

Machine Learning

10-701/15-781, Fall 2011

Expectation Maximization

Eric Xing

Lecture 9, October 10, 2011

Reading: Chap. 9, C.B book



Eric Xing

© Eric Xing @ CMU, 2006-2011

1

Poster data?

Thursday, 12/8

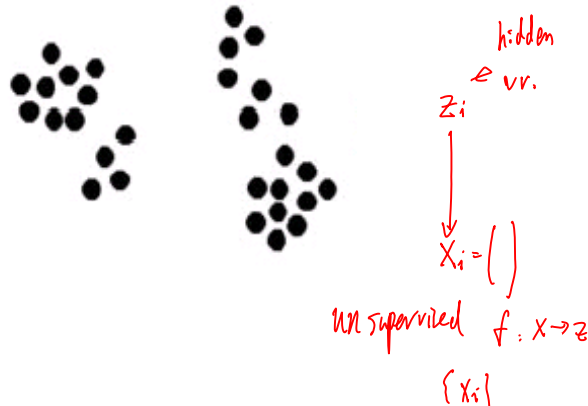


Eric Xing

© Eric Xing @ CMU, 2006-2011

2

Clustering and partially observable probabilistic models

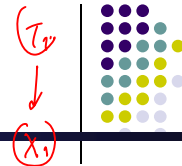


Eric Xing

© Eric Xing @ CMU, 2006-2011

3

Unobserved Variables



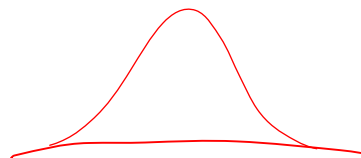
- A variable can be unobserved (latent) because:
 - it is an imaginary quantity meant to provide some simplified and abstractive view of the data generation process
 - e.g., speech recognition models, mixture models ...
 - it is a real-world object and/or phenomena, but difficult or impossible to measure
 - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
 - it is a real-world object and/or phenomena, but sometimes wasn't measured, because of faulty sensors; or was measure with a noisy channel, etc.
 - e.g., traffic radio, aircraft signal on a radar screen,
- Discrete latent variables can be used to partition/cluster data into sub-groups (mixture models, forthcoming).
- Continuous latent variables (factors) can be used for dimensionality reduction (factor analysis, etc., later lectures).

Eric Xing

© Eric Xing @ CMU, 2006-2011

4

Uni-modal and multi-modal distributions



$$P(x) = \mathcal{N}(x | \mu, \Sigma)$$

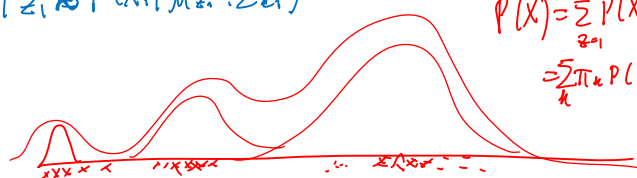
$$X_i | z_i \sim \mathcal{P}(X_i | \mu_{z_i}, \Sigma_{z_i})$$

$$P(x|z)P(z)$$

$$= P(x, z)$$

$$P(x) = \sum_{z=1}^K P(x, z)$$

$$= \sum_k \pi_k P(x | z_i = k)$$



$$X_i \rightarrow z_i$$

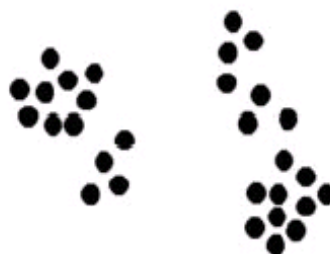
$$\pi_k = P(z_i = k)$$

Eric Xing

© Eric Xing @ CMU, 2006-2011

5

Mixture Models



$$X_i | z_i \sim \mathcal{N}(\mu, \Sigma)$$

$$mix \rightarrow X_i \sim \sum_k \pi_k \mathcal{N}(\mu_k, \Sigma_k)$$

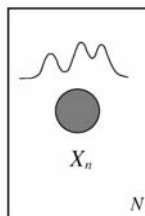
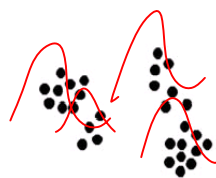
Eric Xing

© Eric Xing @ CMU, 2006-2011

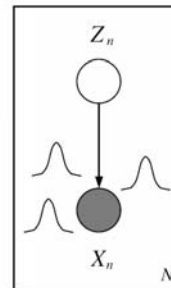
6

Mixture Models, con'd

- A density model $p(x)$ may be multi-modal.
- We may be able to model it as a mixture of uni-modal distributions (e.g., Gaussians).
- Each mode may correspond to a different sub-population (e.g., male and female).



(a)



(b)

Eric Xing

© Eric Xing @ CMU, 2006-2011

7

Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:
 - Z is a latent class indicator vector:

$$p(Z_n) = \text{multi}(Z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

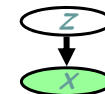
- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$p(x_n) = \sum_{z_n} p(x_n, z_n) = \sum_{z_n} p(x_n | z_n) p(z_n)$$

$$p(x_n | \mu, \Sigma) = \sum_k p(z^k = 1 | \pi) p(x_n | z^k = 1, \mu, \Sigma) = \sum_k \pi_k \sum_{z_n} \left((\pi_k)^{z_n^k} \mathcal{N}(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$$



$$Z_n = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

mixture component

mixture proportion

Eric Xing

© Eric Xing @ CMU, 2006-2011

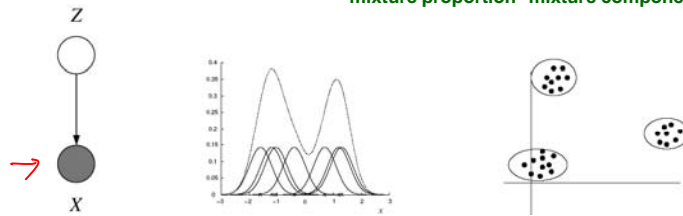
8

Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$$

mixture proportion mixture component



- This model can be used for unsupervised clustering.
 - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

Eric Xing

© Eric Xing @ CMU, 2006-2011

9

Learning mixture models

Goal: $\theta = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_K \\ \mu_1 \\ \vdots \\ \mu_K \\ \Sigma_1 \\ \vdots \\ \Sigma_K \end{bmatrix}$

$z_n^k \rightarrow \langle z_n^k \rangle_{p(z_n^k | x_n)}$

θ - est. or random

$\Sigma_k \quad \forall k$

Data: $X = [x_1, \dots, x_N]$

Suppose Data $\equiv \begin{matrix} X = [x_1, \dots, x_N] \\ Z = [z_1, \dots, z_N] \end{matrix}$

MLE:

$\pi_k \equiv \frac{\sum_n \mathbb{1}(z_n^k)}{N}$
 $\mu_k \equiv \frac{\sum_n x_n \mathbb{1}(z_n^k)}{\sum_n \mathbb{1}(z_n^k)}$
 $\Sigma_k \equiv \frac{\sum_n (x_n - \mu_k)(x_n - \mu_k)^T \mathbb{1}(z_n^k)}{\sum_n \mathbb{1}(z_n^k)}$

Eric Xing

© Eric Xing @ CMU, 2006-2011

10

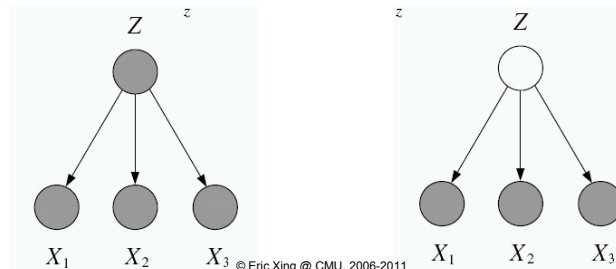
Why is Learning Harder?

- In fully observed iid settings, the log likelihood decomposes into a sum of local terms.

$$\mathcal{L}_c(\theta; D) = \log p(x, z | \theta) = \log p(z | \theta_z) + \log p(x | z, \theta_x)$$

- With latent variables, all the parameters become coupled together via *marginalization*

$$\mathcal{L}_c(\theta; D) = \log \sum_z p(x, z | \theta) = \log \sum_z p(z | \theta_z) p(x | z, \theta_x)$$



Eric Xing

© Eric Xing @ CMU, 2006-2011

11

Gradient Learning for mixture models

- We can learn mixture densities using gradient descent on the log likelihood. The gradients are quite interesting:

$$\mathcal{L}(\theta) = \log p(x | \theta) = \log \sum_k \pi_k p_k(x | \theta_k)$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{1}{p(x | \theta)} \sum_k \pi_k \frac{\partial p_k(x | \theta_k)}{\partial \theta}$$

π_k is u_k
 \sum_k

$$= \sum_k \frac{\pi_k}{p(x | \theta)} p_k(x | \theta_k) \frac{\partial \log p_k(x | \theta_k)}{\partial \theta}$$

$$= \sum_k \pi_k \frac{p_k(x | \theta_k)}{p(x | \theta)} \frac{\partial \log p_k(x | \theta_k)}{\partial \theta_k} = \sum_k r_k \frac{\partial \mathcal{L}_k}{\partial \theta_k}$$

- In other words, the gradient is the responsibility weighted sum of the individual log likelihood gradients.
- Can pass this to a conjugate gradient routine.

Eric Xing

© Eric Xing @ CMU, 2006-2011

12

Parameter Constraints

$$\theta^{(0)} \rightarrow \theta^{(1)} \rightarrow \theta^{(2)}$$



- Often we have constraints on the parameters, e.g. $\sum_k \pi_k = 1$, Σ being symmetric positive definite (hence $\Sigma_{ii} > 0$).
- We can use constrained optimization, or we can reparameterize in terms of unconstrained values.

$$\pi_k = \frac{e^{\lambda_k}}{\sum_j e^{\lambda_j}}$$

- For normalized weights, use the softmax transform:
- For covariance matrices, use the Cholesky decomposition:

$$\Sigma^{-1} = \mathbf{A}^T \mathbf{A}$$

where \mathbf{A} is upper diagonal with positive diagonal:

$$\mathbf{A}_{ii} = \exp(\lambda_i) > 0 \quad \mathbf{A}_{ij} = \eta_{ij} \quad (j > i) \quad \mathbf{A}_{ij} = 0 \quad (j < i)$$

the parameters $\gamma, \lambda, \eta_{ij} \in \mathbb{R}$ are unconstrained.

- Use chain rule to compute $\frac{\partial \ell}{\partial \pi}, \frac{\partial \ell}{\partial \mathbf{A}}$.

Eric Xing

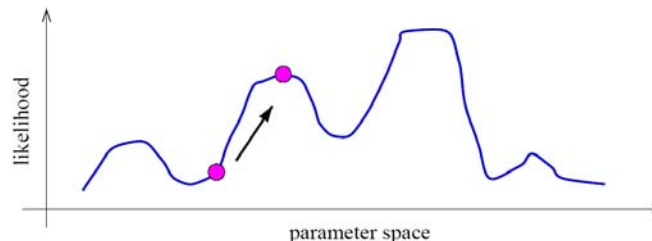
© Eric Xing @ CMU, 2006-2011

13

Identifiability



- A mixture model induces a multi-modal likelihood.
- Hence gradient ascent can only find a local maximum.
- Mixture models are unidentifiable, since we can always switch the hidden labels without affecting the likelihood.
- Hence we should be careful in trying to interpret the “meaning” of latent variables.



Eric Xing

© Eric Xing @ CMU, 2006-2011

14

Identifiability



$$\pi = \begin{bmatrix} 0.1 \\ 0.2 \\ 0.7 \end{bmatrix}$$

$$\pi = \begin{bmatrix} 0.2 \\ 0.7 \\ 0.1 \end{bmatrix}$$

Toward the EM algorithm



- E.g., A mixture of K Gaussians:

- Z is a latent class indicator vector

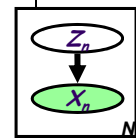
$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z_n^k = 1 | \pi) p(x_n | z_n^k = 1, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k) \end{aligned}$$



Toward the EM algorithm



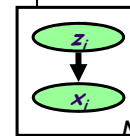
- Recall MLE for **completely observed data**
- Data log-likelihood**

$$\begin{aligned}\mathcal{L}(\theta; D) &= \log \prod_n p(z_n, x_n) = \log \prod_n p(z_n | \pi) p(x_n | z_n, \mu, \sigma) \\ &= \sum_n \log \prod_k \pi_k^{z_n^k} + \sum_n \log \prod_k N(x_n; \mu_k, \sigma)^{z_n^k} \\ &= \sum_n \sum_k z_n^k \log \pi_k - \sum_n \sum_k \frac{1}{2\sigma^2} (x_n - \mu_k)^2 + C\end{aligned}$$

- MLE**

$$\begin{aligned}\hat{\pi}_{k,MLE} &= \arg \max_{\pi} \mathcal{L}(\theta; D), \\ \hat{\mu}_{k,MLE} &= \arg \max_{\mu} \mathcal{L}(\theta; D) \\ \hat{\sigma}_{k,MLE} &= \arg \max_{\sigma} \mathcal{L}(\theta; D)\end{aligned}\quad \Rightarrow \quad \hat{\mu}_{k,MLE} = \frac{\sum_n z_n^k x_n}{\sum_n z_n^k}$$

- What if we do not know z_n ?**



Expectation-Maximization (EM) Algorithm



Expectation-Maximization (EM) Algorithm



- EM is an optimization strategy for objective functions that can be interpreted as likelihoods in the presence of missing data.
- It is much simpler than gradient methods:
 - No need to choose step size.
 - Enforces constraints automatically.
 - Calls inference and fully observed learning as subroutines.
- EM is an iterative algorithm with two linked steps:
 - E-step: fill-in hidden values using inference, $p(z|x, \theta)$.
 - M-step: update parameters $t+1$ using standard MLE/MAP method applied to completed data
- We will prove that this procedure monotonically improves (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.

Eric Xing

© Eric Xing @ CMU, 2006-2011

19

K-means

$$z_n = \arg \min_k (x_n - \mu_k)^2$$

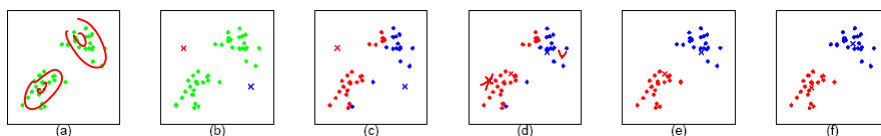


- Start:
 - "Guess" the centroid μ_k and covariance Σ_k of each of the K clusters
- Loop
 - For each point $n=1$ to N , compute its cluster label:

$$z_n^{(t)} = \arg \max_k (x_n - \mu_k^{(t)})^T \Sigma_k^{-1(t)} (x_n - \mu_k^{(t)})$$

- For each cluster $k=1:K$

$$\mu_k^{(t+1)} = \frac{\sum_n \delta(z_n^{(t)}, k) x_n}{\sum_n \delta(z_n^{(t)}, k)} \quad \Sigma_k^{(t+1)} = \dots$$



Eric Xing

© Eric Xing @ CMU, 2006-2011

20

Expectation-Maximization

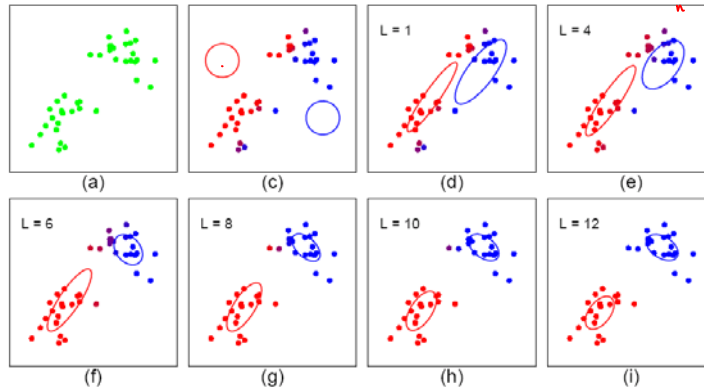
$$E: z_n = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\langle z_n \rangle = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

- Start:
 - "Guess" the centroid μ_k and coverage Σ_k of each of the K clusters
- Loop

fractional / safe

$$M: \mu_k = \frac{\sum_n \langle z_n \rangle x_n}{\sum_n \langle z_n \rangle}$$



Eric Xing

© Eric Xing @ CMU, 2006-2011

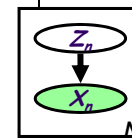
21

How is EM derived?

- A mixture of K Gaussians:
 - Z is a latent class indicator vector

$$p(z_n) = \text{multi}(z_n; \pi) = \prod_k (\pi_k)^{z_n^k}$$
 - X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$
- The likelihood of a sample:



- The "complete" likelihood

$$p(x_n, z_n^k = 1 | \mu, \Sigma) = p(z_n^k = 1 | \pi) p(x_n | z_n^k = 1, \mu, \Sigma) = \pi_k N(x_n | \mu_k, \Sigma_k)$$

$$p(x_n, z_n | \mu, \Sigma) = \prod_k [\pi_k N(x_n | \mu_k, \Sigma_k)]^{z_n^k}$$

But this is itself a random variable! Not good as objective function

Eric Xing

© Eric Xing @ CMU, 2006-2011

22

How is EM derived?

- The complete log likelihood:

$$\begin{aligned}
 \mathcal{L}(\theta; D) &= \log \prod_n p(z_n, x_n) = \log \prod_n p(z_n | \pi) p(x_n | z_n, \mu, \sigma) \\
 &= \sum_n \log \prod_k \pi_k^{z_n^k} + \sum_n \log \prod_k N(x_n; \mu_k, \sigma)^{z_n^k} \\
 &= \sum_n \sum_k \underbrace{z_n^k}_{\text{r.v.}} \log \pi_k - \sum_n \sum_k \underbrace{z_n^k}_{\text{r.v.}} \frac{1}{2\sigma^2} (x_n - \mu_k)^2 + C
 \end{aligned}$$

- The expected complete log likelihood

$$\begin{aligned}
 \langle \mathcal{L}_c(\theta; \mathbf{x}, \mathbf{z}) \rangle &= \sum_n \langle \log p(z_n | \pi) \rangle_{p(z|\mathbf{x})} + \sum_n \langle \log p(x_n | z_n, \mu, \Sigma) \rangle_{p(z|\mathbf{x})} \\
 &= \sum_n \sum_k \langle z_n^k \rangle \log \pi_k - \frac{1}{2} \sum_n \sum_k \langle z_n^k \rangle ((x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log |\Sigma_k| + C)
 \end{aligned}$$

$q(z_n) \propto p(z_n | x_n) \in p(z_1, z_2)$
 $?$

Eric Xing

© Eric Xing @ CMU, 2006-2011

23

E-step

- We maximize $\langle \mathcal{L}_c(\theta) \rangle$ iteratively using the following iterative procedure:

- Expectation step:** computing the expected value of the sufficient statistics of the hidden variables (i.e., \mathbf{z}) given current est. of the parameters (i.e., π and μ).

$$\tau_n^{k(t)} = \langle z_n^k \rangle_{q^{(t)}} = p(z_n^k = 1 | \mathbf{x}_n, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(\mathbf{x}_n | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

\uparrow
 $p(z_n | x_n)$

- Here we are essentially doing **inference**

Eric Xing

© Eric Xing @ CMU, 2006-2011

24

M-step

- We maximize $\langle J_c(\theta) \rangle$ iteratively using the following iterative procedure:

- Maximization step:** compute the parameters under current results of the expected value of the hidden variables

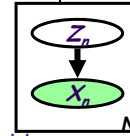
$$\pi_k^* = \arg \max \langle J_c(\theta) \rangle, \quad \Rightarrow \quad \frac{\partial}{\partial \pi_k} \langle J_c(\theta) \rangle = 0, \forall k, \quad \text{s.t.} \quad \sum_k \pi_k = 1$$

$$\Rightarrow \pi_k^* = \frac{\sum_n \langle z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \tau_n^{k(t)}}{N} = \frac{\langle n_k \rangle}{N}$$

$$\mu_k^* = \arg \max \langle J(\theta) \rangle, \quad \Rightarrow \quad \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

$$\Sigma_k^* = \arg \max \langle J(\theta) \rangle, \quad \Rightarrow \quad \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

- This is isomorphic to **MLE** except that the variables that are hidden are replaced by their expectations (in general they will be replaced by their corresponding "**sufficient statistics**")



Fact:

$$\frac{\partial \log |A^{-1}|}{\partial A^{-1}} = A^T$$

$$\frac{\partial x^T A x}{\partial A} = x x^T$$

Eric Xing

© Eric Xing @ CMU, 2006-2011

25

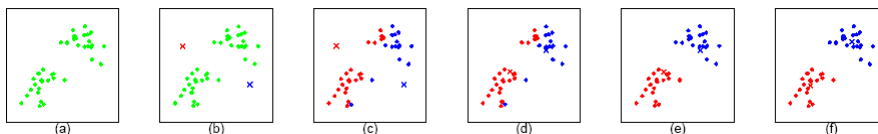
Compare: K-means

- The EM algorithm for mixtures of Gaussians is like a "soft version" of the K-means algorithm.
- In the K-means "E-step" we do hard assignment:

$$z_n^{(t)} = \arg \max_k (x_n - \mu_k^{(t)})^T \Sigma_k^{-1(t)} (x_n - \mu_k^{(t)})$$

- In the K-means "M-step" we update the means as the weighted sum of the data, but now the weights are 0 or 1:

$$\mu_k^{(t+1)} = \frac{\sum_n \delta(z_n^{(t)}, k) x_n}{\sum_n \delta(z_n^{(t)}, k)}$$



Eric Xing

© Eric Xing @ CMU, 2006-2011

26

Theory underlying EM



- What are we doing?
- Recall that according to MLE, we intend to learn the model parameter that would have maximize the likelihood of the data.
- But we do not observe z , so computing

$$\ell_c(\theta; D) = \log \sum_z p(x, z | \theta) = \log \sum_z p(z | \theta_z) p(x | z, \theta_x)$$

is difficult!

- What shall we do?

Eric Xing

© Eric Xing @ CMU, 2006-2011

27

Complete & Incomplete Log Likelihoods



- Complete log likelihood
Let X denote the observable variable(s), and Z denote the latent variable(s).
If Z could be observed, then

$$\ell_c(\theta; \mathbf{x}, \mathbf{z}) \stackrel{\text{def}}{=} \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Usually, optimizing $\ell_c()$ given both z and x is straightforward (c.f. MLE for fully observed models).
- Recalled that in this case the objective for, e.g., MLE, decomposes into a sum of factors, the parameter for each factor can be estimated separately.
- **But given that Z is not observed, $\ell_c()$ is a random quantity, cannot be maximized directly.**

- Incomplete log likelihood

With z unobserved, our objective becomes the log of a marginal probability:

$$\ell_c(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_z p(\mathbf{x}, \mathbf{z} | \theta)$$

- **This objective won't decouple**

Eric Xing

© Eric Xing @ CMU, 2006-2011

28

Expected Complete Log Likelihood

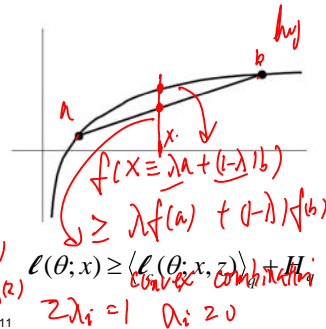
- For **any** distribution $q(z)$, define **expected complete log likelihood**:

$$\langle \ell_c(\theta; x, z) \rangle_q \stackrel{\text{def}}{=} \sum_z q(z | x, \theta) \log p(x, z | \theta)$$

- A deterministic function of θ *any posterior of z*
- Linear in $\ell_c()$ --- inherit its factorizability
- Does maximizing this surrogate yield a maximizer of the likelihood?

- Jensen's inequality

$$\begin{aligned} \ell(\theta; x) &= \log p(x | \theta) \\ &= \log \sum_z p(x, z | \theta) \\ &= \log \sum_z q(z | x) \frac{p(x, z | \theta)}{q(z | x)} \\ &\geq \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\ &= \sum_z q(z | x) \log p(x, z) - \sum_z q(z | x) \log q(z) \end{aligned}$$



Eric Xing

© Eric Xing @ CMU, 2006-2011

29

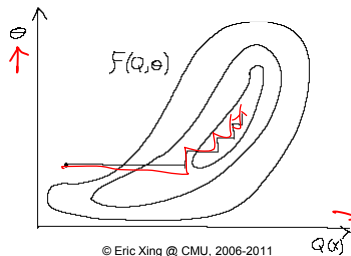
Lower Bounds and Free Energy

- For fixed data x , define a functional called the free energy:

$$F(q, \theta) \stackrel{\text{def}}{=} \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \leq \ell(\theta; x)$$

- The EM algorithm is coordinate-ascent on F :

- E-step:** $q^{t+1} = \arg \max_q F(q, \theta^t)$
- M-step:** $\theta^{t+1} = \arg \max_{\theta} F(q^{t+1}, \theta)$



Eric Xing

© Eric Xing @ CMU, 2006-2011

30

E-step: maximization of expected ℓ_c w.r.t. q



- Claim:

$$q^{t+1} = \arg \max_q F(q, \theta^t) = p(z | x, \theta^t)$$

- This is the posterior distribution over the latent variables given the data and the parameters. Often we need this at test time anyway (e.g. to perform classification).

- Proof (easy): this setting attains the bound $\ell(\theta; x) \geq F(q, \theta)$

$$\begin{aligned} F(p(z|x, \theta'), \theta') &= \sum_z p(z|x, \theta') \log \frac{p(x, z | \theta')}{p(z|x, \theta')} \\ &= \sum_z p(z|x, \theta') \log p(x | \theta') \\ &= \log p(x | \theta') = \ell(\theta'; x) \end{aligned}$$

- Can also show this result using variational calculus or the fact that $\ell(\theta; x) - F(q, \theta) = \text{KL}(q \| p(z | x, \theta))$

Eric Xing

© Eric Xing @ CMU, 2006-2011

31

E-step \equiv plug in posterior expectation of latent variables



- Without loss of generality: assume that $p(x, z | \theta)$ is a generalized exponential family distribution:

$$p(x, z | \theta) = \frac{1}{Z(\theta)} h(x, z) \exp \left\{ \sum_i \theta_i f_i(x, z) \right\}$$

$$\begin{aligned} & \hat{p}(z) \\ & = p(z|x) \end{aligned}$$

- Special cases: if $p(x|z)$ are GLIMs, then

$$f_i(x, z) = \eta_i^T(z) \xi_i(x) \quad \langle \eta \rangle p(z|x)$$

- The expected complete log likelihood under $q^{t+1} = p(z | x, \theta^t)$ is

$$\begin{aligned} \langle \ell_c(\theta^t; x, z) \rangle_{q^{t+1}} &= \sum_z q(z | x, \theta^t) \log p(x, z | \theta^t) - A(\theta) \\ &= \sum_i \theta_i^t \langle f_i(x, z) \rangle_{q(z|x, \theta^t)} - A(\theta) \\ &\stackrel{p\text{-GLIM}}{=} \sum_i \theta_i^t \langle \eta_i(z) \rangle_{q(z|x, \theta^t)} \xi_i(x) - A(\theta) \end{aligned}$$

Eric Xing

© Eric Xing @ CMU, 2006-2011

32

M-step: maximization of expected ℓ_c w.r.t. θ



- Note that the free energy breaks into two terms:

$$\begin{aligned} F(q, \theta) &= \sum_z q(z | \mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})} \\ &= \sum_z q(z | \mathbf{x}) \log p(\mathbf{x}, \mathbf{z} | \theta) - \sum_z q(z | \mathbf{x}) \log q(\mathbf{z} | \mathbf{x}) \\ &= \langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_q + H_q \end{aligned}$$

- The first term is the expected complete log likelihood (energy) and the second term, which does not depend on θ , is the entropy.
- Thus, in the M-step, maximizing with respect to θ for fixed q we only need to consider the first term:

$$\theta^{t+1} = \arg \max_{\theta} \langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_{q^{t+1}} = \arg \max_{\theta} \sum_z q(z | \mathbf{x}) \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Under optimal q^{t+1} , this is equivalent to solving a standard MLE of fully observed model $p(\mathbf{x}, \mathbf{z} | \theta)$, with the sufficient statistics involving \mathbf{z} replaced by their expectations w.r.t. $p(\mathbf{z} | \mathbf{x}, \theta)$.

Eric Xing

© Eric Xing @ CMU, 2006-2011

33

Summary: EM Algorithm



- A way of maximizing likelihood function for latent variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces:
 - Estimate some “missing” or “unobserved” data from observed data and current parameters.
 - Using this “complete” data, find the maximum likelihood parameter estimates.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
 - E-step: $q^{t+1} = \arg \max_q F(q, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \max_{\theta} F(q^{t+1}, \theta^t)$
- In the M-step we optimize a lower bound on the likelihood. In the E-step we close the gap, making bound=likelihood.

Eric Xing

© Eric Xing @ CMU, 2006-2011

34

EM Variants

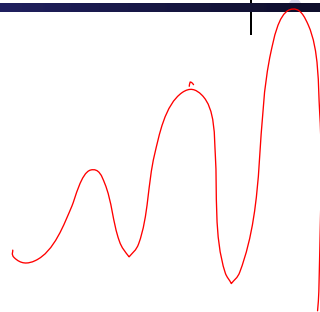


- Sparse EM:
Do not re-compute exactly the posterior probability on each data point under all models, because it is almost zero. Instead keep an “active list” which you update every once in a while.
- Generalized (Incomplete) EM:
It might be hard to find the ML parameters in the M-step, even given the completed data. We can still make progress by doing an M-step that improves the likelihood a bit (e.g. gradient step). Recall the IRLS step in the mixture of experts model.

A Report Card for EM



- Some good things about EM:
 - no learning rate (step-size) parameter
 - automatically enforces parameter constraints
 - very fast for low dimensions
 - each iteration guaranteed to improve likelihood
- Some bad things about EM:
 - can get stuck in local minima
 - can be slower than conjugate gradient (especially near convergence)
 - requires expensive inference step
 - is a maximum likelihood/MAP method



$z?$ $p(z|x)$