

# Expectation Maximization

Machine Learning – CSE546

Carlos Guestrin

University of Washington

November 13, 2014

©Carlos Guestrin 2005-2014

1

## E.M.: The General Case

- E.M. widely used beyond mixtures of Gaussians
  - The recipe is the same...
- Expectation Step: Fill in missing data, given current values of parameters,  $\theta^{(t)}$ 
  - If variable  $y$  is missing (could be many variables)
  - Compute, for each data point  $\mathbf{x}^i$ , for each value  $i$  of  $z$ :
    - $P(z=i|\mathbf{x}^i, \theta^{(t)})$
- Maximization step: Find maximum likelihood parameters for (weighted) “completed data”:
  - For each data point  $\mathbf{x}^i$ , create  $k$  weighted data points
    -
  - Set  $\theta^{(t+1)}$  as the maximum likelihood parameter estimate for this weighted data
- Repeat

©Carlos Guestrin 2005-2014

2

# The general learning problem with missing data

- Marginal likelihood –  $\mathbf{x}$  is observed,  $\mathbf{z}$  is missing:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log \prod_{j=1}^m P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j | \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} | \theta)\end{aligned}$$

©Carlos Guestrin 2005-2014

3

## E-step

- $\mathbf{x}$  is observed,  $\mathbf{z}$  is missing
- Compute probability of missing data given current choice of  $\theta$ 
  - $Q(\mathbf{z} | \mathbf{x}_j)$  for each  $\mathbf{x}_j$ 
    - e.g., probability computed during classification step
    - corresponds to “classification step” in K-means

$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

©Carlos Guestrin 2005-2014

4

## Jensen's inequality

$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{z} | \mathbf{x}_j) P(\mathbf{x}_j | \theta)$$

- **Theorem:**  $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

## Applying Jensen's inequality

- Use:  $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\theta^{(t)} : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)}$$

## The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Corresponds to weighted dataset:

- $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_1)$
- $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_1)$
- $\langle \mathbf{x}_1, \mathbf{z}=3 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_1)$
- $\langle \mathbf{x}_2, \mathbf{z}=1 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=1 | \mathbf{x}_2)$
- $\langle \mathbf{x}_2, \mathbf{z}=2 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=2 | \mathbf{x}_2)$
- $\langle \mathbf{x}_2, \mathbf{z}=3 \rangle$  with weight  $Q^{(t+1)}(\mathbf{z}=3 | \mathbf{x}_2)$
- ...

©Carlos Guestrin 2005-2014

7

## The M-step

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

- Use expected counts instead of counts:

- If learning requires  $\text{Count}(\mathbf{x}, \mathbf{z})$
- Use  $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{z})]$

©Carlos Guestrin 2005-2014

8

# Convergence of EM

- Define potential function  $F(\theta, Q)$ :

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on  $F$ 
  - Thus, maximizes lower bound on marginal log likelihood
  - We saw that M-step corresponds to fixing  $Q$ , max  $\theta$
  - E-step fix  $\theta$  and max  $Q$

©Carlos Guestrin 2005-2014

9

# 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
- General case for EM

©Carlos Guestrin 2005-2014

16



# Dimensionality Reduction PCA

Machine Learning – CSE4546

Carlos Guestrin

University of Washington

November 13, 2014

©Carlos Guestrin 2005-2014

17

## 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

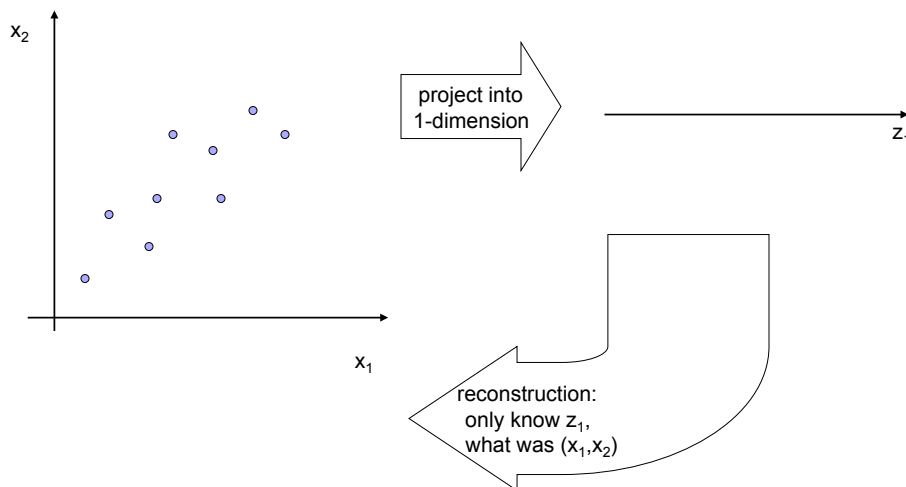
©Carlos Guestrin 2005-2014

## 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 2005-2014

## Linear projection and reconstruction



©Carlos Guestrin 2005-2014

## 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 2005-2014

## 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 2005-2014



## PCA finds projection that minimizes reconstruction error

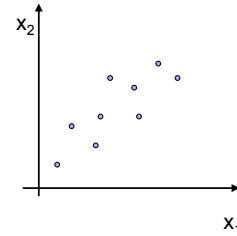
- Given N data points:  $\mathbf{x}^i = (x_1^i, \dots, x_d^i)$ ,  $i=1 \dots N$
- Will represent each point as a projection:

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

- 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 2005-2014

## 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 2005-2014

## 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 2005-2014

## 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 2005-2014

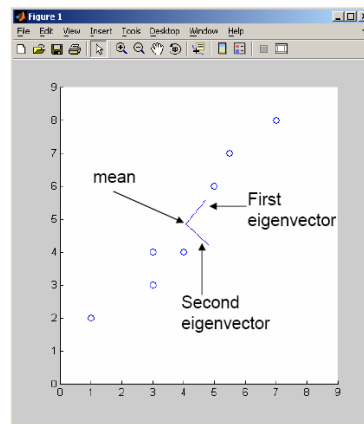
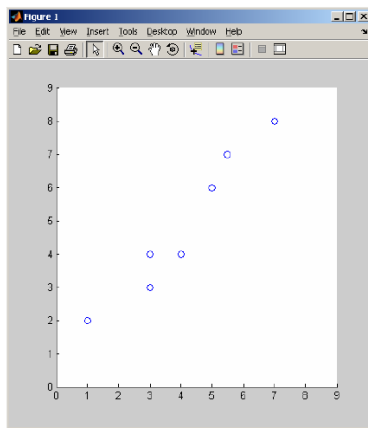
# Basic PCA algorithm

- 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  $\Sigma$
- **Principal components**:  $k$  eigen vectors with highest eigen values

©Carlos Guestrin 2005-2014

# PCA example

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

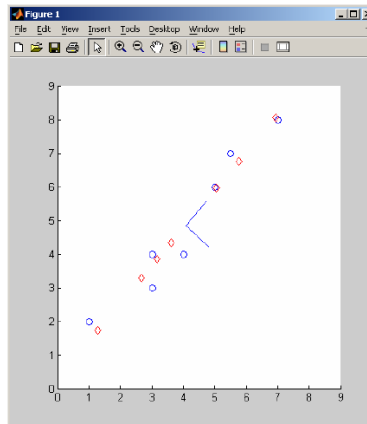
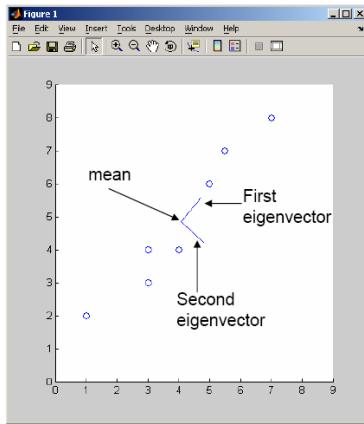


©Carlos Guestrin 2005-2014

# PCA example – reconstruction

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

only used first principal component



©Carlos Guestrin 2005-2014

# Eigenfaces [Turk, Pentland '91]

■ Input images:



■ Principal components:



©Carlos Guestrin 2005-2014

## Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



©Carlos Gue@tin 2005-2014

## Scaling up

- Covariance matrix can be really big!
  - $\Sigma$  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 Gue@tin 2005-2014

# 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  $\lambda_j$
  - $\mathbf{V}^T$  ← singular vector matrix
    - in our setting each row is eigenvector  $\mathbf{v}_j$

©Carlos Gue#in 2005-2014

# 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 2005-2014

# 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