

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

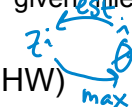
- Example: MoG (derivation soon... + HW)

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.



$$\hat{\theta} = \{\hat{\mu}_k, \hat{\Sigma}_k\}$$

2. Optimize parameters

max w.r.t. π_k :

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

max w.r.t. μ_k, Σ_k :

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

$$\hat{\Sigma}_k^{(t)} = \frac{1}{r_k} \sum r_{ik} x_i x_i^T - \hat{\mu}_k^{(t)} \hat{\mu}_k^{(t)T}$$

©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
 θ parameters

- Interested in maximizing (wrt θ):

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

- 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

what we have
what we wish we had

introduce \rightarrow

Expectation Maximization (EM) – Derivation

- Step 1

- Rewrite desired likelihood in terms of complete data terms

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

$$\Rightarrow \underbrace{\log p(x|\theta)}_{L_x(\theta)} = \log p(y|\theta) - \log p(y|x, \theta)$$

- Step 2

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

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

“ $E[\cdot | x, \hat{\theta}]$ ”

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

Expectation Maximization (EM) – Derivation

Step 3

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

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

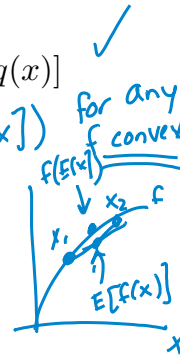
- Aside: Gibbs Inequality** $E_p[\log p(x)] \geq E_p[\log q(x)]$ ✓

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

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] + \underbrace{[V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]}_{\geq 0}$$

Step 4

- Determine conditions under which log likelihood at θ 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

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

choosing θ 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 $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y|\theta) | x, \hat{\theta}^{(t)}]$

- **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

- **E-Step** Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y | \theta) | x, \hat{\theta}^{(t)}]$
- **M-Step** Compute $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$

- Consider $y^i = \{z^i, x^i\}$ i.i.d.

$$p(x^i, z^i | \theta) = \pi_{z^i} p(x^i | \phi_{z^i}) =$$

$$E_{q_t}[\log p(y | \theta)] = \sum_i E_{q_t}[\log p(x^i, z^i | \theta)] =$$

©Emily Fox 2013

8

Coordinate Ascent Behavior

- Bound log likelihood:

$$\begin{aligned} L_x(\theta) &= U(\theta, \hat{\theta}^{(t)}) + V(\theta, \hat{\theta}^{(t)}) \\ &\geq \\ L_x(\hat{\theta}^{(t)}) &= U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) + V(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) \end{aligned}$$

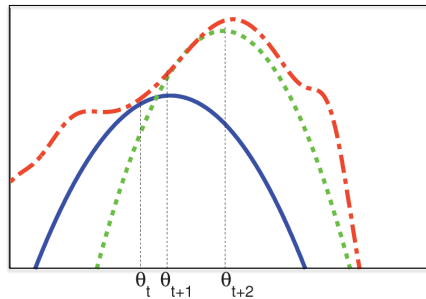


Figure from
KM textbook

©Emily Fox 2013

9

Comments on EM

- Since Gibbs inequality is satisfied with equality only if $p=q$, any step that changes θ should strictly **increase likelihood**
- In practice, can replace the **M-Step** with increasing U instead of maximizing it (**Generalized EM**)
- Under certain conditions (e.g., in exponential family), can show that EM **converges to a stationary point** of $L_x(\theta)$
- Often there is a **natural choice for y** ... has physical meaning
- 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:
 - 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

©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

13

Dimensionality reduction



- Input data may have thousands or millions of dimensions!
 - e.g., text data has
- **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

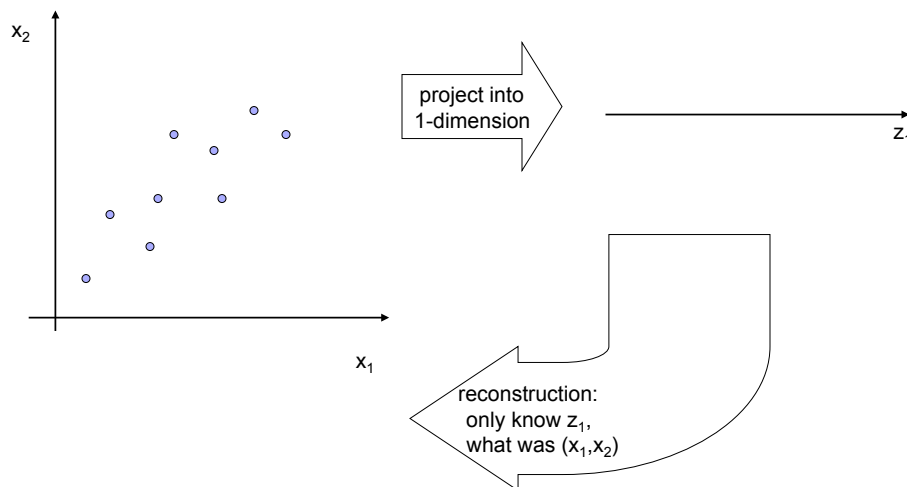
Lower dimensional projections

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

- Let's see this in the unsupervised setting
 - just \mathbf{X} , but no \mathbf{Y}

©Carlos Guestrin 2005-2013

Linear projection and reconstruction



©Carlos Guestrin 2005-2013

Principal component analysis – basic idea

- Project n-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

Linear projections, a review

- Project a point into a (lower dimensional) space:
 - **point**: $\mathbf{x} = (x_1, \dots, x_d)$
 - **select a basis** – set of basis vectors – $(\mathbf{u}_1, \dots, \mathbf{u}_k)$
 - we consider orthonormal basis:
 - $\mathbf{u}_i \bullet \mathbf{u}_i = 1$, and $\mathbf{u}_i \bullet \mathbf{u}_j = 0$ for $i \neq j$
 - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
 - **best coordinates** in lower dimensional space defined by dot-products: (z_1, \dots, z_k) , $z_i = (\mathbf{x} - \bar{\mathbf{x}}) \bullet \mathbf{u}_i$
 - minimum squared error

©Carlos Guestrin 2005-2013

PCA finds projection that minimizes reconstruction error

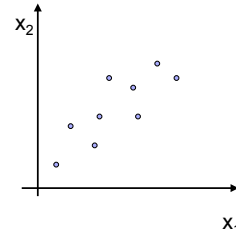
- Given N data points: $\mathbf{x}^i = (x_1^i, \dots, x_d^i)$, $i=1 \dots N$
- Will represent each point as a projection:

$$\square \hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j \quad \text{where: } \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \quad \text{and} \quad z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- PCA:

- Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

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



©Carlos Guestrin 2005-2013

Understanding the reconstruction error

- Note that \mathbf{x}^i can be represented exactly by d-dimensional projection:

$$\mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^d z_j^i \mathbf{u}_j$$

- Rewriting error:

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

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

©Carlos Guestrin 2005-2013

Reconstruction error and covariance matrix

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

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

©Carlos Guestrin 2005-2013

Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_d)$ minimizing:

$$error_k = \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

- Eigen vector:
- Minimizing reconstruction error equivalent to picking $(\mathbf{u}_{k+1}, \dots, \mathbf{u}_d)$ to be eigen vectors with smallest eigen values

©Carlos Guestrin 2005-2013

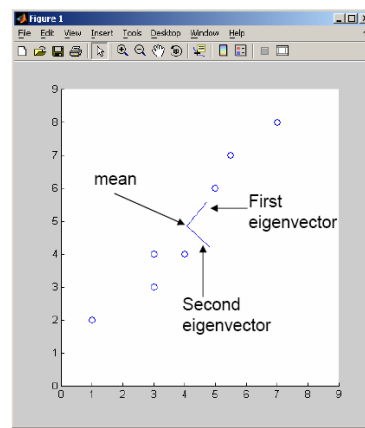
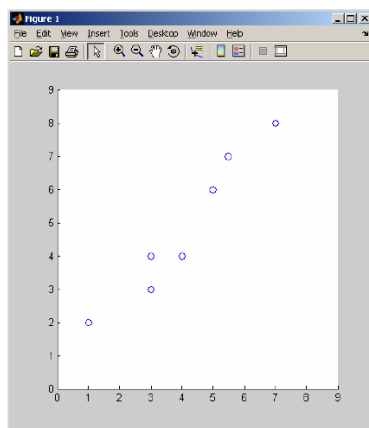
Basic PCA algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Compute covariance matrix**:
 - $\Sigma \leftarrow 1/N \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of Σ
- **Principal components**: k eigen vectors with highest eigen values

©Carlos Guestrin 2005-2013

PCA example

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

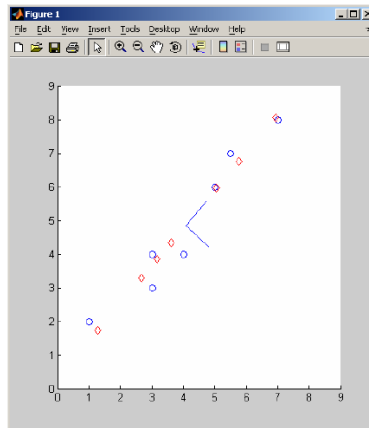
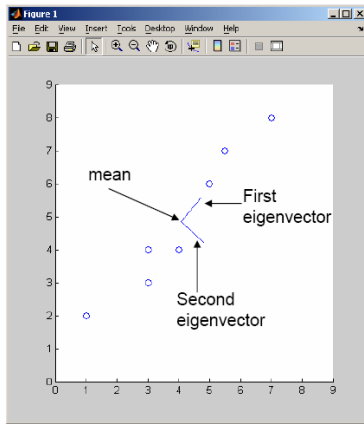


©Carlos Guestrin 2005-2013

PCA example – reconstruction

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

only used first principal component



©Carlos Guestrin 2005-2013

Eigenfaces [Turk, Pentland '91]

■ Input images:



■ Principal components:



©Carlos Guestrin 2005-2013

Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



©Carlos Guestrin 2005-2013

Scaling up

- Covariance matrix can be really big!
 - Σ is d by d
 - Say, only 10000 features
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - finds to k eigenvectors
 - great implementations available, e.g., python, R, Matlab svd

©Carlos Guestrin 2005-2013

SVD

- Write $\mathbf{X} = \mathbf{W} \mathbf{S} \mathbf{V}^T$
 - \mathbf{X} ← data matrix, one row per datapoint
 - \mathbf{W} ← weight matrix, one row per datapoint – coordinate of \mathbf{x}^i in eigenspace
 - \mathbf{S} ← singular value matrix, diagonal matrix
 - in our setting each entry is eigenvalue λ_j
 - \mathbf{V}^T ← singular vector matrix
 - in our setting each row is eigenvector \mathbf{v}_j

©Carlos Gue@in 2005-2013

PCA using SVD algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- Call SVD algorithm on \mathbf{X}_c – ask for k singular vectors
- **Principal components**: k singular vectors with highest singular values (rows of \mathbf{V}^T)
 - **Coefficients** become:

©Carlos Gue@in 2005-2013

What you need to know

- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD