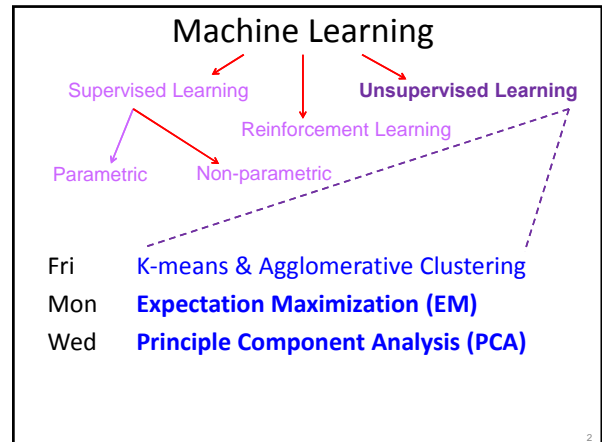


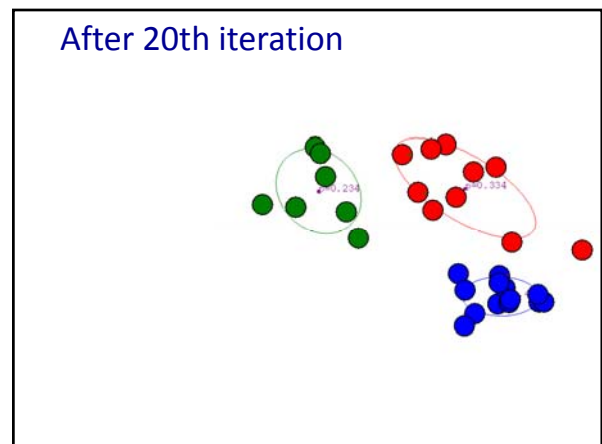
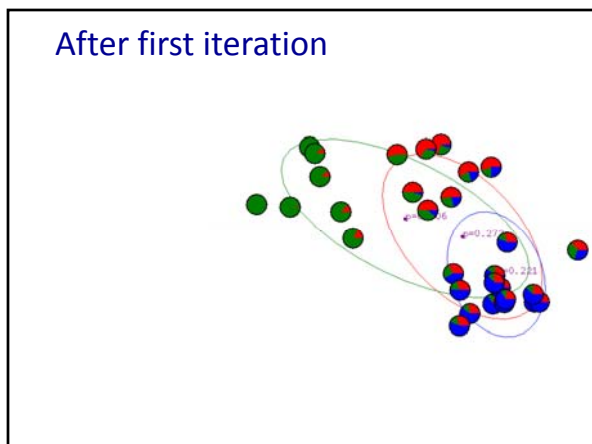
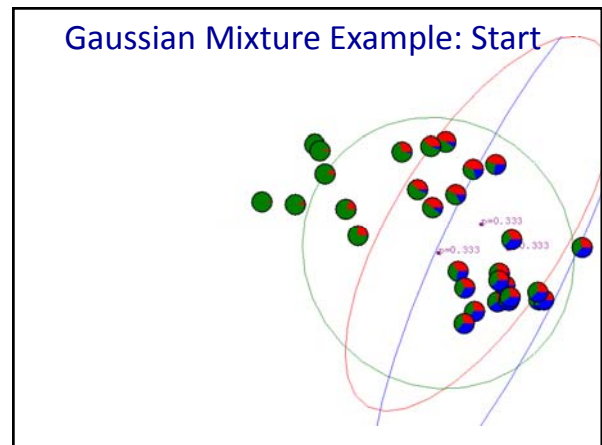
CSE 446  
Dimensionality Reduction and PCA  
Winter 2012

Slides adapted from Carlos Guestrin & Luke Zettlemoyer



**EM**  
another iterative clustering algorithm

- Pick K random cluster models
- Alternate:
  - Assign data instances *proportionately* to different models
  - Revise each cluster model based on its (proportionately) assigned points
- Stop when no changes



## Can we do EM with hard assignments?

(univariate case for simplicity)

**Iterate:** On the  $t$ 'th iteration let our estimates be

$$\theta_t = \{\mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}\}$$

**E-step**

Compute "expected" classes of all datapoints

$$p(y=i|x_j, \mu_1, \dots, \mu_k) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) p(y=i)$$

**M-step**

Compute most likely new  $\mu$ s given class expectations

$$\mu_i = \frac{\sum_{j=1}^m P(y=i|x_j) x_j}{\sum_{j=1}^m P(y=i|x_j)}$$

$$\mu_i = \frac{\sum_{j=1}^m \delta(y=i, x_j) x_j}{\sum_{j=1}^m \delta(y=i, x_j)}$$

$\delta$  represents hard assignment to "most likely" or nearest cluster

Equivalent to k-means clustering algorithm!!!

## Summary

- K-means for clustering:
  - algorithm
  - converges because it's coordinate ascent
- Hierarchical agglomerative clustering
- EM for mixture of Gaussians:
  - Viewed as coordinate ascent in parameter space  $(\mu, \Sigma, \pi)$
  - How to "learn" maximum likelihood parameters (locally max. like.) in the case of unlabeled data
  - Relation to K-means
    - Hard / soft clustering
    - Probabilistic model
- Remember, E.M. can (& *does*) get stuck in local minima

## Dimensionality Reduction

Input data may have thousands or millions of dimensions!

e.g., images have??? And text data has ???,

Dnl S. Wild s Thms J. Cbl / WRF Prfssr f Cmpt  
Scnc nd ngngng t th nvrsty f Wshngtn. ftr gng t schl  
t Philips cdmry, h rcvcd bchlr's dgrs n bth Cmpt  
Scnc nd Bchmstry t YI nvrsty n 1982.

Are all those dimensions (letters, pixels) necessary?

## Dimensionality Reduction

- Represent data with fewer dimensions! 😊
- Easier learning – fewer parameters
  - |Features| >> |training examples| ??
- Better visualization
  - Hard to understand more than 3D or 4D
- Discover "intrinsic dimensionality" of data
  - High dimensional data that is truly lower dimensional

## Feature Selection

- Want to learn  $f: X \rightarrow Y$ 
  - $X = \langle X_1, \dots, X_n \rangle$
  - but some features are more important than others
- **Approach:** select subset of features to be used by learning algorithm
  - **Score** each feature (or sets of features)
  - **Select** set of features with best score

## Greedy forward feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- **Greedy:** Start from empty (or simple) set of features  $F_0 = \emptyset$ 
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select **next best feature**  $X_i$ 
    - e.g.,  $X_i$  giving lowest **held-out error** when learning with  $F_t \cup \{X_i\}$
  - $F_{t+1} \leftarrow F_t \cup \{X_i\}$
  - Repeat

## Greedy **backward** feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- Greedy: Start with all features  $F_0 = F$ 
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select **next worst feature**  $X_i$ 
    - $X_i$  giving lowest held out-error learner when learning with  $F_t - \{X_i\}$
  - $F_{t+1} \leftarrow F_t - \{X_i\}$
  - Repeat

## Impact of feature selection on classification of fMRI data [Pereira et al. '05]

Accuracy classifying category of word read by subject

#voxels	mean	subjects							
		233B	329B	332B	424B	474B	696B	775B	867B
50	0.735	0.783	0.817	0.55	0.783	0.75	0.8	0.65	0.75
100	0.742	0.767	0.8	0.533	0.817	0.85	0.783	0.6	0.783
200	0.737	0.783	0.783	0.517	0.817	0.883	0.75	0.583	0.783
<b>300</b>	<b>0.75</b>	<b>0.8</b>	<b>0.817</b>	<b>0.567</b>	<b>0.833</b>	<b>0.883</b>	<b>0.75</b>	<b>0.583</b>	<b>0.767</b>
400	0.742	0.8	0.783	0.583	0.85	0.833	0.75	0.583	0.75
800	0.735	0.833	0.817	0.567	0.833	0.833	0.7	0.55	0.75
1600	0.698	0.8	0.817	0.45	0.783	0.833	0.633	0.5	0.75
all (~2500)	0.638	0.767	0.767	0.25	0.75	0.833	0.567	0.433	0.733

Table 1: Average accuracy across all pairs of categories, restricting the procedure to use a certain number of voxels for each subject. The highlighted line corresponds to the best mean accuracy, obtained using 300 voxels.

Voxels scored by p-value of regression to predict voxel value from the task

## Feature Selection through Regularization

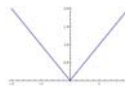
- Previously, we discussed regularization with a squared norm:

$$\hat{\theta} = \arg \min_{\theta} \text{Loss}(\theta; \mathcal{D}) + \lambda \sum_i \theta_i^2$$



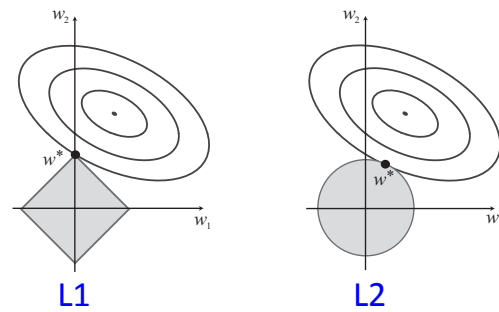
- What if we used an  $L_1$  norm instead?

$$\hat{\theta} = \arg \min_{\theta} \text{Loss}(\theta; \mathcal{D}) + \lambda \sum_i |\theta_i|$$



- What about  $L_\infty$ ?
- These norms work, but are harder to optimize! And, it can be tricky to set  $\lambda$ !!!

## Regularization



16

## Lower Dimensional Projections

- Rather than picking a subset of the features, we can **make new ones** by combining existing features  $x_1 \dots x_n$

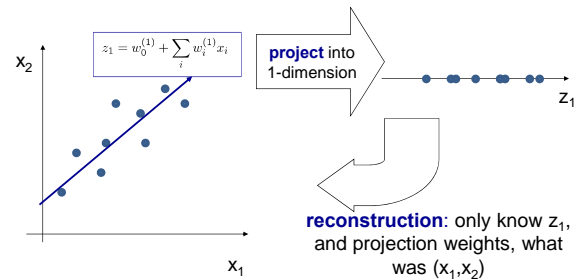
$$z_1 = w_0^{(1)} + \sum_i w_i^{(1)} x_i$$

...

$$z_k = w_0^{(k)} + \sum_i w_i^{(k)} x_i$$

- New features are linear combinations of old ones
- Reduces dimension when  $k < n$
- Let's see this in the **unsupervised setting**
  - just  $X$ , but no  $Y$

## Linear projection and reconstruction



## 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

## Linear projections, a review

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

- Project a point into a (lower dimensional) space:
  - point**:  $x = (x_1, \dots, x_n)$
  - select a basis** – set of unit (length 1) basis vectors  $(u_1, \dots, u_k)$ 
    - we consider orthonormal basis:
      - $u_i \cdot u_i = 1$ , and  $u_i \cdot u_j = 0$  for  $i \neq j$
  - select a center** –  $\bar{x}$ , defines offset of space
  - best coordinates** in lower dimensional space defined by dot-products:  $(z_1, \dots, z_k)$ ,  $z_i = (x - \bar{x}) \cdot u_i$

## 3D Data

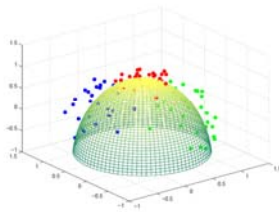


FIGURE 14.15. Simulated data in three classes, near the surface of a half-sphere.

21

## Projection to 2D

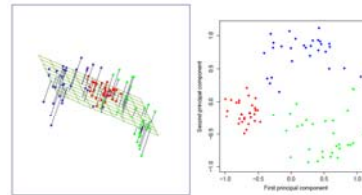


FIGURE 14.21. The best rank-two linear approximation to the half-sphere data. The right panel shows the projected points with coordinates given by  $U_2 D_2$ , the first two principal components of the data.

22

## PCA finds projection that minimizes reconstruction error

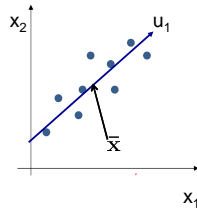
- Given m data points:  $x^i = (x_1^i, \dots, x_n^i)$ ,  $i=1 \dots m$
- Will represent each point as a projection:

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

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x^i$$

- PCA:  $z_j^i = (x^i - \bar{x}) \cdot u_j$ 
  - Given  $k < n$ , find  $(u_1, \dots, u_k)$  minimizing reconstruction error:

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



## Understanding the reconstruction error

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

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

Given  $k < n$ , find  $(u_1, \dots, u_k)$

minimizing reconstruction error:

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

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

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

- Rewriting error:

$$\begin{aligned} error_k &= \sum_{i=1}^m \left( x^i - \left[ \bar{x} + \sum_{j=1}^k z_j^i u_j \right] \right)^2 = \sum_{i=1}^m \left( \left[ \bar{x} + \sum_{j=1}^n z_j^i u_j \right] - \left[ \bar{x} + \sum_{j=1}^k z_j^i u_j \right] \right)^2 \\ &= \sum_{i=1}^m \left( \sum_{j=k+1}^n z_j^i u_j \right)^2 = \sum_{i=1}^m \sum_{j=k+1}^n z_j^i u_j \cdot z_j^i u_j + \sum_{i=1}^m \sum_{j=k+1}^n \sum_{l=k+1, l \neq j}^n z_j^i u_j \cdot u_l z_l^i \\ &= \sum_{i=1}^m \sum_{j=k+1}^n (z_j^i)^2 \end{aligned}$$

Aha! Error is sum of squared weights that would have been used for dimensions that are cut!!!!

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

### Reconstruction error and covariance matrix

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

Now, to find the  $u_j$ , we minimize:

$$u^T \Sigma u + \lambda(1 - u^T u)$$

Lagrange multiplier to ensure orthonormal

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

Take derivative, set equal to 0, ..., solutions are eigenvectors

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

$$error_k = \sum_{j=k+1}^n u_j^T \Sigma u_j$$

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

$$\Sigma u_i = \lambda_i u_i$$

### Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis  $(u_1, \dots, u_n)$  minimizing:

$$error_k = m \sum_{j=k+1}^n u_j^T \Sigma u_j$$

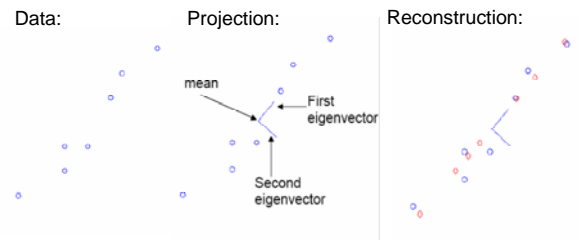
- Solutions: eigen vectors  $\Sigma u_i = \lambda_i u_i$
- So, minimizing reconstruction error equivalent to picking  $(u_{k+1}, \dots, u_n)$  to be eigen vectors with smallest eigen values
- And, our projection should be onto the  $(u_1, \dots, u_k)$  with the largest values

### Basic PCA algorithm

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

### PCA example

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



### Handwriting

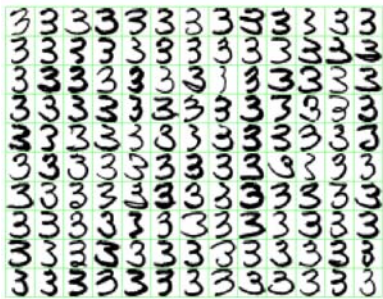
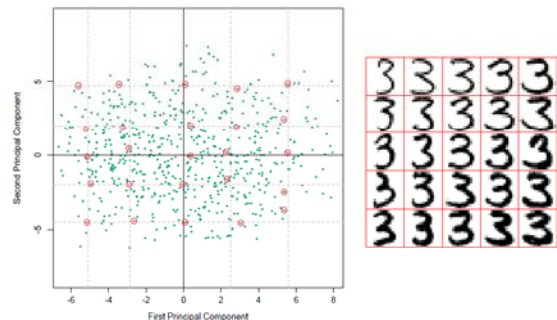


FIGURE 14.22. A sample of 130 handwritten 3's shows a variety of writing styles.

29

### Two Principle Eigenvectors



30

## Eigenfaces [Turk, Pentland '91]

### Input images:

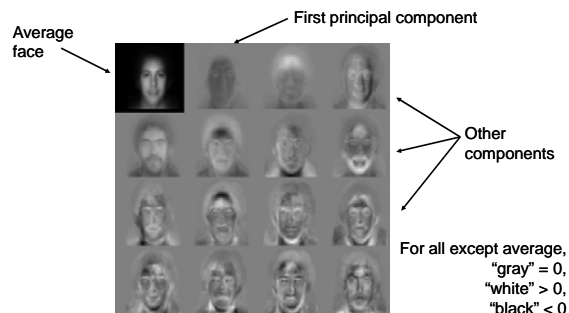
- N images
- Each 50x50 pixels
- 2500 features



Misleading figure.

Best to think of as an  $N \times 2500$  matrix: [Examples]  $\times$  [Features]

## Reduce Dimensionality 2500 $\rightarrow$ 15



## Using PCA for Compression

- Store each face as coefficients of projection onto first few principal components

$$\text{image} = \sum_{i=0}^{i_{\max}} a_i \text{Eigenface}_i$$

## Using PCA for Recognition

- Compute projections of target image, compare to database ("nearest neighbor classifier")

## Eigenfaces reconstruction

- Each image corresponds to adding together the principal components:



## Scaling up

- Covariance matrix can be really big!
  - $\Sigma$  is  $n \times n$
  - 10000 features can be common!
  - Finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
  - Finds top k eigenvectors
  - Great implementations available, e.g., Matlab svd

## SVD

Write  $\mathbf{X} = \mathbf{W} \mathbf{S} \mathbf{V}^T$

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

## SVD

notation change ☹

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U} \end{pmatrix} \begin{pmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{pmatrix} \begin{pmatrix} \mathbf{V} \end{pmatrix}^T$$

- Treat as black box: code widely available  
In Matlab: `[U,W,V]=svd(A,0)`

## 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**: project each point onto the new vectors

## Singular Value Decomposition (SVD)

- Handy mathematical technique that has application to many problems
- Given any  $m \times n$  matrix  $\mathbf{A}$ , algorithm to find matrices  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{W}$  such that
$$\mathbf{A} = \mathbf{U}\mathbf{W}\mathbf{V}^T$$
  - $\mathbf{U}$  is  $m \times n$  and orthonormal
  - $\mathbf{W}$  is  $n \times n$  and diagonal
  - $\mathbf{V}$  is  $n \times n$  and orthonormal

## SVD

- The  $w_i$  are called the **singular values** of  $\mathbf{A}$
- If  $\mathbf{A}$  is singular, some of the  $w_i$  will be 0
- In general  $rank(\mathbf{A}) =$  number of nonzero  $w_i$
- SVD is mostly unique (up to permutation of singular values, or if some  $w_i$  are equal)

## What you need to know

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