

Reading:  
Vapnik 1998  
Joachims 1999 (see class website)

# Transductive SVMs

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April <sup>24</sup>~~17~~<sup>th</sup>, 2006

# Semi-supervised learning and discriminative models

- We have seen semi-supervised learning for generative models

- EM

$$\max \log P(X) = \max \log \sum_y P(X, y)$$

*X is visible  
Y hidden*

- What can we do for discriminative models

- Not regular EM

*$P(Y|X) \rightarrow$  can't compute  $p(X)$*

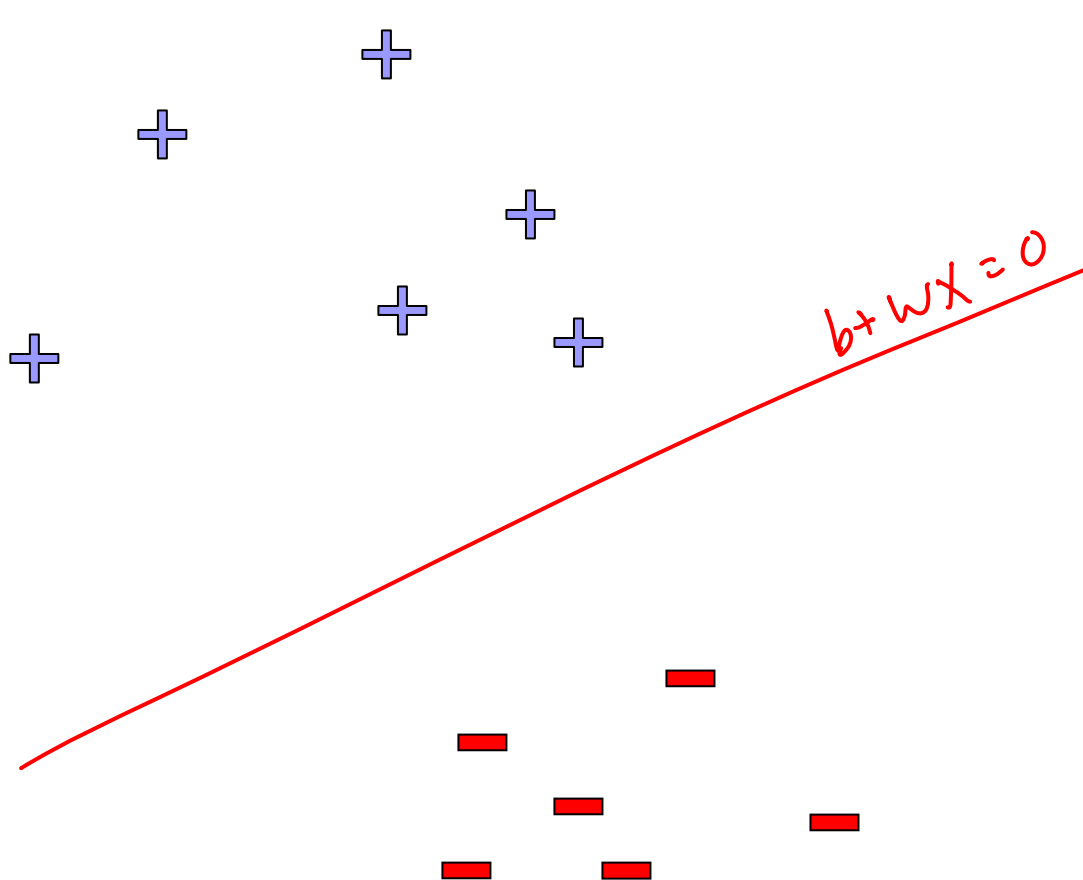
- we can't compute  $P(x)$

- But there are discriminative versions of EM

- Co-Training!

- Many other tricks... let's see an example

# Linear classifiers – Which line is better?



**Data:**

*features*

*class*

$$\begin{aligned} &\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle \\ &\vdots \\ &\langle x_n^{(1)}, \dots, x_n^{(m)}, y_n \rangle \end{aligned}$$

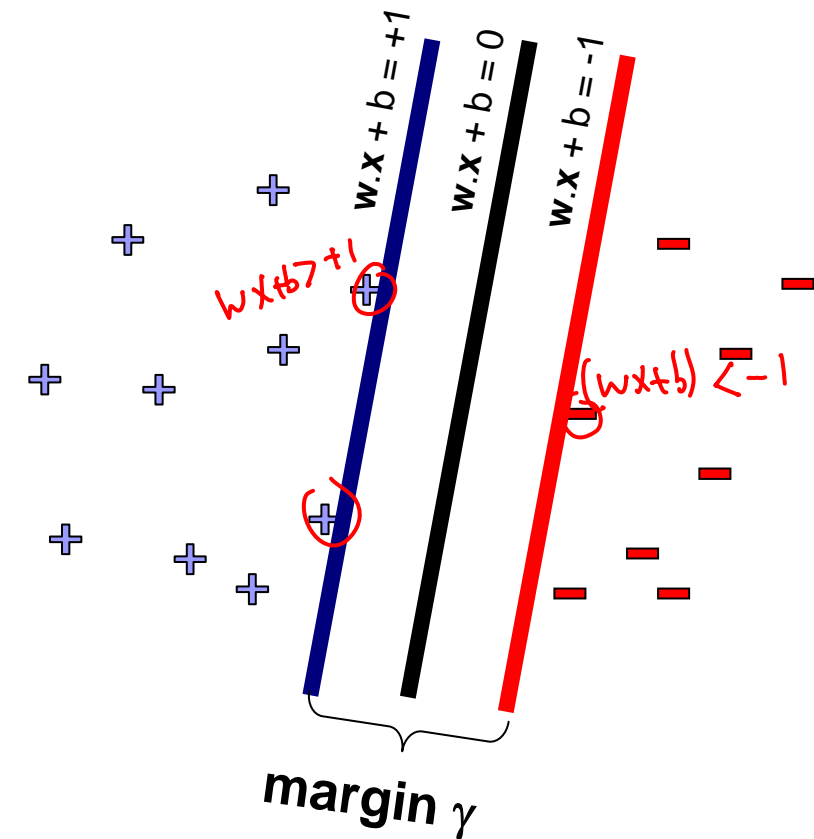
**Example i:**

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \text{ — } m \text{ features}$$

$$y_i \in \{-1, +1\} \text{ — class}$$

$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

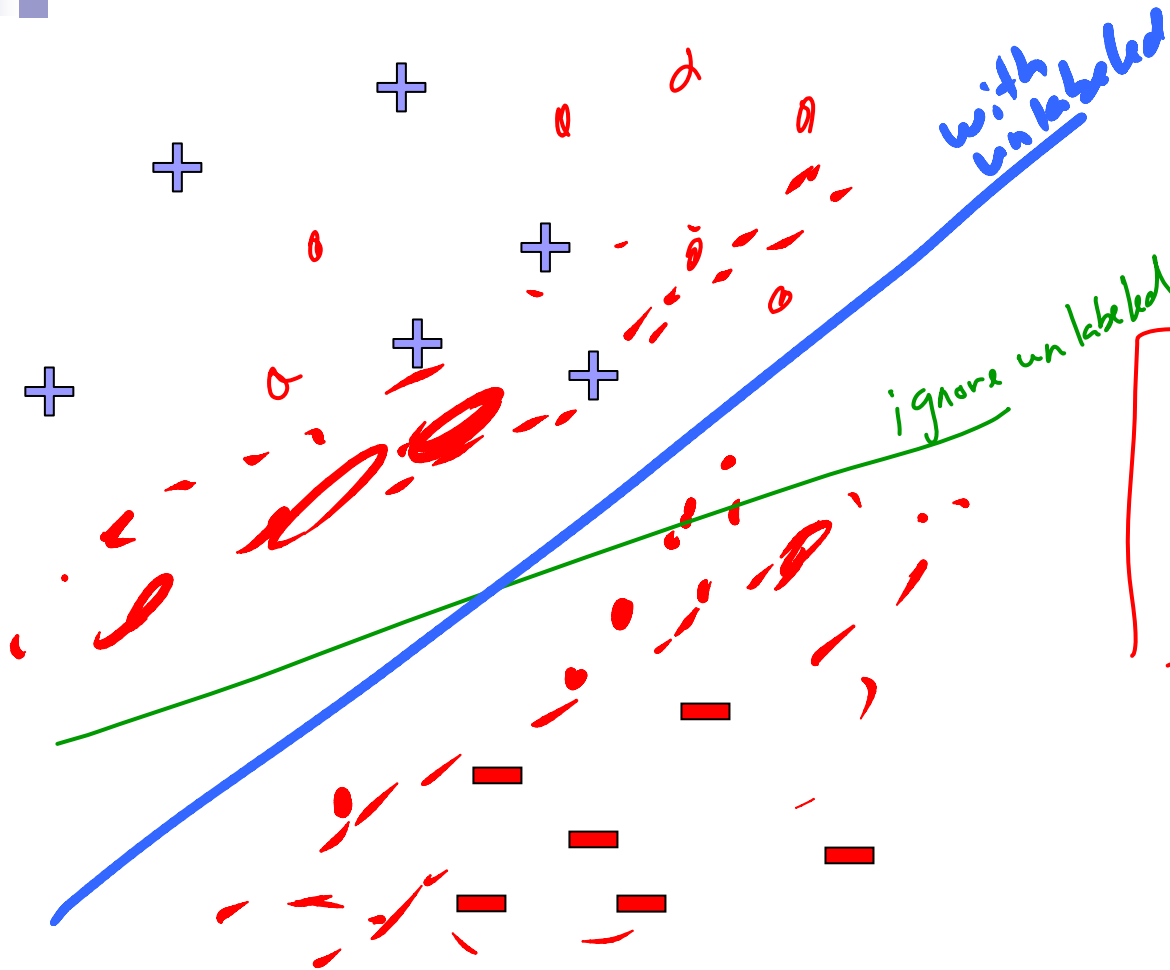
# Support vector machines (SVMs)



$$\text{minimize}_w \quad w \cdot w$$
$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j$$

- Solve efficiently by quadratic programming (QP)
  - Well-studied solution algorithms
- Hyperplane defined by support vectors

# What if we have unlabeled data?



**$n_L$  Labeled Data:**

$$\begin{aligned} &\langle x_1^{(1)}, \dots, x_1^{(m)}, y_1 \rangle \\ &\vdots \\ &\langle x_n^{(1)}, \dots, x_n^{(m)}, y_{n_L} \rangle \end{aligned}$$

**Example i:**

$$\langle x_i^{(1)}, \dots, x_i^{(m)} \rangle \text{ — } m \text{ features}$$

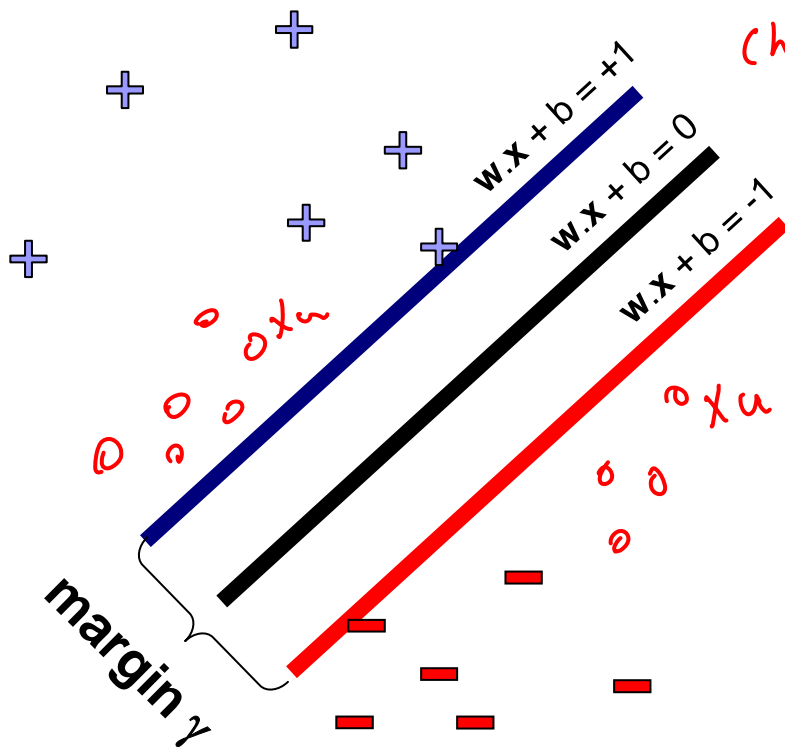
$$y_i \in \{-1, +1\} \text{ — class}$$

**$n_U$  Unlabeled Data:**

$$\begin{aligned} &\langle x_1^{(1)}, \dots, x_1^{(m)}, ? \rangle \\ &\vdots \\ &\langle x_n^{(1)}, \dots, x_n^{(m)}, ? \rangle \end{aligned}$$

$$\mathbf{w} \cdot \mathbf{x} = \sum_j w^{(j)} x^{(j)}$$

# Transductive support vector machines (TSVMs)



minimize  $\underline{w}$   $w \cdot w$

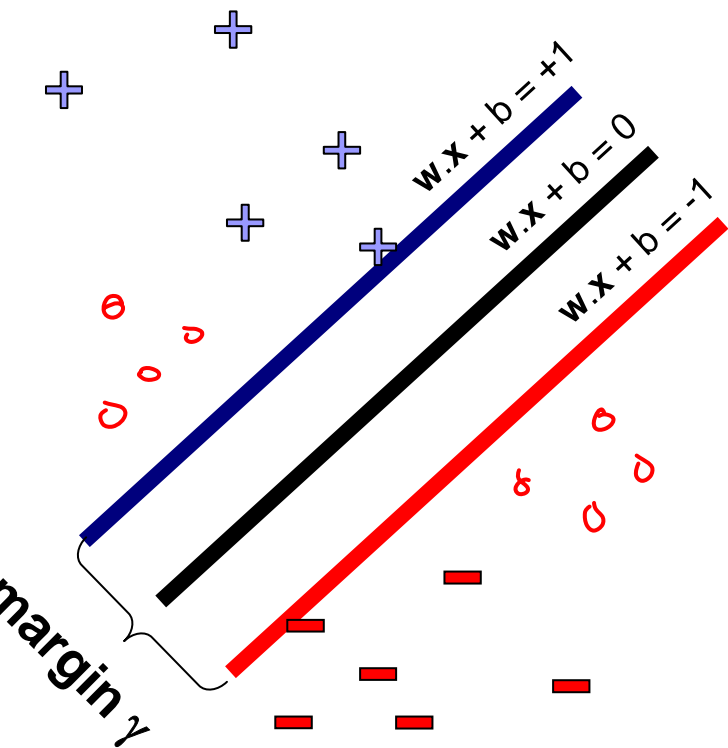
choose  $\{ \hat{y}_1, \dots, \hat{y}_n \}$

$$(w \cdot x_j + b) y_j \geq 1, \quad \forall j \in L$$

$$(w \cdot x_u + b) \hat{y}_u \geq 1 \quad \forall u \in U$$

$$\hat{y}_u \in \{+1, -1\} \quad \forall u \in U$$

# Transductive support vector machines (TSVMs)



$$\text{minimize}_{\mathbf{w}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}} \quad \mathbf{w} \cdot \mathbf{w}$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1, \quad \forall u = 1, \dots, n_U$$

*bilinear constraint*

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_U$$

*integer constraint - hard!*

# What's the difference between transductive learning and semi-supervised learning?

- Not much, and
- A lot!!!

## ■ Semi-supervised learning:

- labeled and unlabeled data  $\rightarrow$  learn  $\mathbf{w}$
- use  $\mathbf{w}$  on test data

## ■ Transductive learning

- same algorithms for labeled and unlabeled data, but...
- unlabeled data is test data!!!

## ■ You are learning on the test data!!!

- OK, because you never look at the labels of the test data
- can get better classification
- but be very very very very very very very very careful!!!

- never use test data prediction accuracy to tune parameters, select kernels, etc.





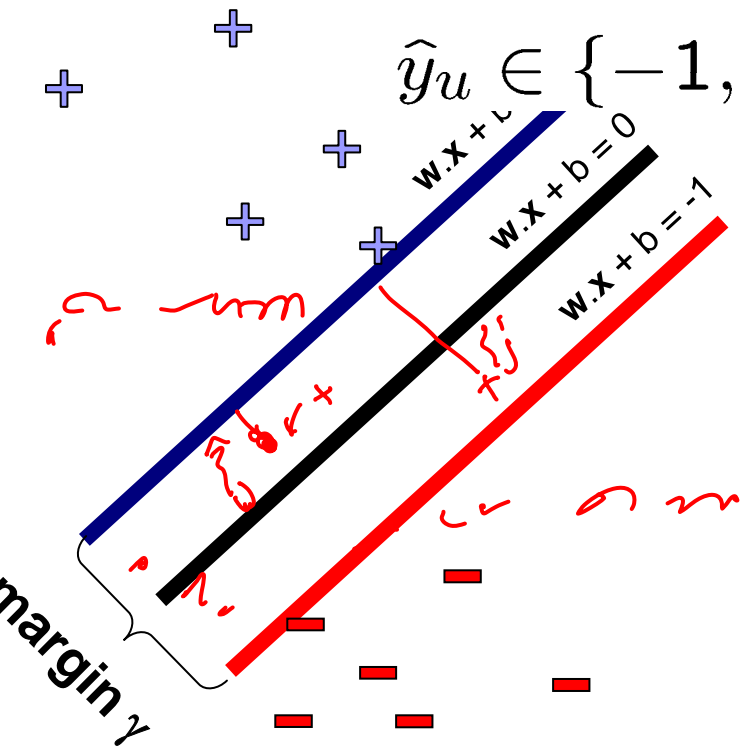
# Adding slack variables

$$\text{minimize}_{\mathbf{w}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}} \quad \mathbf{w} \cdot \mathbf{w} + C \sum \xi_j + \hat{C} \sum \hat{\xi}_u$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u \quad \forall u = 1, \dots, n_U$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_U$$



if  $\hat{C} = 0$   
learn only on labeled data  
(ignore unlabeled data)

# Transductive SVMs – now with slack variables!

[Vapnik 98]

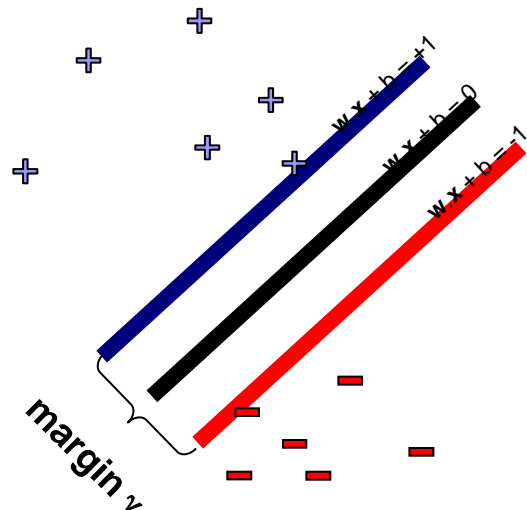
Optimize  $\underset{\text{weights}}{w}, \{\underset{\text{labeled slacks}}{\xi_1, \dots, \xi_{n_L}}\}, \{\underset{\text{class of unlabeled}}{\hat{y}_1, \dots, \hat{y}_{n_U}}\}, \{\underset{\text{unlabeled slacks}}{\hat{\xi}_1, \dots, \hat{\xi}_{n_U}}\}$

$$\text{minimize } w \cdot w + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$$

$$(w \cdot x_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$$

$$(w \cdot x_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$$



# Learning Transductive SVMs is hard!

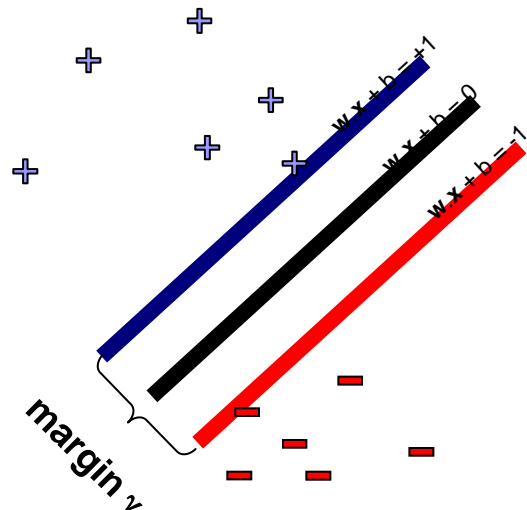
Optimize  $w, \{\xi_1, \dots, \xi_{n_L}\}, \{\hat{y}_1, \dots, \hat{y}_{n_U}\}, \{\hat{\xi}_1, \dots, \hat{\xi}_{n_U}\}$

minimize  $w \cdot w + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$

$(w \cdot x_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$

$(w \cdot x_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$

$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$



## ■ Integer Program

□ NP-hard!!!

□ Well-studied solution algorithms, but will not scale up to very large problems

*100 000 datapoints*

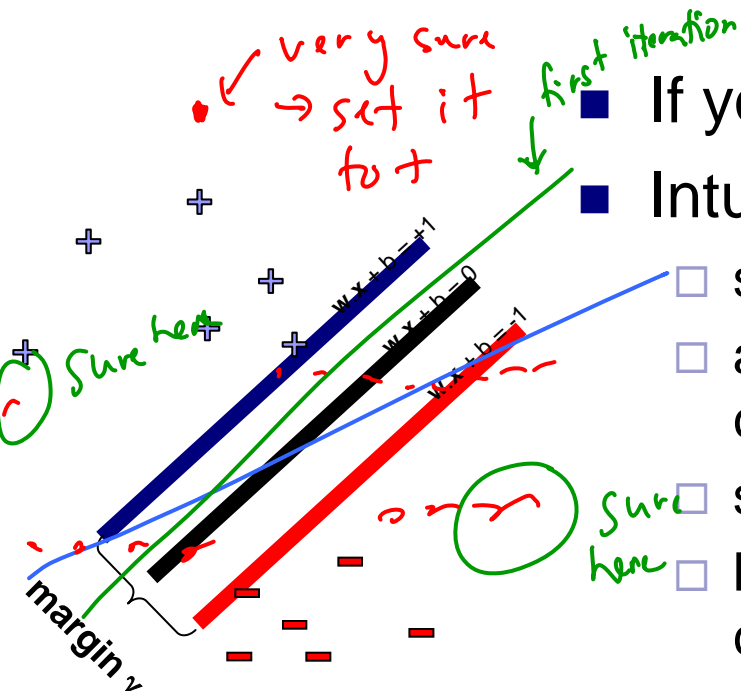
# A (~~heuristic~~) learning algorithm for Transductive SVMs [Joachims 99]

$$\text{minimize} \quad \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j + \hat{C} \sum_u \hat{\xi}_u$$

$$(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1 - \xi_j, \quad \forall j = 1, \dots, n_L$$

$$(\mathbf{w} \cdot \mathbf{x}_u + b) \hat{y}_u \geq 1 - \hat{\xi}_u, \quad \forall u = 1, \dots, n_u$$

$$\hat{y}_u \in \{-1, +1\}, \quad \forall u = 1, \dots, n_u$$



■ If you set  $\hat{C}$  to zero  $\rightarrow$  ignore unlabeled data

■ Intuition of algorithm:

- start with small  $\hat{C}$
- add labels to some unlabeled data based on classifier prediction *if classifier very sure about  $x_u \rightarrow$  label it*
- slowly increase  $\hat{C}$
- keep on labeling unlabeled data and re-running classifier

# Some results classifying news articles – from [Joachims 99]

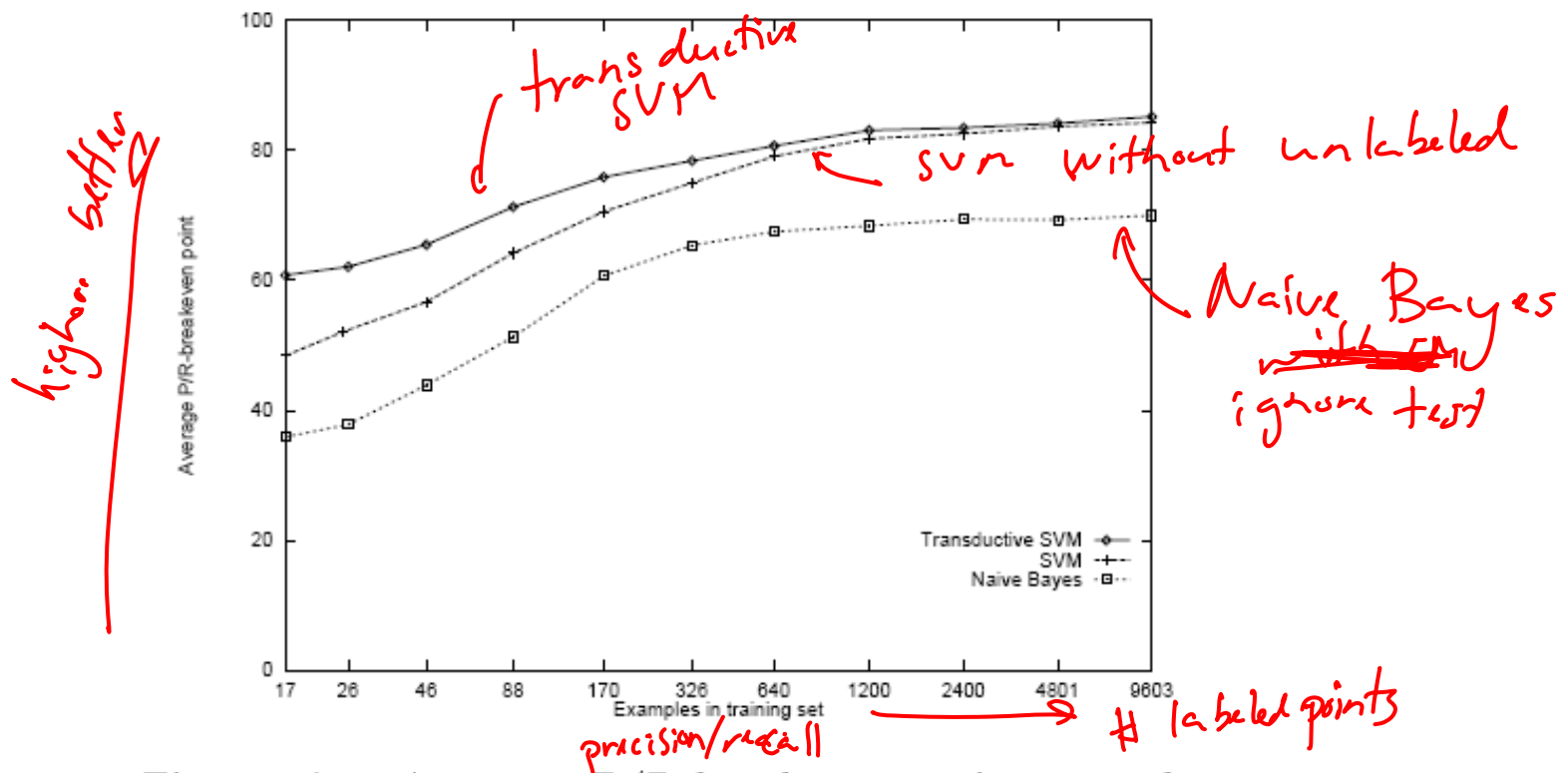


Figure 6: Average P/R-breakeven point on the Reuters dataset for different training set sizes and a test set size of 3,299. unlabeled

# What you need to know about transductive SVMs

- What is transductive v. semi-supervised learning
- Formulation for transductive SVM
  - can also be used for semi-supervised learning
- Optimization is hard!
  - Integer program
- There are simple heuristic solution methods that work well here

Recommended reading:

Bishop, Chapters 3.6, 8.6

Shlens PCA tutorial

Wall et al. 2003 (PCA applied to gene expression data)



# Dimensionality reduction

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 24<sup>th</sup>, 2006

# Dimensionality reduction

- Input data may have thousands or millions of dimensions!

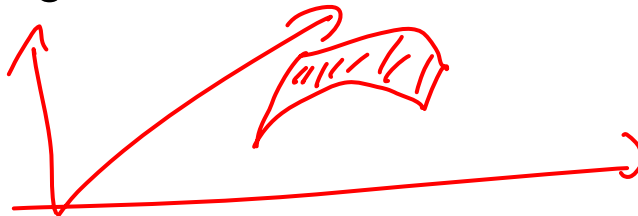
- e.g., text data has

about 10,000 — 100,000  
single words

pairs of words → 100,000<sup>2</sup>  
features

- **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





# Feature selection

- Want to learn  $f: \mathbf{X} \mapsto Y$  *for classification*
  - $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
  - 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

# Simple greedy forward feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression

$\{1, x, x^2, \dots\}$

- Greedy heuristic:

$\{1, x\}$

- Start from empty (or simple) set of features  $F_0 = \emptyset$
- Run learning algorithm for current set of features  $F_t$ 
  - Obtain  $h_t$

hypothesis

- Select **next best feature**  $X_i$ 
  - e.g.,  $X_i$  that results in lowest cross-validation error learner when learning with  $F_t \cup \{X_i\}$
- $F_{t+1} \leftarrow F_t \cup \{X_i\}$

add in best feature

- Recurse

there other criteria  
- mutual info.  
- X-squared test  
- ...  
- ...

# Simple greedy **backward** feature selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- Greedy heuristic:
  - Start from all features  $F_0 = F$
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select next worst feature  $X_i$ 
    - e.g.,  $X_i$  that results in lowest cross-validation error learner when learning with  $F_t - \{X_i\}$  *remove*
  - $F_{t+1} \leftarrow F_t - \{X_i\}$  *remove it.*
  - Recurse

# 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	496B	77B	86B
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
300	<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

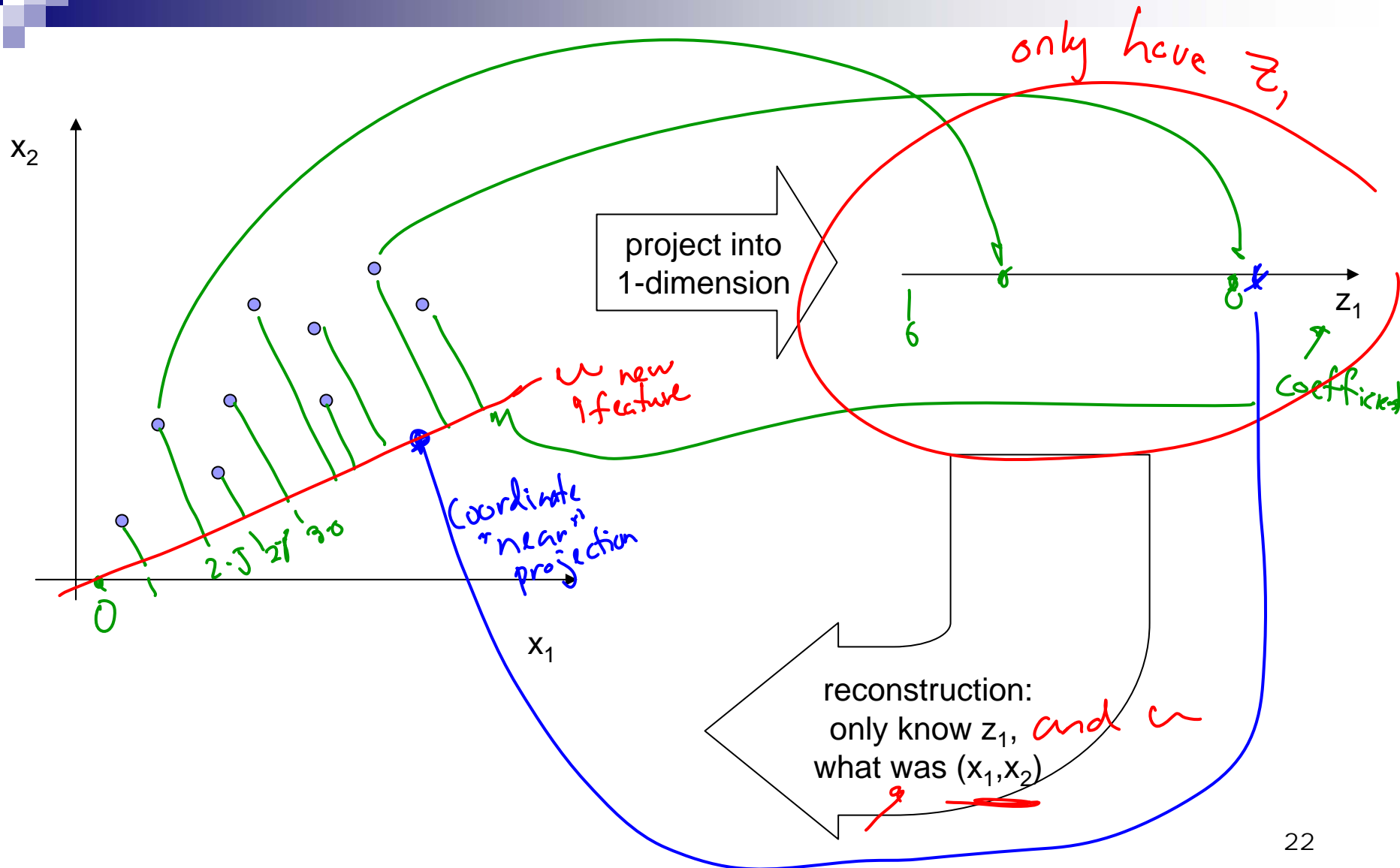
# Lower dimensional projections

- Rather than picking a subset of the features, we can <sup>generate</sup> new features that are combinations of existing features

$$X_{\text{new}} = 10X^2 + 17X^{17} - 28X^{-12}$$
$$X_{\text{new}} = g(X)$$

- Let's see this in the unsupervised setting
  - just X, but no Y  
*ignore / don't have*

# 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

## ■ Project a point into a (lower dimensional) space:

□ point:  $\mathbf{x} = (x_1, \dots, x_n)$

□ **select a basis** – set of basis vectors –  $(\mathbf{u}_1, \dots, \mathbf{u}_k)$

■ we consider orthonormal basis:  $u_i = \begin{matrix} \uparrow \\ \beta_i \end{matrix}$ ,  $u_i = \begin{matrix} \uparrow \\ \beta_i \end{matrix} \rightarrow z_{ij}$

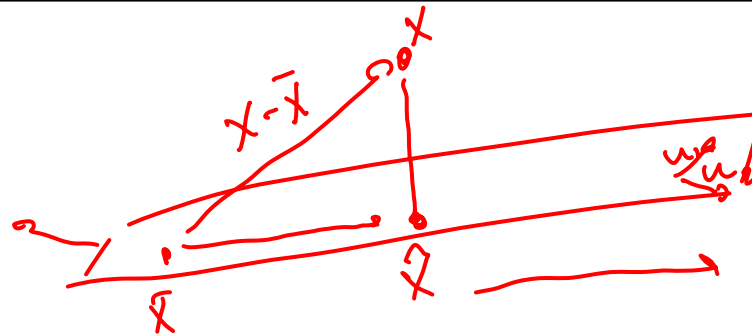
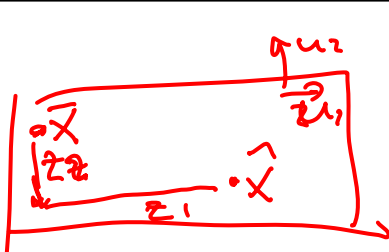
□  $\mathbf{u}_i \cdot \mathbf{u}_i = 1$ , and  $\mathbf{u}_i \cdot \mathbf{u}_j = 0$  for  $i \neq j$

□ **select a center** –  $\bar{\mathbf{x}}$ , defines offset of space

□ **best coordinates** in lower dimensional space defined by dot-products:  $(z_1, \dots, z_k)$ ,  $z_i = (\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{u}_i$

■ minimum squared error

$$a \cdot b = \sum_{i=1}^n a_i b_i$$



$$\hat{\mathbf{x}} = \bar{\mathbf{x}} + \sum_{i=1}^k z_i \mathbf{u}_i$$



# PCA finds projection that minimizes reconstruction error

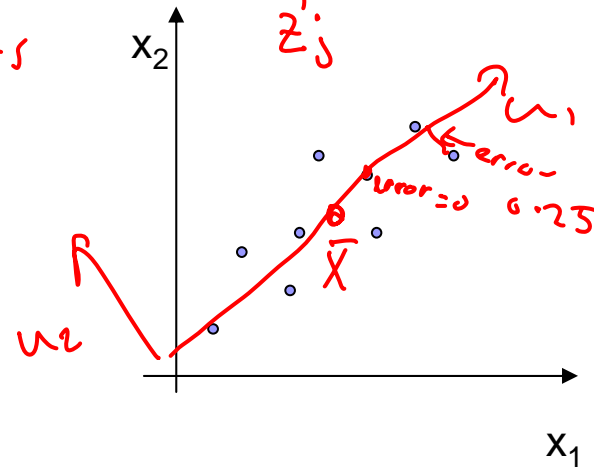
- Given m data points:  $\mathbf{x}^i = (x_1^i, \dots, x_n^i)$ ,  $i=1 \dots m$
- Will represent each point as a projection:

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

## ■ PCA:

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

$$\text{error}_k = \sum_{i=1}^m (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$



# Understanding the reconstruction error

only use first  $K$

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

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

define  $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$

exactly

$$\text{error}_K = \sum_{i=1}^m (\mathbf{x}_i - \hat{\mathbf{x}}_i)^2$$

- Rewriting error:

$$= \sum_{i=1}^m \left[ \bar{\mathbf{x}} + \sum_{j=1}^n z_j^i \mathbf{u}_j - \left( \bar{\mathbf{x}} + \sum_{j=1}^K z_j^i \mathbf{u}_j \right) \right]^2 = \sum_{i=1}^m \left[ \sum_{j=K+1}^n z_j^i \mathbf{u}_j \right]^2$$

$$= \sum_{i=1}^m \sum_{j=K+1}^n (z_j^i \mu_j)^2 + \sum_{i=1}^m \sum_{j=K+1}^n \sum_{v=K+1, v \neq j}^n \cancel{z_j^i \mu_j \cdot z_v^i \mu_v} \left[ z_j^i \mu_j \cdot z_v^i \mu_v \right]$$

$\mathbf{u}_j \cdot \mathbf{u}_v = 0$  (orthogonal)

$$= \sum_{i=1}^m \sum_{j=K+1}^n (\mu_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}}) \mu_j)^2$$

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

$$\rightarrow \hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^K z_j^i \mathbf{u}_j \quad z_j^i = \mathbf{x}^i \cdot \mathbf{u}_j$$

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

$$\text{error}_K = \sum_{i=1}^m (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$

$$(a+b)^2 = a^2 + b^2 + 2ab$$

# Reconstruction error and covariance matrix

$$(a \cdot b)^2 = a^T b \quad b^T a$$

$$a^2 = a^T a$$

$u_i$  is  $n$  by 1

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

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

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

$m \cdot \Sigma$  ← covariance matrix of data

$$\begin{matrix} n \times 1 & 1 \times n \\ (x^i - \bar{x}) & (x^i - \bar{x})^T \end{matrix} = \begin{matrix} n \times 1 \\ \vdots \\ (x_j^i - \bar{x}_j) \\ \vdots \\ (x^i_e - \bar{x}_e) \end{matrix} \left[ \begin{matrix} \sigma_{je} \end{matrix} \right]$$

$$\begin{aligned} \text{error}_k &= \sum_{i=1}^m \sum_{j=k+1}^n [u_j \cdot (x^i - \bar{x})] \underbrace{u_j}_{z_j^i} \\ &\quad \begin{matrix} \text{scalar} \\ \downarrow \\ (z_j^i u_j)^2 \end{matrix} \quad \begin{matrix} n \text{ vector} \\ \downarrow \\ u_j \end{matrix} \\ &= z_j^i u_j^T u_j z_j^i \quad \left| \begin{matrix} u_j^T u_j \\ = 1 \\ \text{orthonormal} \end{matrix} \right. \\ &= z_j^i \cdot z_j^i \\ &= z_j^i = u_j \cdot (x^i - \bar{x}) \end{aligned}$$

# Minimizing reconstruction error and eigen vectors

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

$$\text{error}_k = \sum_{j=k+1}^n \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

← consider ignored vectors in error

- Eigen vector:

$$\mathbf{u}: \quad \underbrace{\Sigma \mathbf{u}}_{\substack{\uparrow \\ \text{eigen} \\ \text{vector}}} = \lambda \underbrace{\mathbf{u}}_{\substack{\uparrow \\ \text{eigen} \\ \text{value}}} \quad \Rightarrow \quad \mathbf{u}^T \Sigma \mathbf{u} = \lambda \mathbf{u}^T \mathbf{u} = \lambda$$

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

take eigen vectors      values of  $\Sigma$

sort!	$\mathbf{u}_1$	$\lambda_1$
	$\mathbf{u}_2$	$\lambda_2$
$\lambda_i \geq \lambda_{i+1}$	$\vdots$	$\vdots$

minimize error<sub>k</sub>

$$\Rightarrow \min_{\mathbf{u}} \sum_{j=k+1}^n \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

$$= \min_{\mathbf{u}} \sum_{j=k+1}^n \lambda_j$$

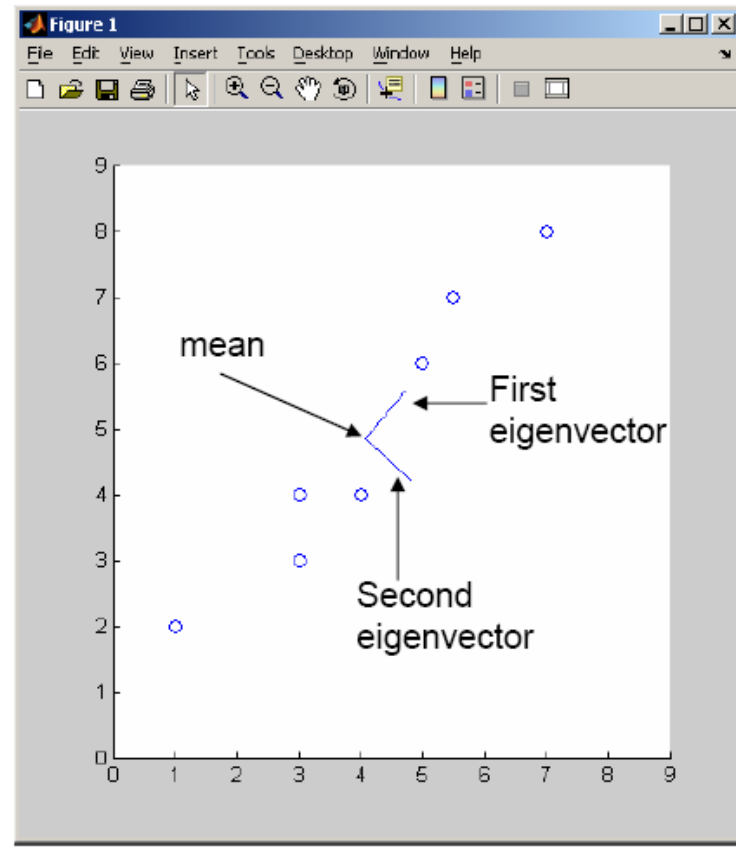
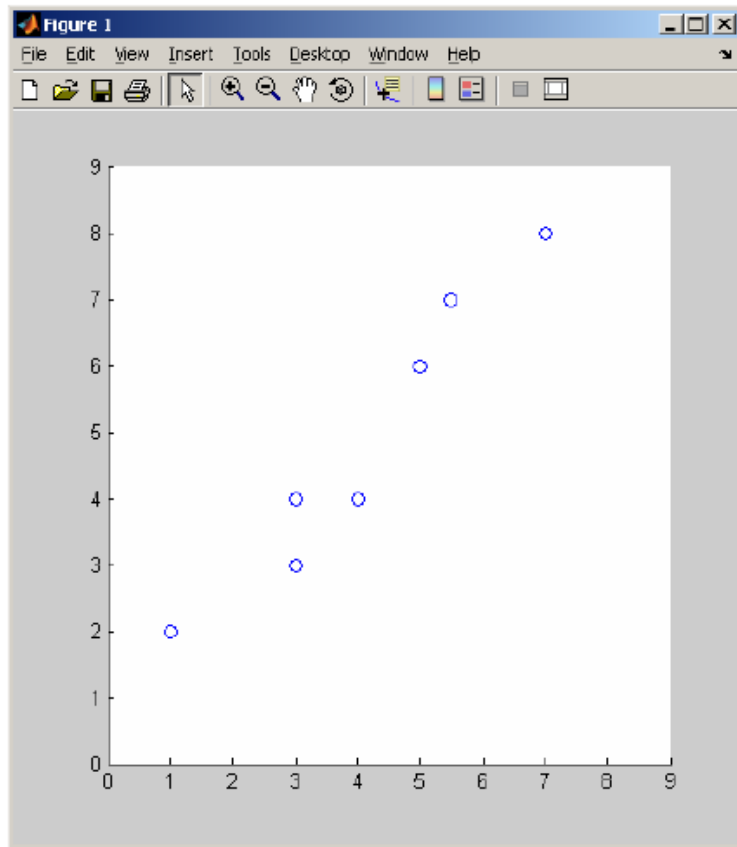
$\Rightarrow$  ignore vectors with low eigen val

# 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**:
  - $\Sigma \leftarrow \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of  $\Sigma$
- **Principal components**: k eigen vectors with highest eigen values

# PCA example

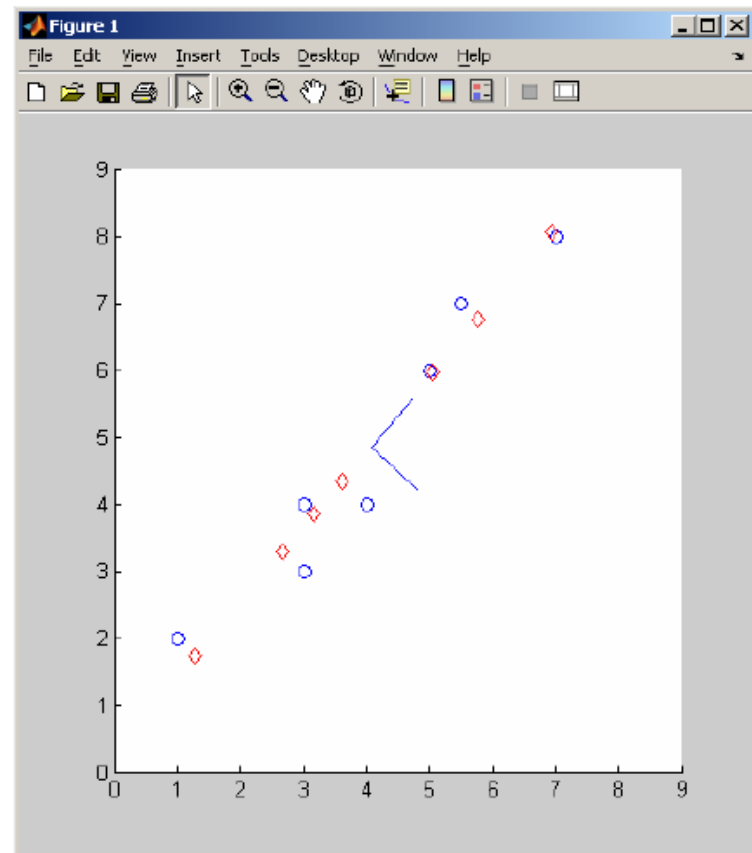
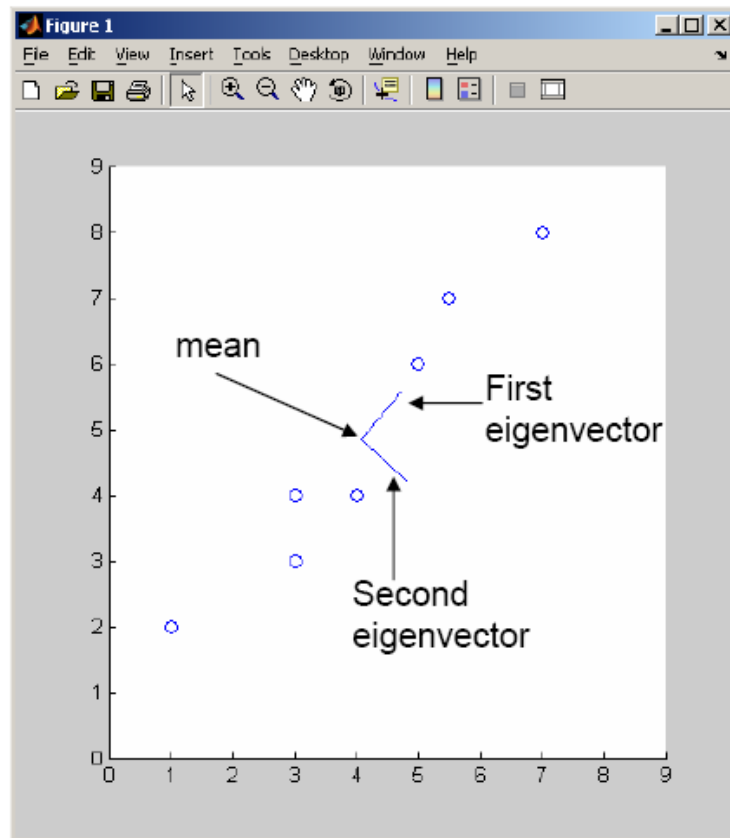
$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$



# PCA example – reconstruction

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

only used first principal component



# Eigenfaces [Turk, Pentland '91]

## ■ Input images:



## ■ Principal components:





# Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



# Relationship to Gaussians

- PCA assumes data is Gaussian

- $\mathbf{x} \sim N(\bar{\mathbf{x}}; \Sigma)$

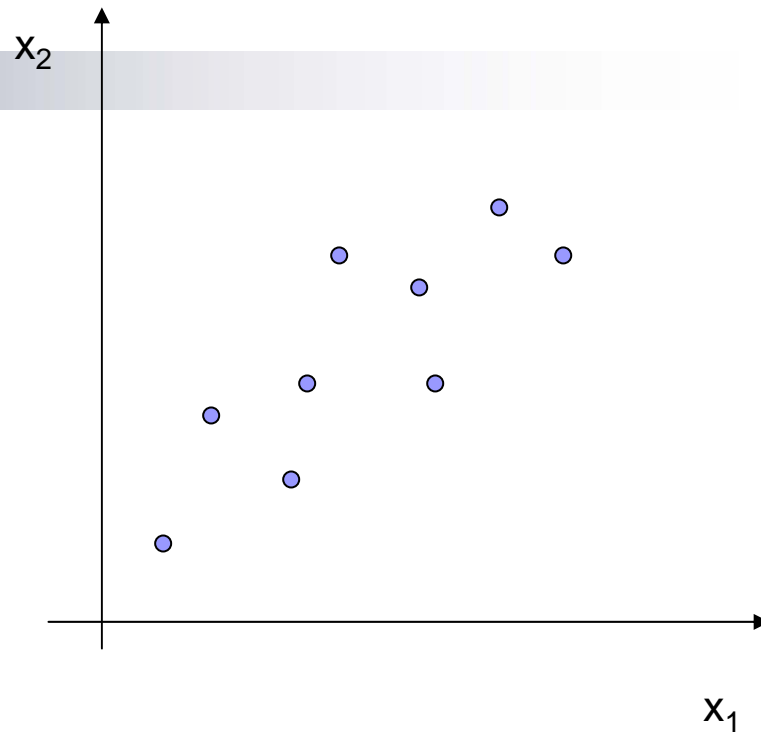
- Equivalent to weighted sum of simple Gaussians:

$$\mathbf{x} = \bar{\mathbf{x}} + \sum_{j=1}^n z_j \mathbf{u}_j; \quad z_j \sim N(0; \sigma_j^2)$$

- Selecting top k principal components equivalent to lower dimensional Gaussian approximation:

$$\mathbf{x} \approx \bar{\mathbf{x}} + \sum_{j=1}^k z_j \mathbf{u}_j + \varepsilon; \quad z_j \sim N(0; \sigma_j^2)$$

- $\varepsilon \sim N(0; \sigma^2)$ , where  $\sigma^2$  is defined by error<sub>k</sub>



# Scaling up

- Covariance matrix can be really big!
  - $\Sigma$  is  $n$  by  $n$
  - 10000 features  $\rightarrow |\Sigma|$
  - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
  - finds to  $k$  eigenvectors
  - great implementations available, e.g., Matlab `svd`

# SVD

- Write  $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$ 
  - $\mathbf{X} \leftarrow$  data matrix, one row per datapoint
  - $\mathbf{U} \leftarrow$  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  $\mathbf{v}_j$

# 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:

# Using PCA for dimensionality reduction in classification

- Want to learn  $f: \mathbf{X} \mapsto Y$ 
  - $\mathbf{X} = \langle X_1, \dots, X_n \rangle$
  - but some features are more important than others
- **Approach:** Use PCA on  $\mathbf{X}$  to select a few important features

# PCA for classification can lead to problems...

- Direction of maximum variation may be unrelated to “discriminative” directions:
- PCA often works very well, but sometimes must use more advanced methods
  - e.g., Fisher linear discriminant

# 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
  - problems with PCA