

*iterative alg. for MLE*

# Expectation Maximization

Machine Learning – CSE546

Emily Fox

University of Washington

November 6, 2013

©Carlos Guestrin 2005-2013

1

## Iterative Algorithm

- Motivates a coordinate ascent-like algorithm:

1. Infer missing values  $z^i$  given estimate of parameters  $\hat{\theta}$
2. Optimize parameters to produce new  $\hat{\theta}$  given "filled in" data  $z^i$
3. Repeat

- Example: MoG (derivation soon... + HW)

1. Infer "responsibilities"

*soft weights*

$$r_{ik} = p(z^i = k | x^i, \hat{\theta}^{(t-1)}) = \frac{\pi_k^{(t-1)} p(x^i | \phi_k^{(t-1)})}{\sum_j \pi_j^{(t-1)} p(x^i | \phi_j^{(t-1)})}$$

*prev. iter.*

2. Optimize parameters

max w.r.t.  $\pi_k$ :

$$\hat{\pi}_k^{(t)} = \frac{1}{N} \sum_i r_{ik} = \frac{r_k}{N} \leftarrow \text{soft counts!}$$

max w.r.t.  $\mu_k, \Sigma_k$ :

$$\hat{\mu}_k^{(t)} = \frac{\sum_i r_{ik} x_i}{\sum_i r_{ik}} \leftarrow \text{weighted mean}$$

$$\hat{\Sigma}_k^{(t)} = \frac{1}{r_k} \left( \sum_i r_{ik} x_i x_i^T - \frac{(\sum_i r_{ik} x_i)^2}{r_k} \right)$$

*weighted covariance*

©Emily Fox 2013

2

# Expectation Maximization (EM) – Setup

- More broadly applicable than just to mixture models considered so far

- Model:  $x$ , observable – “incomplete” data  
 $y$  not (fully) observable – “complete” data  
 $\theta$  parameters

- Interested in maximizing (wrt  $\theta$ ):

$$\max_{\theta} p(x | \theta) = \sum_y p(x, y | \theta) = \sum_y p(x|y, \theta) p(y|\theta)$$

- Special case:

$$x = g(y)$$

$$\text{e.g. } y = \begin{bmatrix} z \\ x \end{bmatrix}$$

non-invertible, deterministic fn  
 class labels  
 obs.  
 in standard mix. models

# Expectation Maximization (EM) – Derivation

- Step 1

- Rewrite desired likelihood in terms of complete data terms

$$p(y | \theta) = p(y | x, \theta) p(x | \theta)$$

$$\Rightarrow \log p(x|\theta) = \log p(y|\theta) - \log p(y|x, \theta)$$

- Step 2

- Assume estimate of parameters  $\hat{\theta}$

- Take expectation with respect to  $p(y | x, \hat{\theta})$

$$L_x(\theta) = E[\log p(y|\theta) | x, \hat{\theta}] + E[-\log p(y|x, \theta) | x, \hat{\theta}]$$

$$U(\theta, \hat{\theta}) \quad V(\theta, \hat{\theta})$$

# Expectation Maximization (EM) – Derivation

## Step 3

- Consider log likelihood of data at any  $\theta$  relative to log likelihood at  $\hat{\theta}$

acknowledge of moving to  $\theta$  when  $\hat{\theta}$   $\leftarrow$  want  $\geq 0$  as long as  $u(\theta, \hat{\theta}) \geq u(\hat{\theta}, \hat{\theta})$  making progress.  $\leq \geq 0$  ✓

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] + [V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]$$

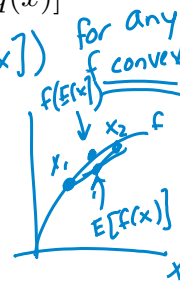
- Aside: Gibbs Inequality**  $E_p[\log p(x)] \geq E_p[\log q(x)]$

Proof: Use Jensen's Ineq.  $E[F(x)] \leq f(E[x])$  for any  $f$  convex

Here:

$$E_p[\log q] - E_p[\log p] = E_p\left[\log \frac{q}{p}\right]$$

$$\leq \log E_p\left[\frac{q}{p}\right] = \log \int_x p(x) \frac{q(x)}{p(x)} dx = 0$$



©Emily Fox 2013

5

# Expectation Maximization (EM) – Derivation

pedal going down hill

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] + [V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]$$

## Step 4

- Determine conditions under which log likelihood at  $\theta$  exceeds that at  $\hat{\theta}$

Using Gibbs inequality:

$$V(\theta, \hat{\theta}) = E[-\log p(y|x, \theta) | x, \hat{\theta}] \geq E[-\log p(y|x, \hat{\theta}) | x, \hat{\theta}]$$

$$= V(\hat{\theta}, \hat{\theta}) \quad \forall \theta$$

If  $U(\theta, \hat{\theta}) \geq U(\hat{\theta}, \hat{\theta})$

Then  $L_x(\theta) \geq L_x(\hat{\theta})$

making progress

choosing  $\theta$  s.t. this is true means we're moving in the right direction (or at least not wrong)

©Emily Fox 2013

6

# Motivates EM Algorithm

- Initial guess:  $\hat{\theta}^{(0)}$
- Estimate at iteration  $t$ :  $\hat{\theta}^{(t)}$
- E-Step**  
 Compute  $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y|\theta) | x, \hat{\theta}^{(t)}]$   
*get function  $U(\theta, \hat{\theta}^{(t)})$*   
*max function  $U$  (pascal forwards)*
- M-Step**  
 Compute  $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$   
*From before,  $U(\hat{\theta}^{(t+1)}, \hat{\theta}^{(t)}) \geq U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)})$*   
 *$\Rightarrow L_x(\hat{\theta}^{(t+1)}) \geq L_x(\hat{\theta}^{(t)})$*

©Emily Fox 2013

7

# Example – Mixture Models

- $\log p(y|\theta) = \log \prod_i p(x^i, z^i|\theta) = \sum_i \log p(x^i, z^i|\theta) = \sum_i \log \prod_k \pi_k P(x^i|\phi_k) P(z^i=k|\pi, \mu, \Sigma)$
- E-Step** Compute  $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y|\theta) | x, \hat{\theta}^{(t)}]$  *Computing  $r_{ik}$*
  - M-Step** Compute  $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$  *estimating  $\pi, \mu, \Sigma$*
  - Consider  $y^i = \{z^i, x^i\}$  i.i.d.  $K$   
 $p(x^i, z^i | \theta) = \pi_{z^i} p(x^i | \phi_{z^i}) = \prod_{k=1}^K (\pi_k P(x^i | \phi_k))^{\mathbb{1}(z^i=k)}$   
 $\max_{\pi, \mu, \Sigma} E_{q_t}[\log p(y|\theta)] = \sum_i E_{q_t}[\log p(x^i, z^i|\theta)] = \sum_i E[\sum_k \mathbb{1}(z^i=k) \log \pi_k P(x^i|\phi_k)]$   
 $= \sum_{i=1}^N \sum_{k=1}^K \underbrace{E[\mathbb{1}(z^i=k) | x^i, \hat{\theta}^{(t)}]}_{P(z^i=k | x^i, \hat{\theta}^{(t)}) = r_{ik}} \log \pi_k P(x^i|\phi_k) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k P(x^i|\phi_k)$   
*weighted log likelihood same as optimizing over weighted data*  
 e.g.,  $\max_{\pi} \Rightarrow \pi_k = \frac{\sum_i r_{ik}}{N}$

©Emily Fox 2013

8

# Coordinate Ascent Behavior

$$V(\theta, \theta^{(t)}) - V(\theta^{(t)}, \theta^{(t)}) \geq 0$$

$$V(\theta, \theta^{(t)}) \geq V(\theta^{(t)}, \theta^{(t)})$$

- Bound log likelihood:

$$\max_{\theta} L_x(\theta) = U(\theta, \hat{\theta}^{(t)}) + V(\theta, \hat{\theta}^{(t)})$$

$$\geq U(\theta, \hat{\theta}^{(t)}) + V(\theta^{(t)}, \hat{\theta}^{(t)}) \equiv LB_x(\theta, \hat{\theta}^{(t)})$$

*max a concave lower bound*

*bound tight at  $\hat{\theta}^{(t)}$*

$$L_x(\hat{\theta}^{(t)}) = U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) + V(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) = LB_x(\hat{\theta}^{(t)}, \hat{\theta}^{(t)})$$

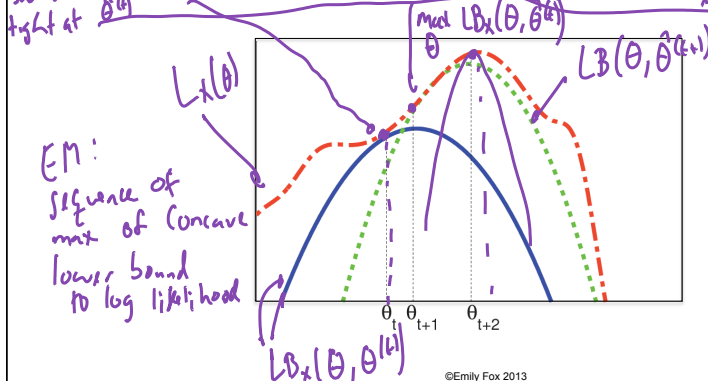


Figure from KM textbook

©Emily Fox 2013

9

# Comments on EM

- Since Gibbs inequality is satisfied with equality only if  $p=q$ , any step that changes  $\theta$  should strictly **increase likelihood**  
*or converged to likelihood doesn't change*
- In practice, can replace the **M-Step** with increasing  $U$  instead of maximizing it (**Generalized EM**) *e.g., a gradient step*
- Under certain conditions (e.g., in exponential family), can show that EM **converges to a stationary point** of  $L_x(\theta)$   
*latent variables*
- Often there is a **natural choice for  $y$**  ... has physical meaning  
 *$y := (\text{cluster assignment}, x)$*
- If you want to choose any  $y$ , not necessarily  $x=g(y)$ , replace  $p(y | \theta)$  in  $U$  with  $p(y, x | \theta)$

©Emily Fox 2013

10

## Initialization

- In mixture model case where  $y^i = \{z^i, x^i\}$  there are many ways to initialize the EM algorithm
- Examples: *K-means*
  - Choose K observations at random to define each cluster. Assign other observations to the nearest "centroid" to form initial parameter estimates
  - Pick the centers sequentially to provide good coverage of data
  - Grow mixture model by splitting (and sometimes removing) clusters until K clusters are formed
- Can be quite important to convergence rates in practice  
*or quality of solution*

©Emily Fox 2013

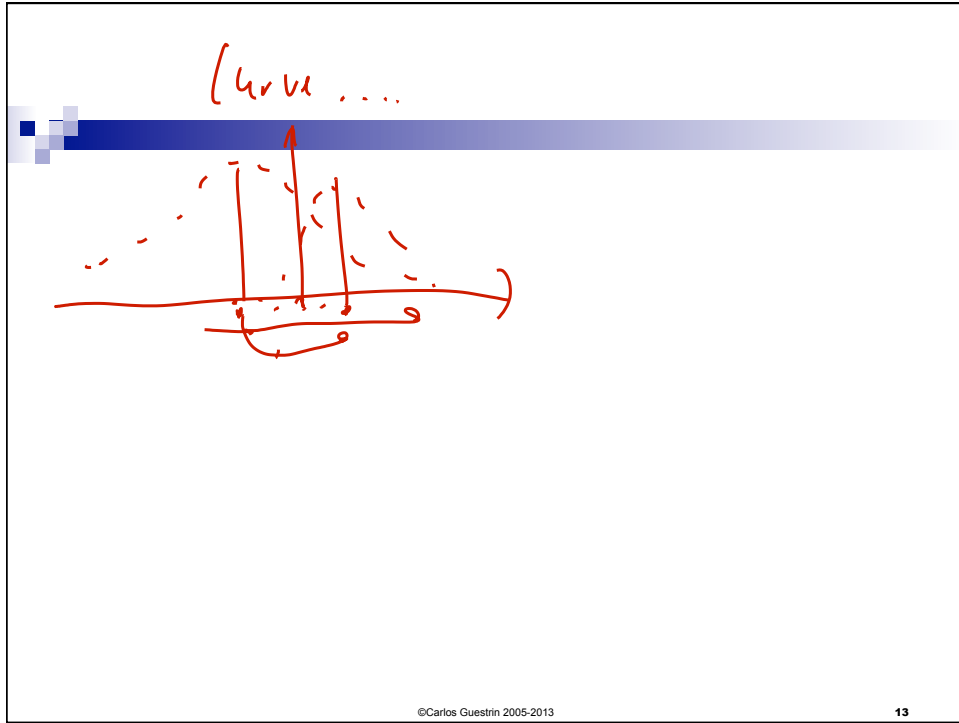
11

## What you should know

- K-means for clustering:
  - algorithm
  - converges because it's coordinate ascent
- EM for mixture of Gaussians:
  - How to "learn" maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent

©Carlos Guestrin 2005-2013

12



# Dimensionality Reduction PCA

Machine Learning – CSE4546  
Carlos Guestrin  
University of Washington

November 13, 2013

©Carlos Guestrin 2005-2013 14

# Dimensionality reduction

- Input data may have thousands or millions of dimensions!

- e.g., text data has

*x*  
↳ 10,000 - 100,000,000 dimensions

- Dimensionality reduction: represent data with fewer dimensions

- easier learning – fewer parameters
  - visualization – hard to visualize more than 3D or 4D
  - discover “intrinsic dimensionality” of data
    - high dimensional data that is truly lower dimensional

©Carlos Guestrin 2005-2013

# Lower dimensional projections

- Rather than picking a subset of the features, we can new features that are combinations of existing features

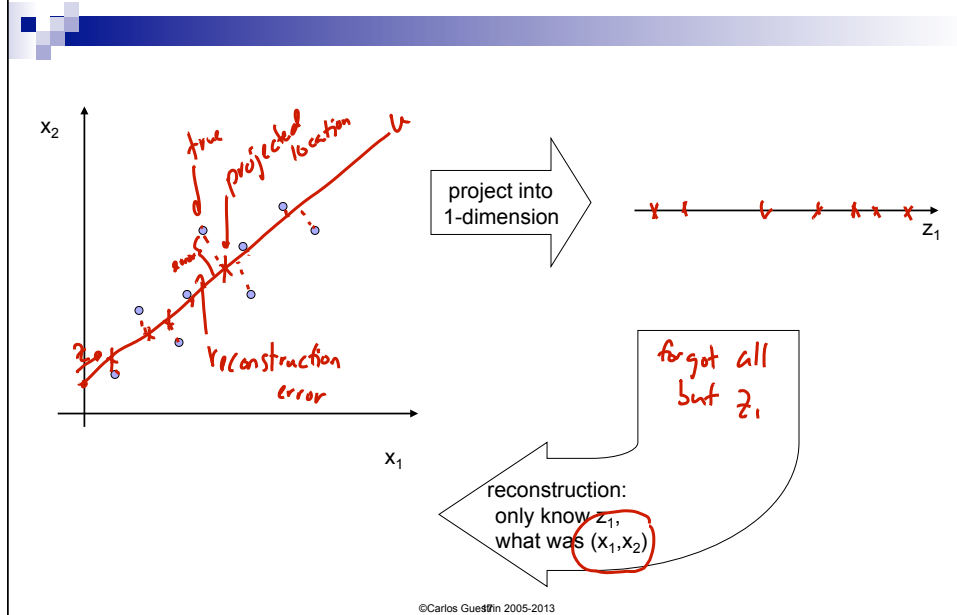
$$z_7 = 2.5x_1 - 2.9x_2 + 3.1x_3 \dots$$
  
$$z = kAx$$
  
*len(A) from data*  
*min reconstruction error*

- Let's see this in the unsupervised setting
  - just  $X$ , but no  $Y$

©Carlos Guestrin 2005-2013



## Linear projection and reconstruction



## Principal component analysis – basic idea

- Project  $d$ -dimensional data into  $k$ -dimensional space while preserving information:
  - e.g., project space of 10000 words into 3-dimensions
  - e.g., project 3-d into 2-d
- Choose projection with minimum reconstruction error

©Carlos Guestrin 2005-2013



# Understanding the reconstruction error

*u will be ordered such that  $u_1$  more "important" than  $u_2$  ...*

*approx  $\rightarrow$*

$$\hat{x}^i = \bar{x} + \sum_{j=1}^k z_j^i u_j \quad \leftarrow k < d$$

$$z_j^i = (x^i - \bar{x}) \cdot u_j$$

*Given  $k \ll d$ , find  $(u_1, \dots, u_k)$  minimizing reconstruction error:*

$$error_k = \sum_{i=1}^N (x^i - \hat{x}^i)^2$$

- Note that  $x^i$  can be represented exactly by  $d$ -dimensional projection:

$$x^i = \bar{x} + \sum_{j=1}^d z_j^i u_j$$

- Rewriting error:

$$error_k = \sum_{i=1}^N (x^i - \hat{x}^i)^2 = \sum_{i=1}^N \left( \bar{x} + \sum_{j=1}^d z_j^i u_j - \left[ \bar{x} + \sum_{j=1}^k z_j^i u_j \right] \right)^2$$

*error is part ignored*

$$= \sum_{i=1}^N \left[ \sum_{j=k+1}^d z_j^i u_j \cdot u_j z_j^i + \sum_{j=k+1}^d \sum_{\ell \neq j} z_j^i u_j \cdot u_\ell z_\ell^i \right]$$

$$= \sum_{i=1}^N \sum_{j=k+1}^d (z_j^i)^2 \leftarrow \text{min error} \equiv \text{min projection onto ignored directions}$$

*$z_j^i = (x^i - \bar{x}) \cdot u_j$*

©Carlos Guestrin 2005-2013

# Reconstruction error and covariance matrix

$$error_k = \sum_{i=1}^N \sum_{j=k+1}^d [u_j \cdot (x^i - \bar{x})]^2$$

$$= \sum_{i=1}^N \sum_{j=k+1}^d u_j^T (x^i - \bar{x}) (x^i - \bar{x})^T u_j$$

$$= \sum_{j=k+1}^d u_j^T \left[ \sum_{i=1}^N (x^i - \bar{x}) (x^i - \bar{x})^T \right] u_j$$

*$N \Sigma$*

min error:

$$\min_u \sum_{j=k+1}^d u_j^T \Sigma u_j$$

*kind  $u_j$  that minimize error*

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (x^i - \bar{x})(x^i - \bar{x})^T$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & & \\ & \sigma_{jj} & \\ & & \end{pmatrix}$$

$$\sigma_{uv} = \frac{1}{N} \sum_{i=1}^N (x_u^i - \bar{x}_u)(x_v^i - \bar{x}_v)$$

in vector format

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (x^i - \bar{x})(x^i - \bar{x})^T$$

©Carlos Guestrin 2005-2013