} \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \frac{p(x_n, z_n | \theta)}{q_n(z_n)}\end{aligned}$$

- The lower bound we derived for $\ell(\theta)$ holds for all choices of $q(\cdot)$
- We want a *tight* lower bound, and given some current estimate θ^t , we will pick $q_n(\cdot)$ such that our lower bound holds *with equality* at θ^t
- $f(\mathbb{E}X) = \mathbb{E}f(X)$? It is sufficient for X to be a constant random variable!
- Choose $q_n(z_n) \propto p(x_n, z_n | \theta^t)$! Since $q_n(\cdot)$ is a distribution, we have

$$q_n(z_n) = \frac{p(x_n, z_n | \theta^t)}{\sum_k p(x_n, z_n = k | \theta^t)} = \frac{p(x_n, z_n | \theta^t)}{p(x_n | \theta^t)} = p(z_n | x_n; \theta^t)$$

- This is the **posterior distribution** of z_n given x_n and θ^t

E and M Steps

Our simplified expression

$$\ell(\theta^t) = \sum_n \sum_{z_n} p(z_n|x_n; \theta^t) \log \frac{p(x_n, z_n|\theta^t)}{p(z_n|x_n; \theta^t)}$$

E and M Steps

Our simplified expression

$$\ell(\theta^t) = \sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log \frac{p(x_n, z_n | \theta^t)}{p(z_n | x_n; \theta^t)}$$

E-Step: For all n , compute $q_n(z_n) = p(z_n | x_n; \theta^t)$

Why is this called the E-Step?

E and M Steps

Our simplified expression

$$\ell(\theta^t) = \sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log \frac{p(x_n, z_n | \theta^t)}{p(z_n | x_n; \theta^t)}$$

E-Step: For all n , compute $q_n(z_n) = p(z_n | x_n; \theta^t)$

Why is this called the E-Step? Because we can view it as computing the expected (complete) log-likelihood:

Q(θ|θ^t) = $\sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log p(x_n, z_n | \theta) = \underbrace{\mathbb{E}_q \sum_n \log p(x_n, z_n | \theta)}_{\text{---}} - \sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log p(z_n | x_n; \theta^t)$

E and M Steps

Our simplified expression

$$\ell(\theta^t) = \sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log \frac{p(x_n, z_n | \theta^t)}{p(z_n | x_n; \theta^t)}$$

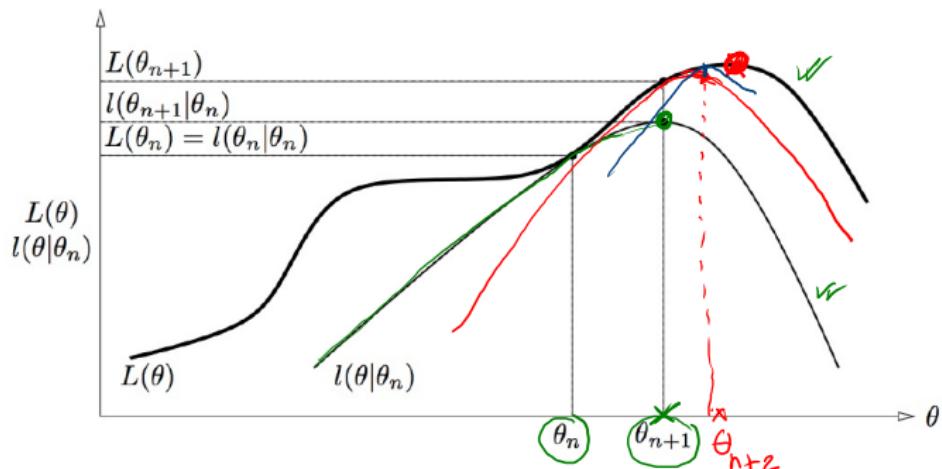
E-Step: For all n , compute $\boxed{q_n(z_n) = p(z_n | x_n; \theta^t)}$

Why is this called the E-Step? Because we can view it as computing the **expected (complete) log-likelihood**:

$$Q(\theta | \theta^t) = \sum_n \sum_{z_n} p(z_n | x_n; \theta^t) \log p(x_n, z_n | \theta) = \mathbb{E}_q \sum_n \log p(x_n, z_n | \theta)$$

M-Step: Maximize $Q(\theta | \theta^t)$, i.e., $\underline{\theta^{t+1}} = \arg \max_{\theta} Q(\theta | \theta^t)$

EM in Pictures



(Figure from tutorial by Sean Borman)

Iterative and monotonic improvement

- We can show that $\ell(\theta^{t+1}) \geq \ell(\theta^t)$

Iterative and monotonic improvement

- We can show that $\ell(\boldsymbol{\theta}^{t+1}) \geq \ell(\boldsymbol{\theta}^t)$
- Recall that we chose $q(\cdot)$ in the E-step such that:

$$\ell(\boldsymbol{\theta}^t) = \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log \frac{p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta}^t)}{q(\mathbf{z}_n)}$$

Iterative and monotonic improvement

- We can show that $\ell(\theta^{t+1}) \geq \ell(\theta^t)$
- Recall that we chose $q(\cdot)$ in the E-step such that:

$$\ell(\theta^t) = \sum_n \sum_{z_n} q(z_n) \log \frac{p(x_n, z_n | \theta^t)}{q(z_n)}$$

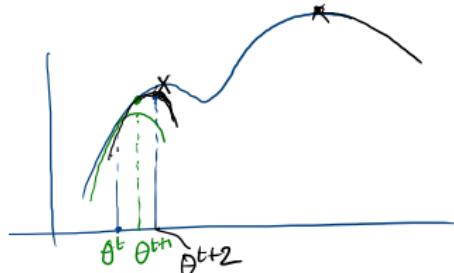
- However, in the M-step, θ^{t+1} is chosen to maximize the right hand side of the equation, thus proving our desired result

$$\begin{aligned}\ell(\theta^{t+1}) &\geq \text{ELBO}(\theta, \theta^{t+1}) \\ &\geq \ell(\theta^t)\end{aligned}$$

- We can show that $\ell(\theta^{t+1}) \geq \ell(\theta^t)$
- Recall that we chose $q(\cdot)$ in the E-step such that:

$$\ell(\theta^t) = \sum_n \sum_{z_n} q(z_n) \log \frac{p(x_n, z_n | \theta^t)}{q(z_n)}$$

- However, in the M-step, θ^{t+1} is chosen to maximize the right hand side of the equation, thus proving our desired result
- Note: the EM procedure converges but only to a local optimum



You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters

posteriori
H

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters
 - (2) Using this “complete” data to find the maximum likelihood parameter estimates

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters
 - (2) Using this “complete” data to find the maximum likelihood parameter estimates
- Pros: Guaranteed to converge, no parameters to tune (e.g., compared to gradient methods)

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters
 - (2) Using this “complete” data to find the maximum likelihood parameter estimates
- Pros: Guaranteed to converge, no parameters to tune (e.g., compared to gradient methods)
- Cons: Can get stuck in local optima, can be expensive

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters
 - (2) Using this “complete” data to find the maximum likelihood parameter estimates
- Pros: Guaranteed to converge, no parameters to tune (e.g., compared to gradient methods)
- Cons: Can get stuck in local optima, can be expensive
- Why is EM useful for unsupervised learning?

You should know . . .

- EM is a general procedure for maximizing a likelihood with *latent (unobserved) variables*
- The two steps of EM:
 - (1) Estimating unobserved data from observed data and current parameters
 - (2) Using this “complete” data to find the maximum likelihood parameter estimates
- Pros: Guaranteed to converge, no parameters to tune (e.g., compared to gradient methods)
- Cons: Can get stuck in local optima, can be expensive
- Why is EM useful for unsupervised learning?
 - EM is a general method to deal with hidden data; we have studied it in the context of hidden *labels* (unsupervised learning)