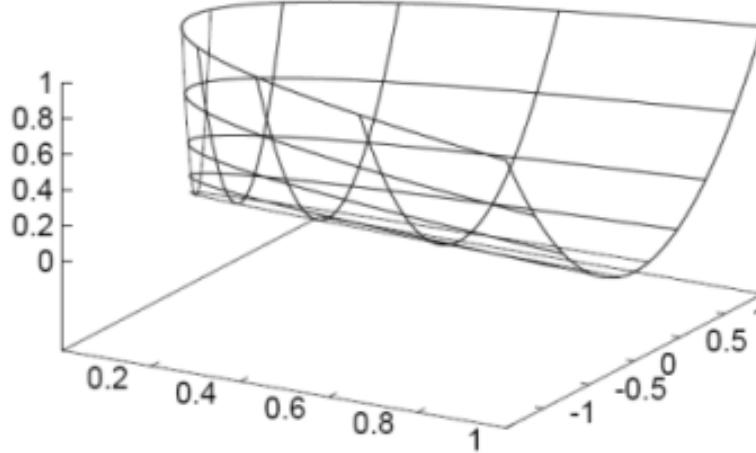


# Introduction to Supervised Learning



Week 3:  
Support Vector Machines

Iasonas Kokkinos

i.kokkinos@cs.ucl.ac.uk

University College London



# Lecture outline

Introduction to Support Vector Machines

Geometric margins

Training criterion & hinge loss

Large margins and generalization

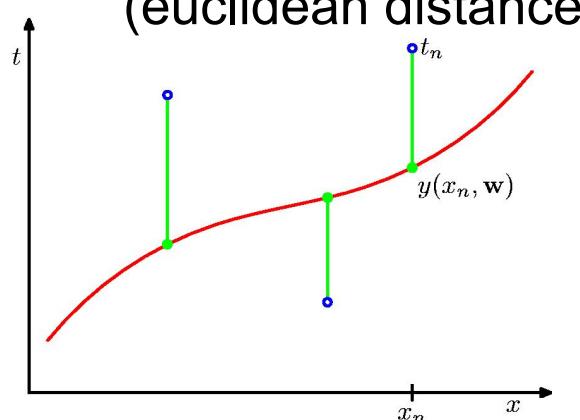
Optimization

Kernels

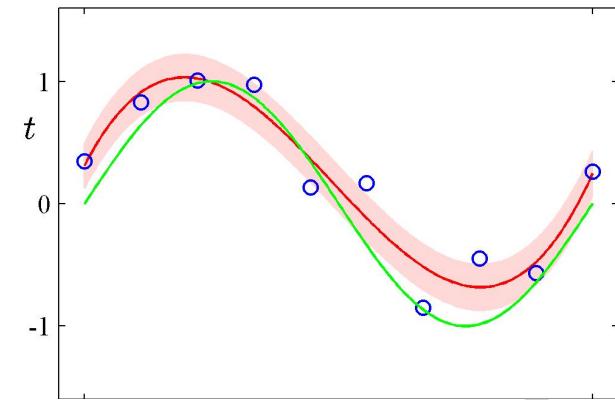
Applications to vision

# Our path so far (week 1-2)

Week 1 - regression: geometric  
(euclidean distance)



Week 2: probabilistic interpretation

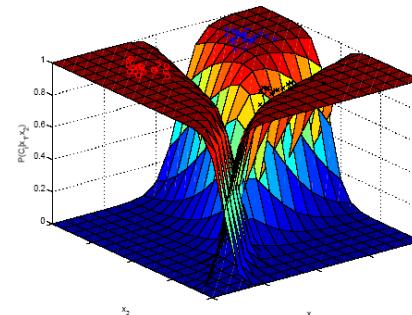
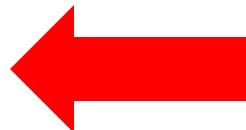


$$P(y^i | \mathbf{x}^i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - \mathbf{w}^T \mathbf{x}^i)^2}{2\sigma^2}\right)$$

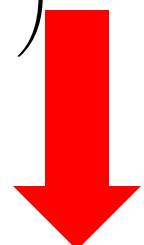
Week 2: switch to classification



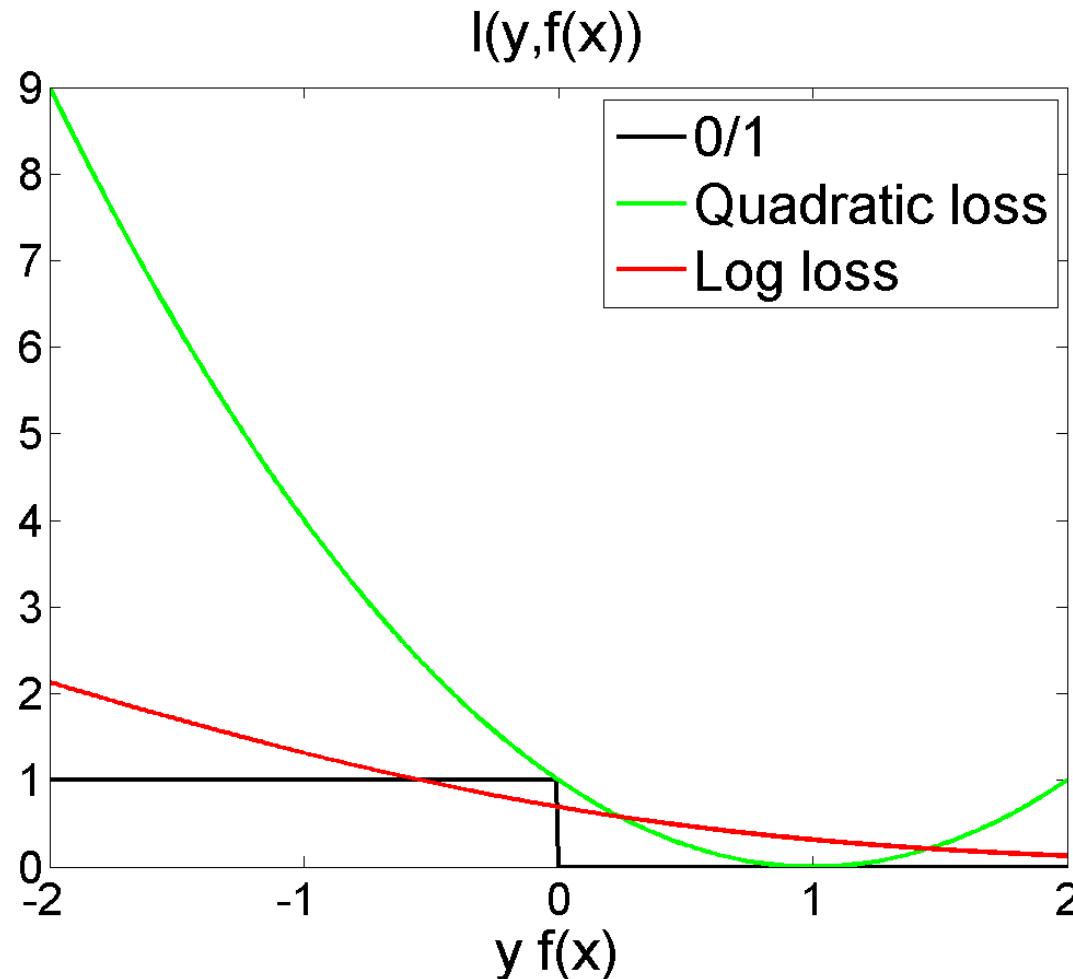
geometry + classification?



$$P(y^i | \mathbf{x}^i) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^T \mathbf{x})}$$



# Week 2: log loss vs. quadratic loss



Quadratic loss

$$l(y, f(x)) = (1 - yf(x))^2$$

Log loss

$$l(y, f(x)) = \log(1 + \exp(-yf(x)))$$

# Do we need the logistic loss?

Week 2: Useful criterion for training classifiers

Maybe we can quickly hack an easy algorithm

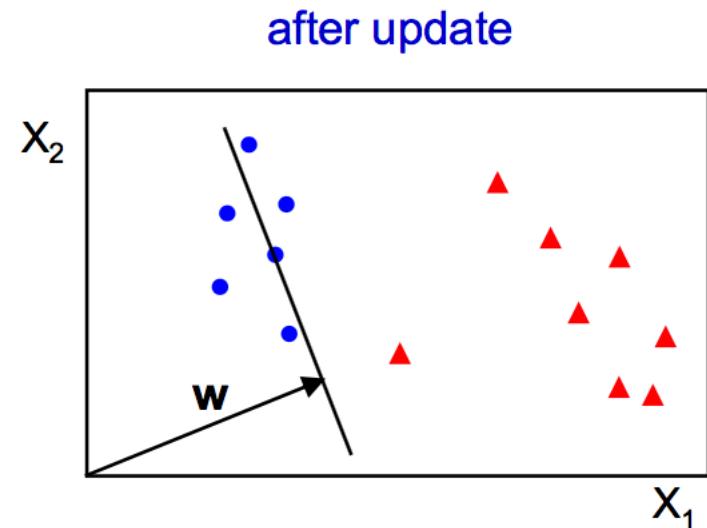
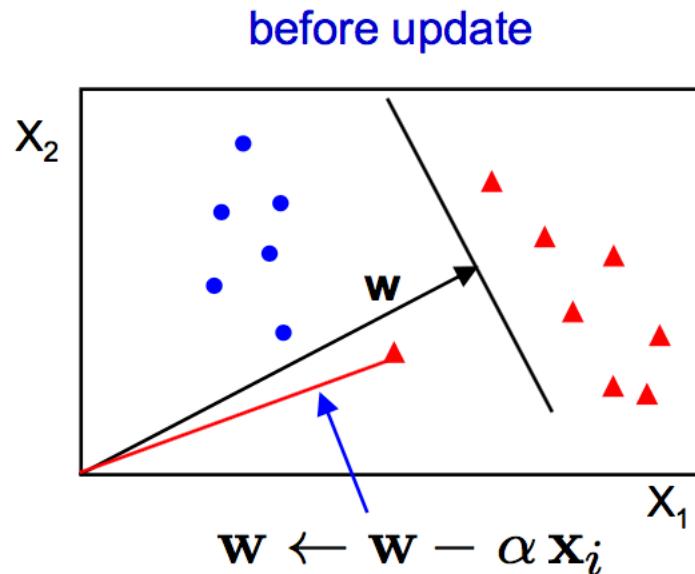
Least squares: Gauss, 1795

Logistic Regression: Cox, 1958

Perceptrons, Minsky & Papert, 1969

# Perceptron algorithm

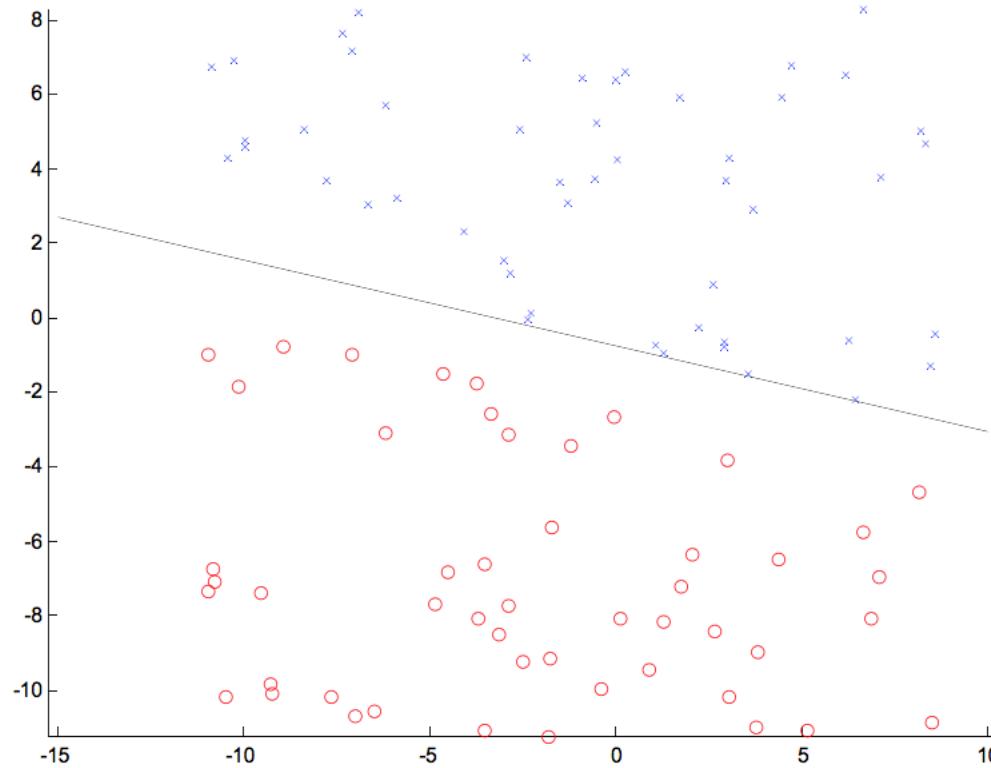
- Initialize  $\mathbf{w} = 0$
- Cycle though the data points  $\{ \mathbf{x}_i, y_i \}$ 
  - if  $\mathbf{x}_i$  is misclassified then  $\mathbf{w} \leftarrow \mathbf{w} + \alpha \text{sign}(f(\mathbf{x}_i)) \mathbf{x}_i$
- Until all the data is correctly classified



after convergence  $\mathbf{w} = \sum_i^N \alpha_i \mathbf{x}_i$

# Perceptron algorithm (first ‘neural network’)

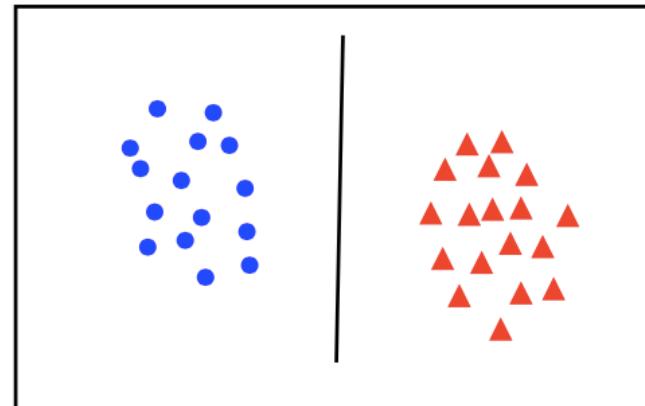
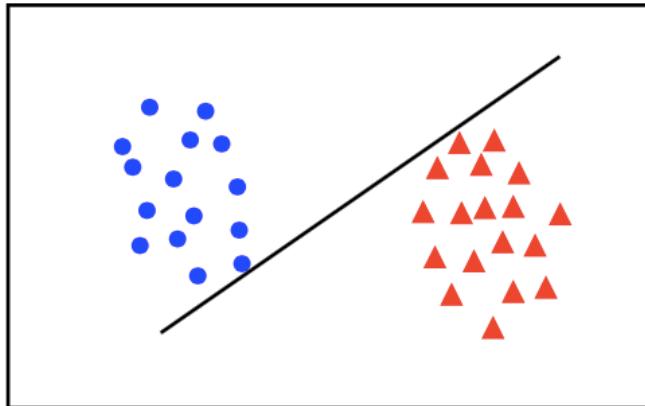
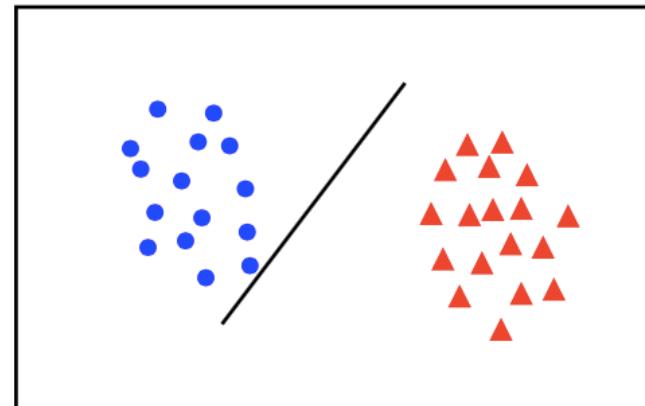
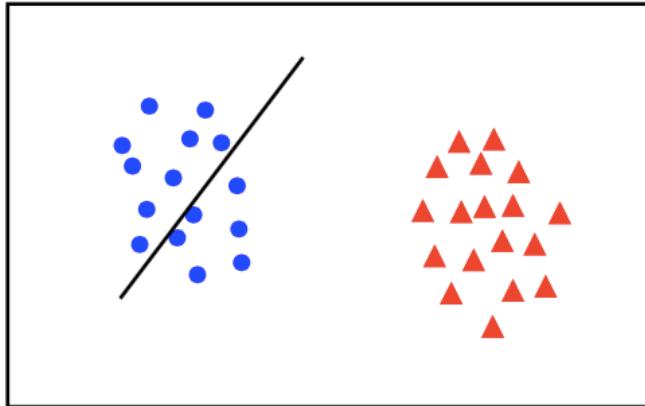
Perceptron example



- if the data is linearly separable, then the algorithm will converge
- convergence can be slow ...
- separating line close to training data

**This lecture: push it far away!**

# Which classifier is best?



All points should lie **clearly** on the correct side of the boundary

How can we quantify this?

How can we enforce this?

# Functional Margins

Consider Logistic Regression:

$$P(y = 1 | \mathbf{x}; \mathbf{w}) = g(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

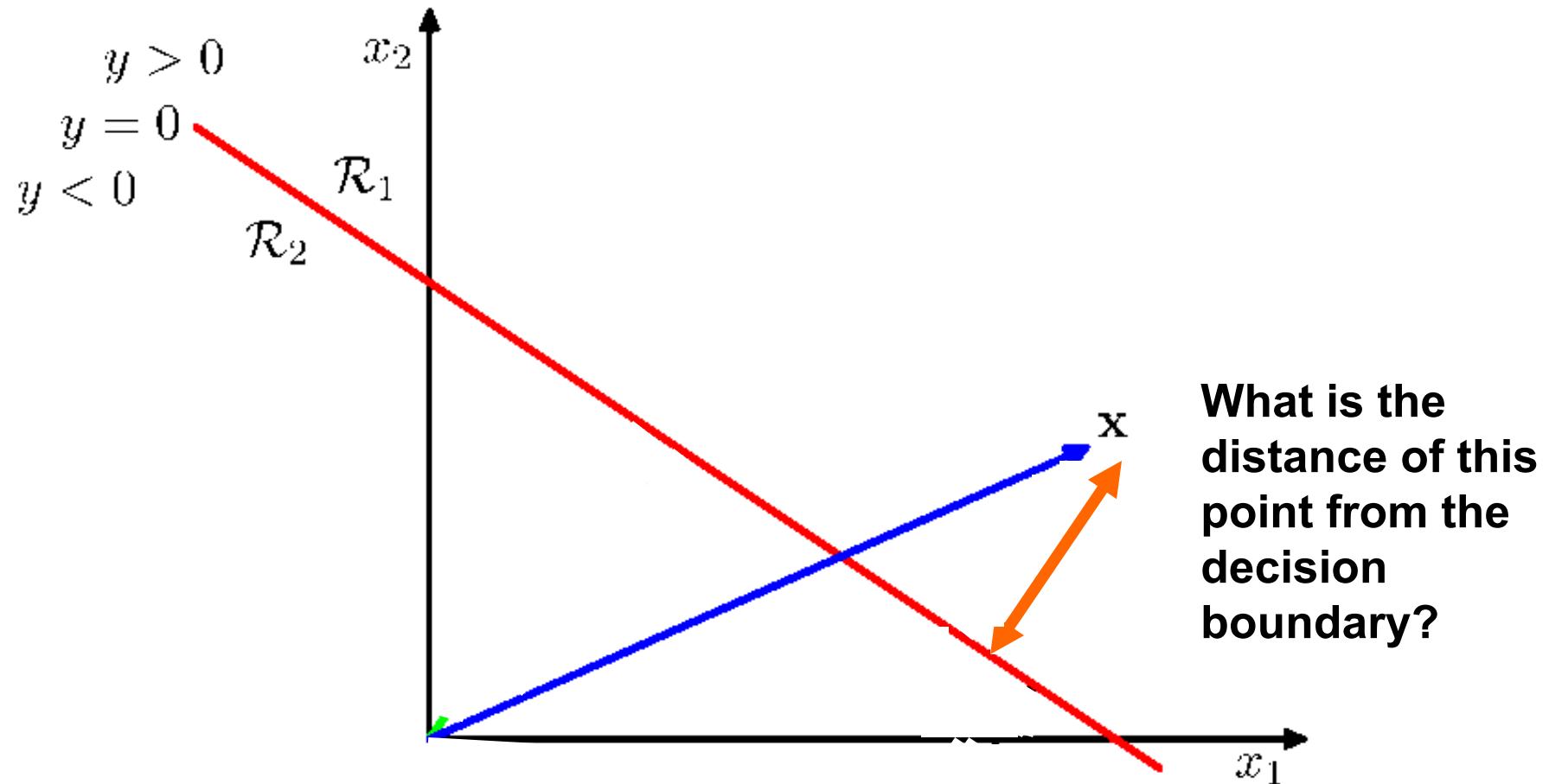
Ideally:

$\mathbf{w}^T \mathbf{x}^i \gg 0,$	if	$y^i = 1$
$\mathbf{w}^T \mathbf{x}^i \ll 0,$	if	$y^i = -1$

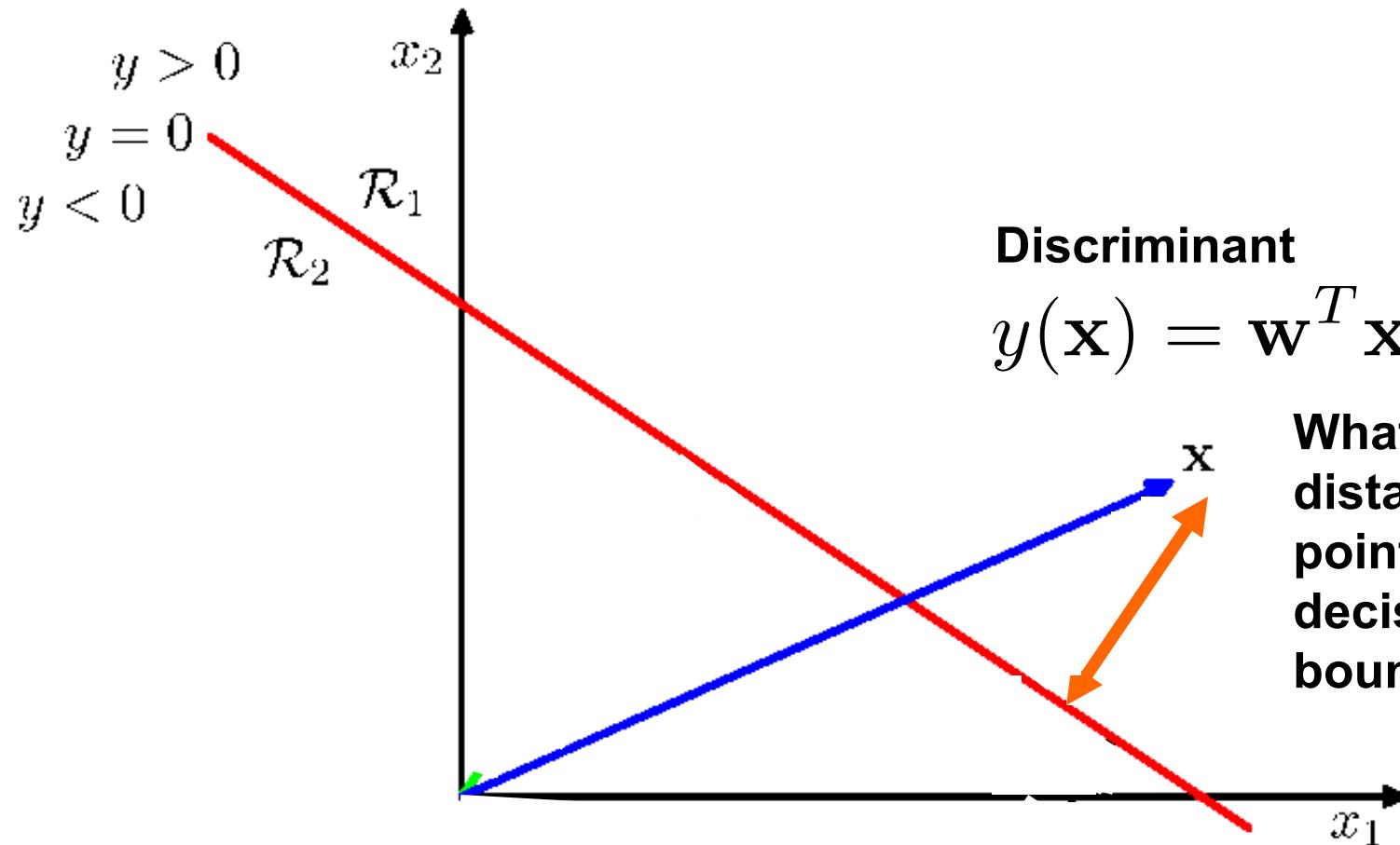
Put together:  $y^i (\mathbf{w}^T \mathbf{x}^i) \gg 0$   
 ‘functional margin’

Problem: scaling  $\mathbf{w}$  changes functional margin, but not decision boundary

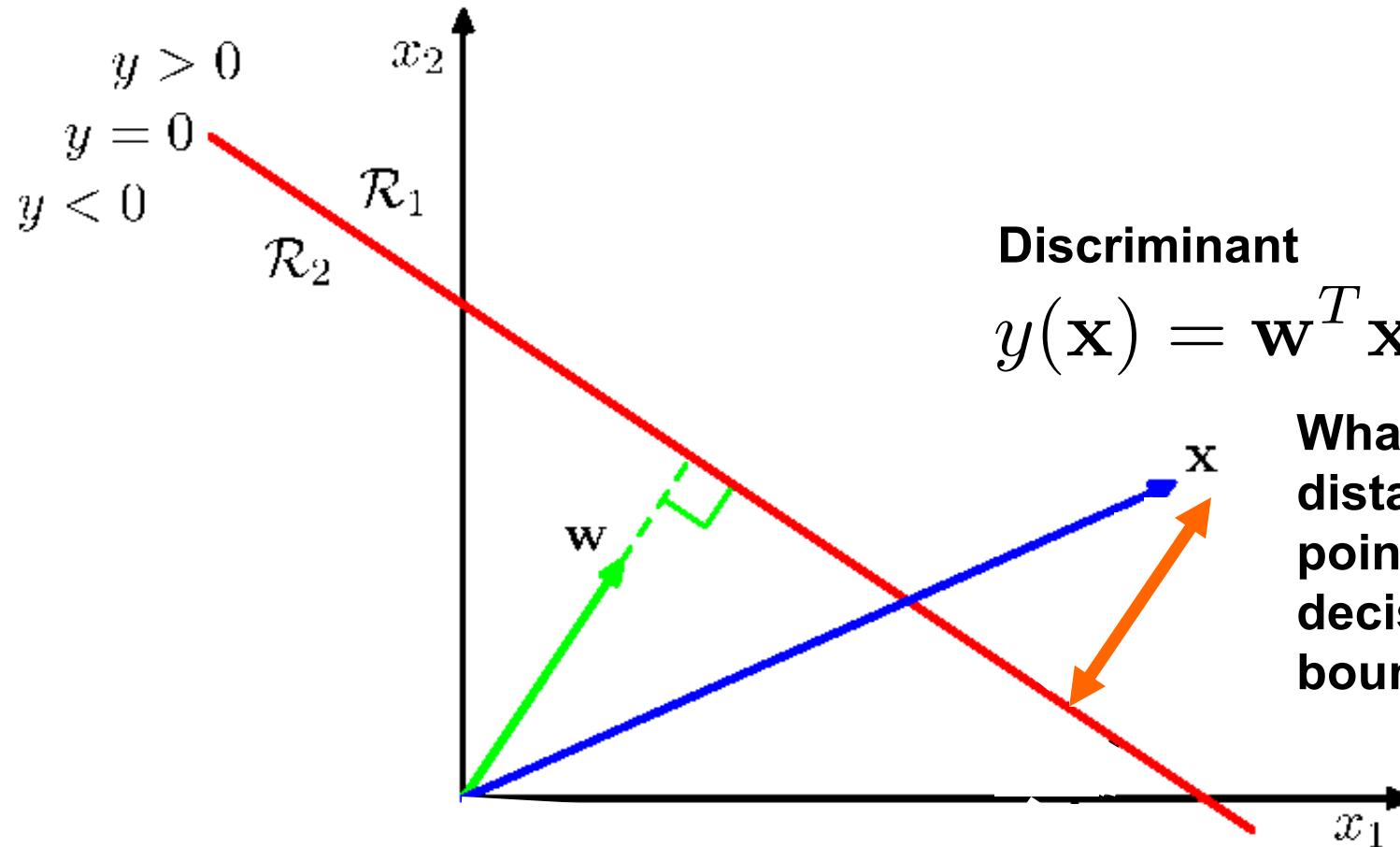
# Geometric Margins



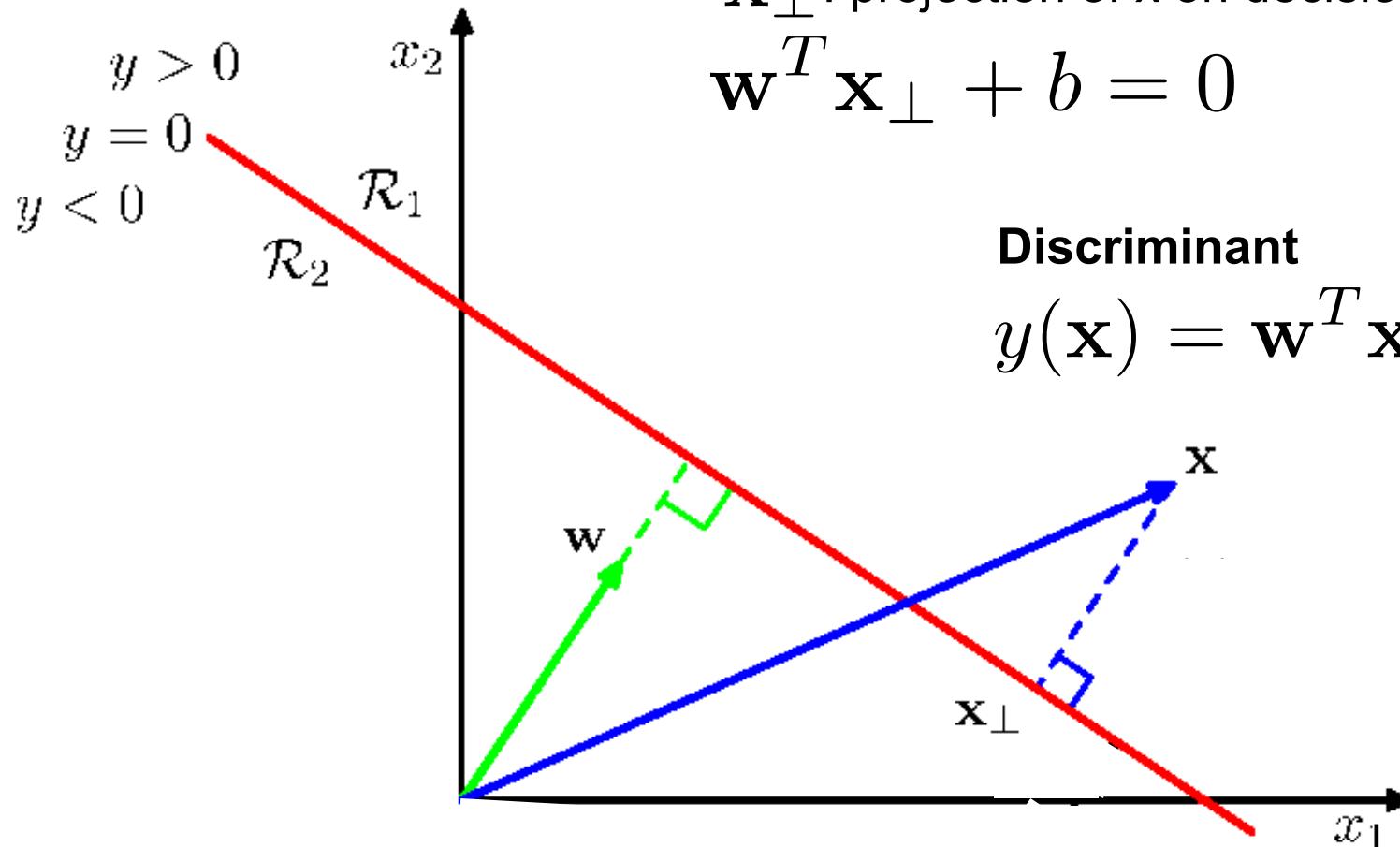
# Geometric Margins



# Geometric Margins



# Geometric Margins



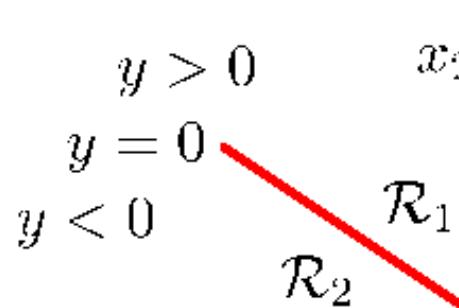
$\mathbf{x}_{\perp}$ : projection of  $\mathbf{x}$  on decision boundary

$$\mathbf{w}^T \mathbf{x}_{\perp} + b = 0$$

**Discriminant**

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

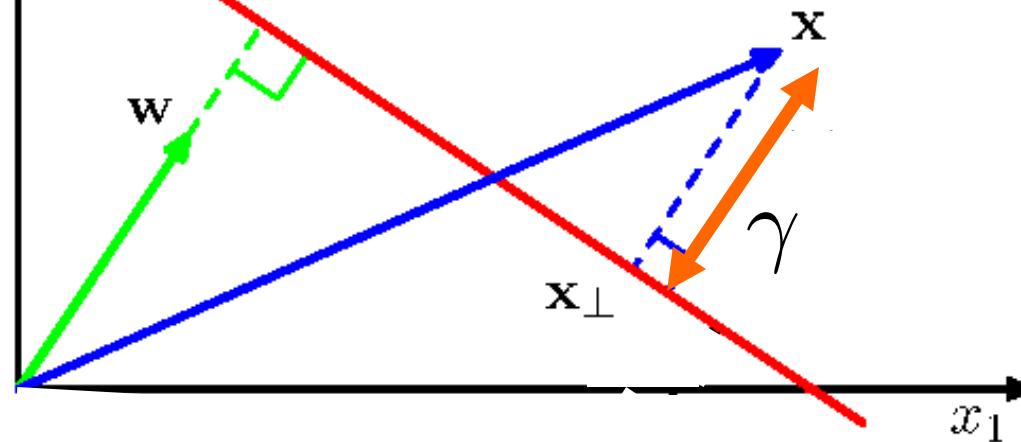
# Geometric Margins



$$\mathbf{x} = \mathbf{x}_{\perp} + \gamma \frac{\mathbf{w}}{|\mathbf{w}|}$$

**Discriminant**

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$



# Geometric Margins

Point = projection + distance\* direction

$$\mathbf{x} = \mathbf{x}_{\perp} + \gamma \frac{\mathbf{w}}{|\mathbf{w}|}$$

Note:  $\gamma$  is independent of  $|\mathbf{w}|$

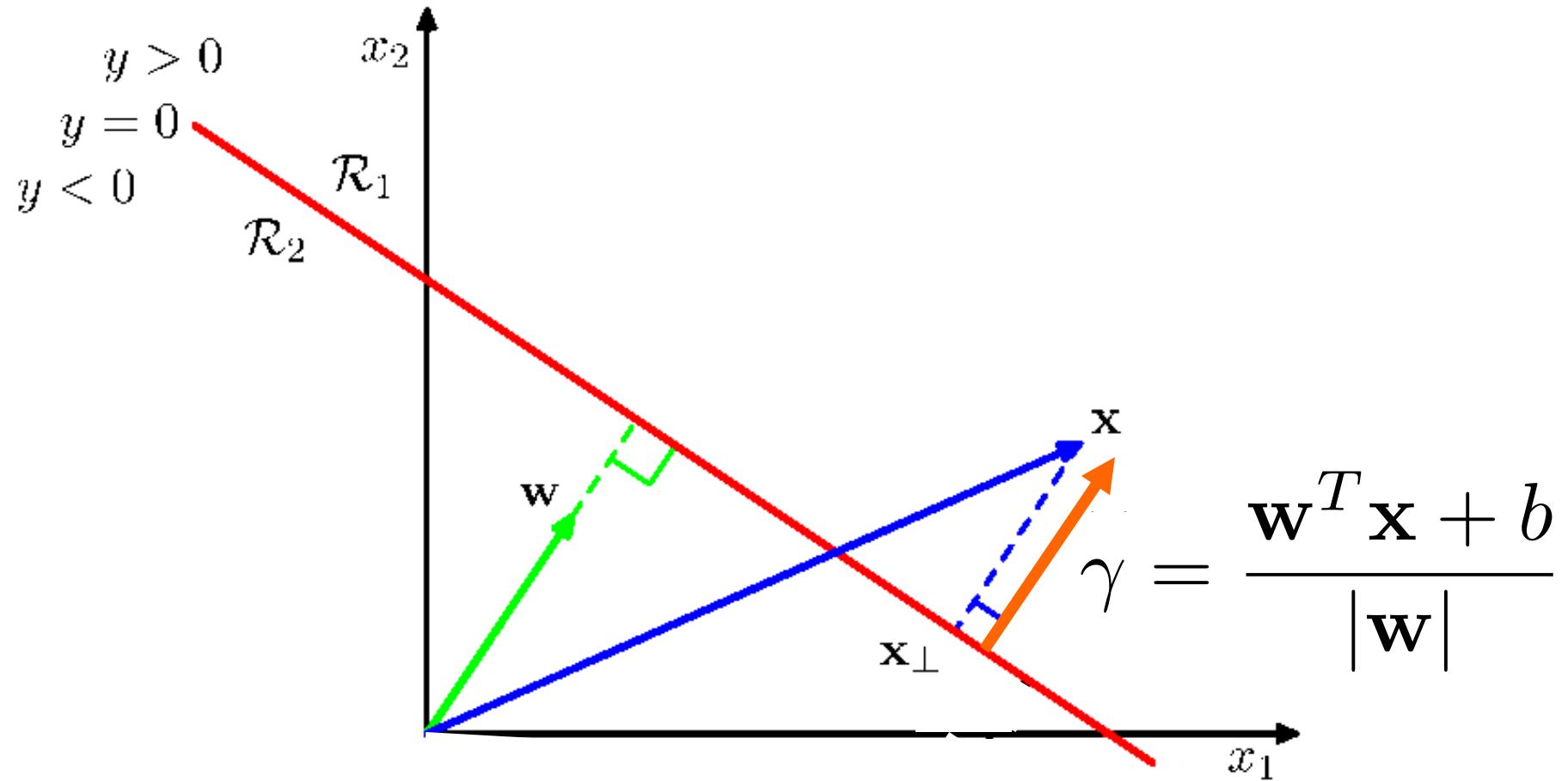
Multiply:  $\mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{x}_{\perp} + \mathbf{w}^T \gamma \frac{\mathbf{w}}{|\mathbf{w}|}$

Rewrite ( $\mathbf{w}^T \mathbf{x}_{\perp} + b = 0$ ):

$$\mathbf{w}^T \mathbf{x} = -b + \gamma |\mathbf{w}|$$

Solve for  $\gamma$ :  $\gamma = \frac{\mathbf{w}^T \mathbf{x} + b}{|\mathbf{w}|} = \frac{\mathbf{w}^T}{|\mathbf{w}|} \mathbf{x} + \frac{b}{|\mathbf{w}|}$

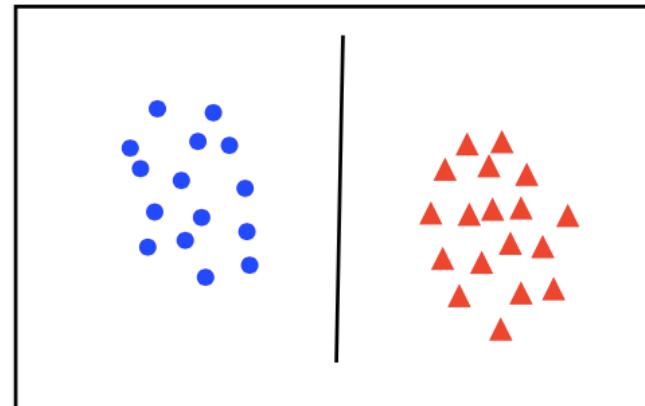
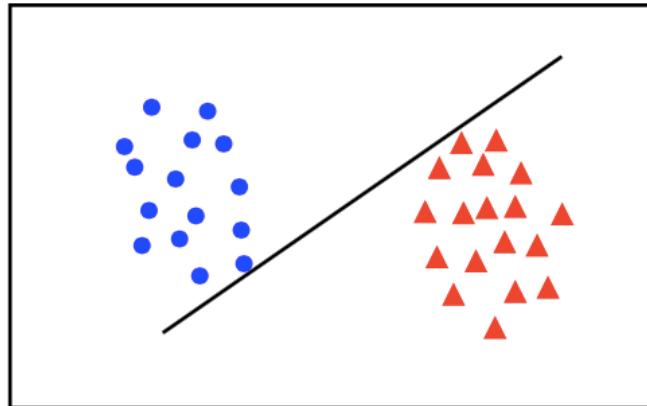
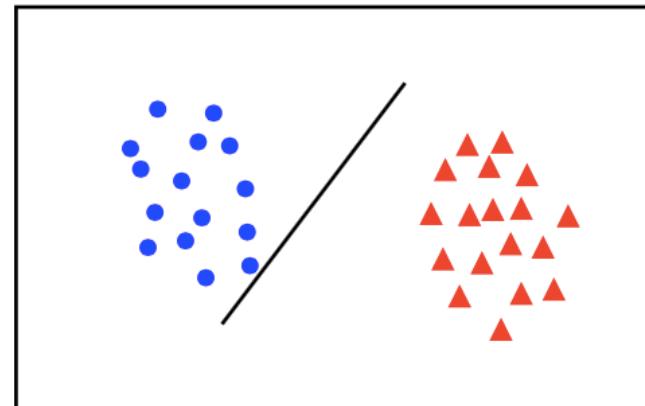
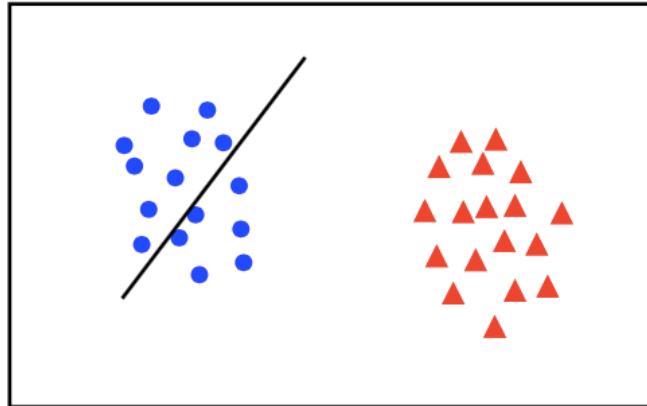
# Geometric Margins



**Geometric Margin:**  $\gamma^i = y^i \left( \frac{\mathbf{w}^T}{\|\mathbf{w}\|} \mathbf{x}^i + \frac{b}{\|\mathbf{w}\|} \right)$

(positive if  $\mathbf{x}$  is on the correct side of the decision boundary)

# Which classifier is best?



All points should lie **clearly** on the correct side of the boundary

How can we quantify this? (large margins!)

How can we enforce this?



# Lecture outline

Introduction to Support Vector Machines

Geometric margins

Training criterion & hinge loss

Large margins and generalization

Optimization

Kernels

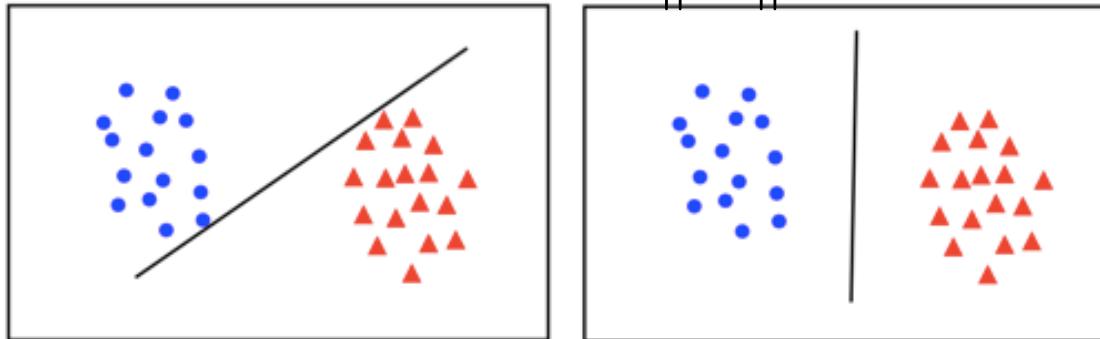
Applications to vision

# What should we be optimizing?

Training set:  $\{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^N, y^N)\}$

Candidate parameter vector:  $(\mathbf{w}, b)$

Related margins:  $\gamma^i = y^i \frac{\mathbf{w}^T \mathbf{x}^i + b}{\|\mathbf{w}\|}$

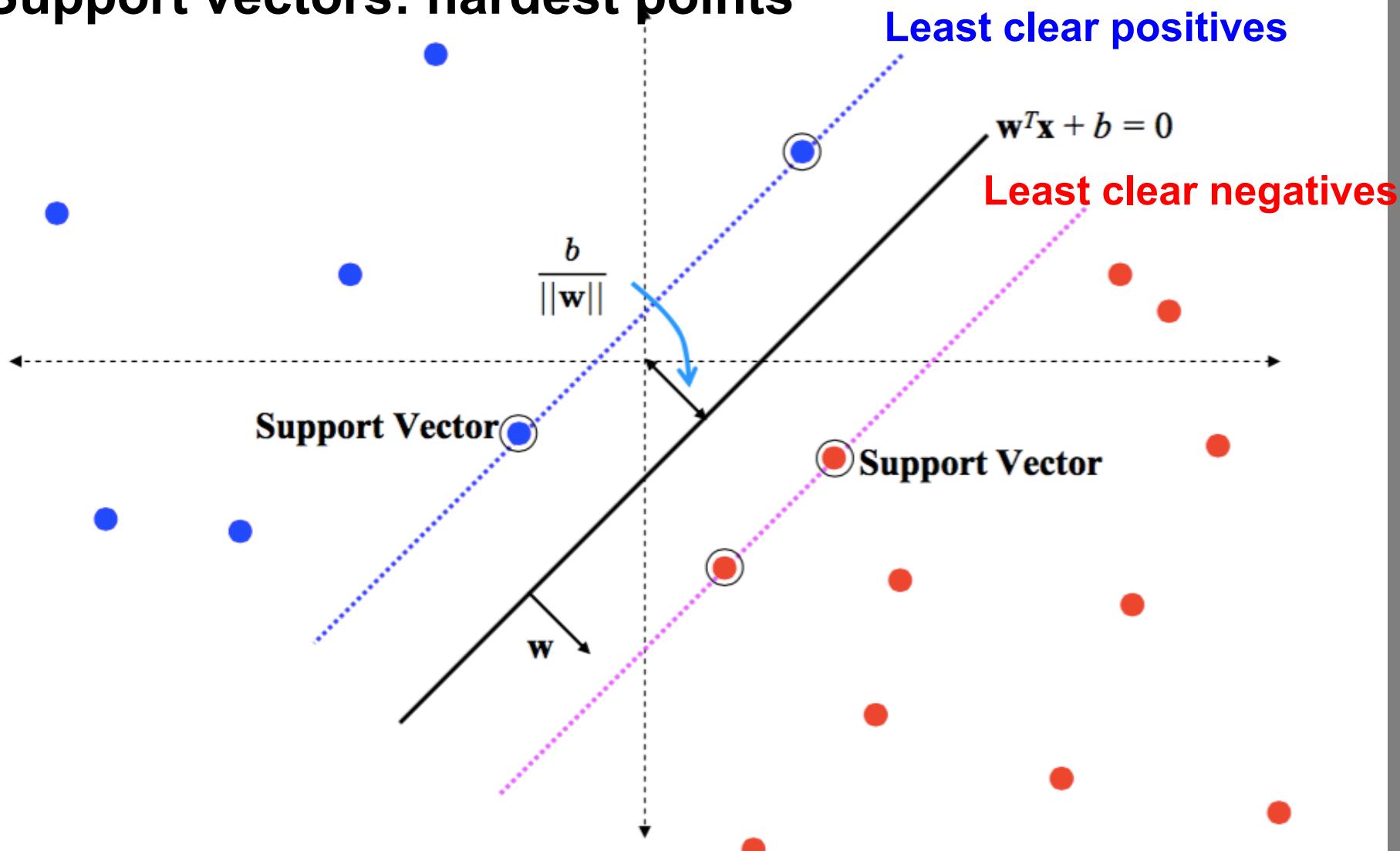


Should we be optimizing the mean, max, min margin?

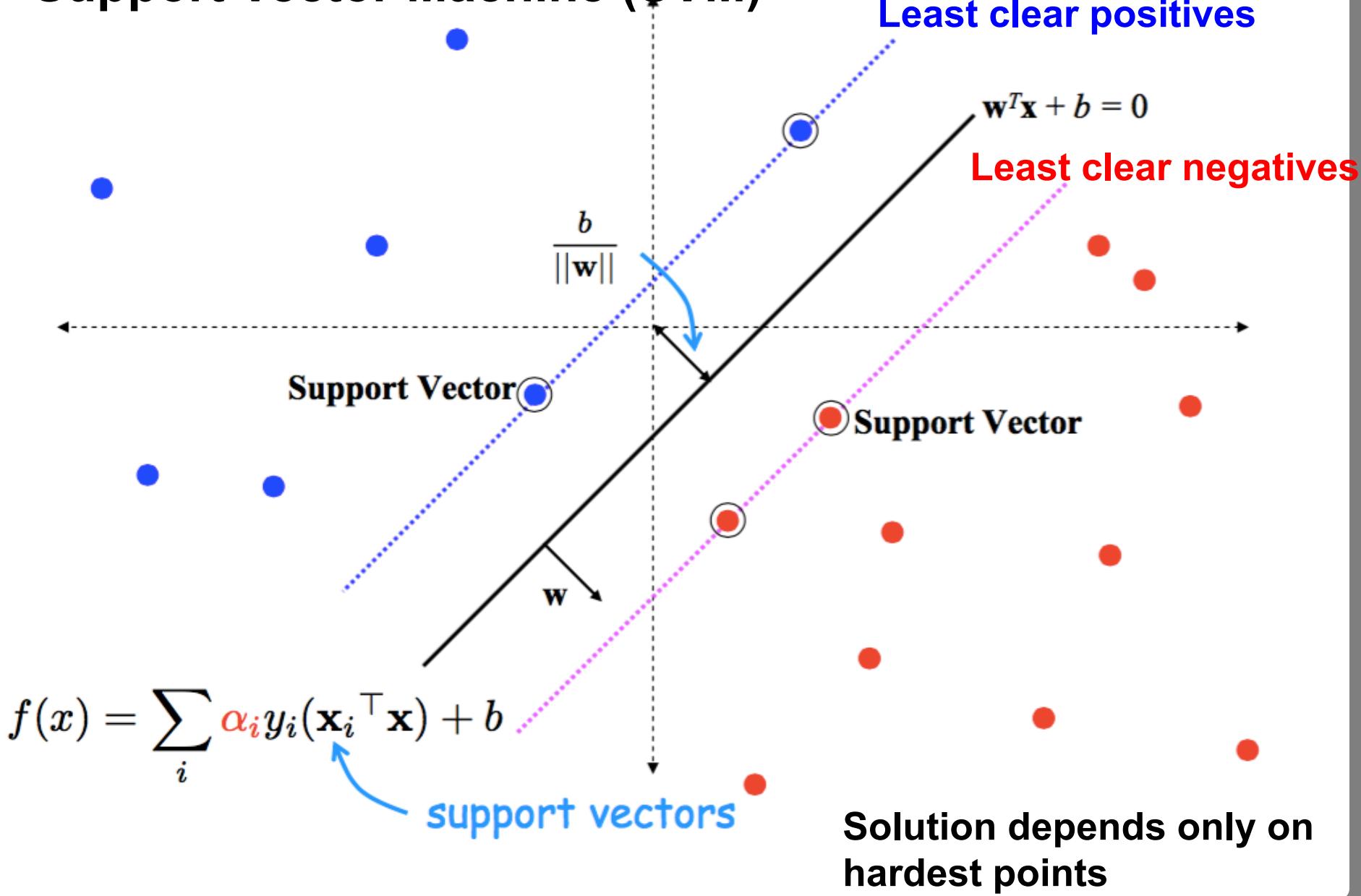
All points should lie **clearly** on the correct side of the boundary

- 1) Take points that do not lie clearly on the correct side
- 2) Make sure they do

# Support vectors: hardest points



# Support Vector Machine (SVM)



# SVM, sketch of derivation

- Since  $\mathbf{w}^\top \mathbf{x} + b = 0$  and  $c(\mathbf{w}^\top \mathbf{x} + b) = 0$  define the same plane, we have the freedom to choose the normalization

# SVM, sketch of derivation

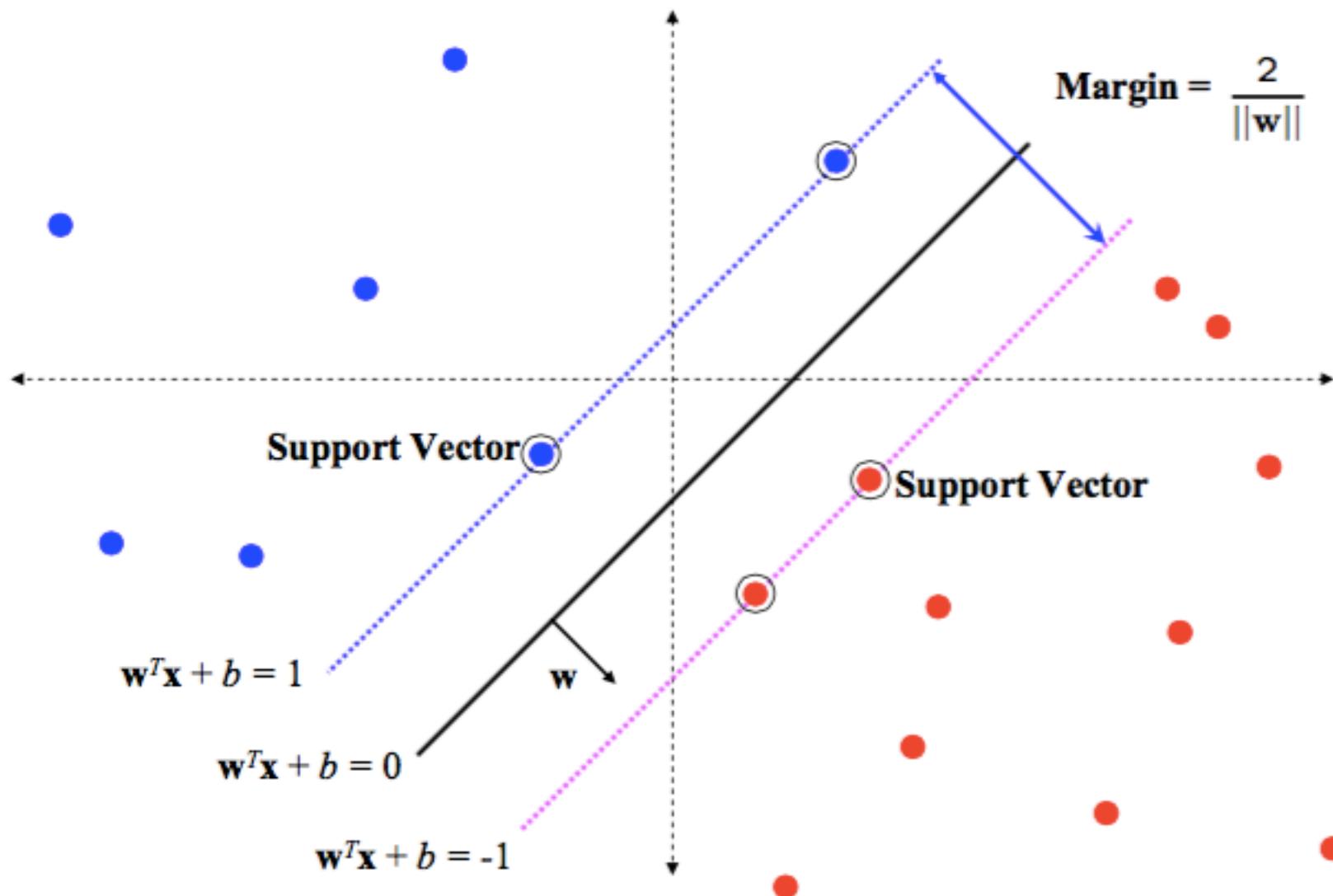
- Since  $\mathbf{w}^\top \mathbf{x} + b = 0$  and  $c(\mathbf{w}^\top \mathbf{x} + b) = 0$  define the same plane, we have the freedom to choose the normalization
- Choose normalization such that  $\mathbf{w}^\top \mathbf{x}_+ + b = +1$  and  $\mathbf{w}^\top \mathbf{x}_- + b = -1$  for the positive and negative support vectors respectively

# SVM, sketch of derivation

- Since  $\mathbf{w}^\top \mathbf{x} + b = 0$  and  $c(\mathbf{w}^\top \mathbf{x} + b) = 0$  define the same plane, we have the freedom to choose the normalization
- Choose normalization such that  $\mathbf{w}^\top \mathbf{x}_+ + b = +1$  and  $\mathbf{w}^\top \mathbf{x}_- + b = -1$  for the positive and negative support vectors respectively
- Then the margin is given by

$$\frac{\mathbf{w}^\top (\mathbf{x}_+ - \mathbf{x}_-)}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

# Support Vector Machine (SVM)



# Representer theorem

Objective: find  $\mathbf{w}$  that maximizes the margin subject to margin constraints

$$\min_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

$$\text{s.t. } y^i (\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \quad \forall i$$

Equivalently:

$$\min_{\mathbf{w}} \|\mathbf{w}\|^2$$

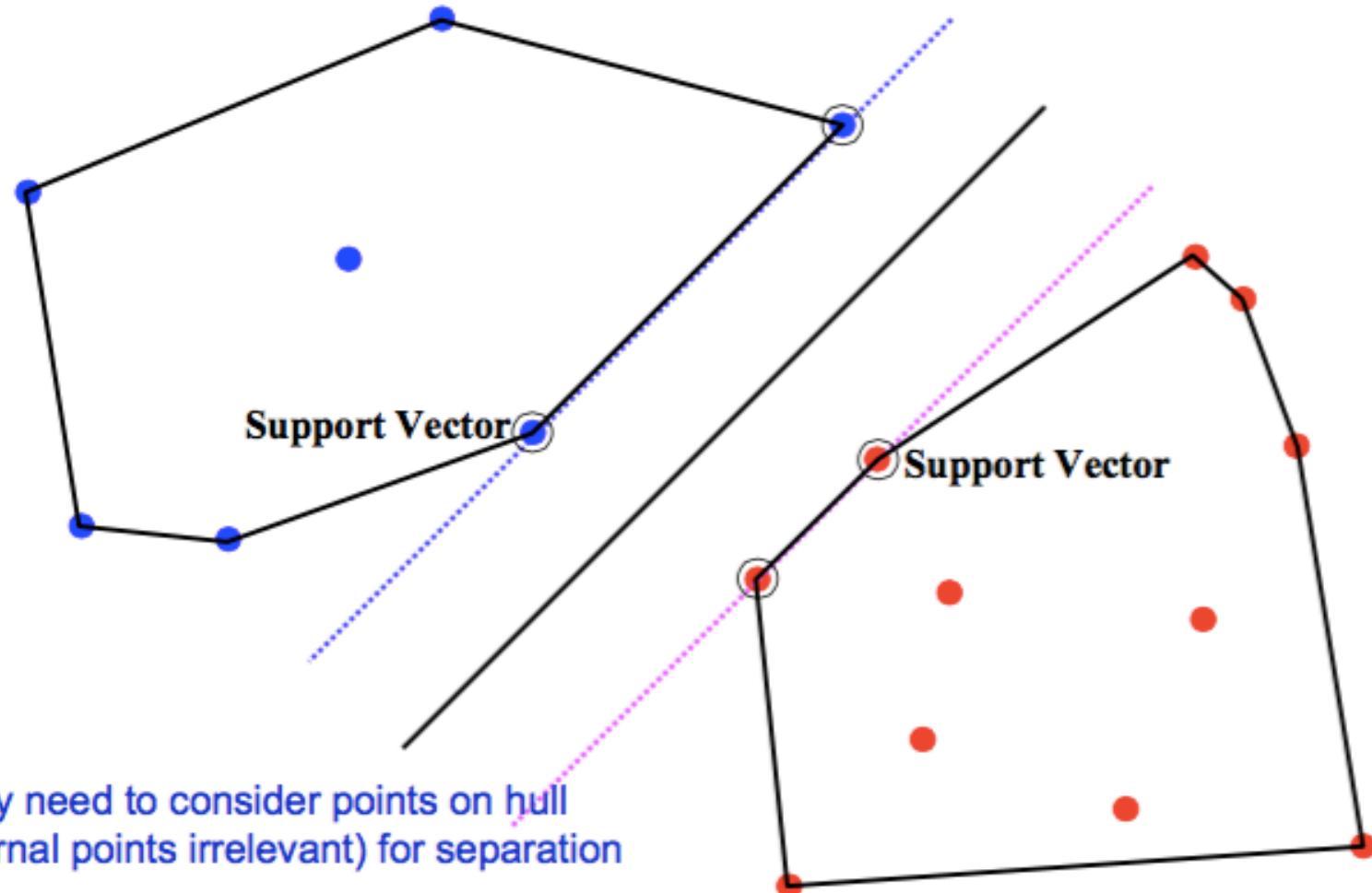
$$\text{s.t. } \underset{N}{\sum} y^i (\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \quad \forall i$$

Representer Theorem:  $\mathbf{w}^* = \sum_{i=1} \alpha^i (y^i \mathbf{x}^i)$

## Geometric algorithm

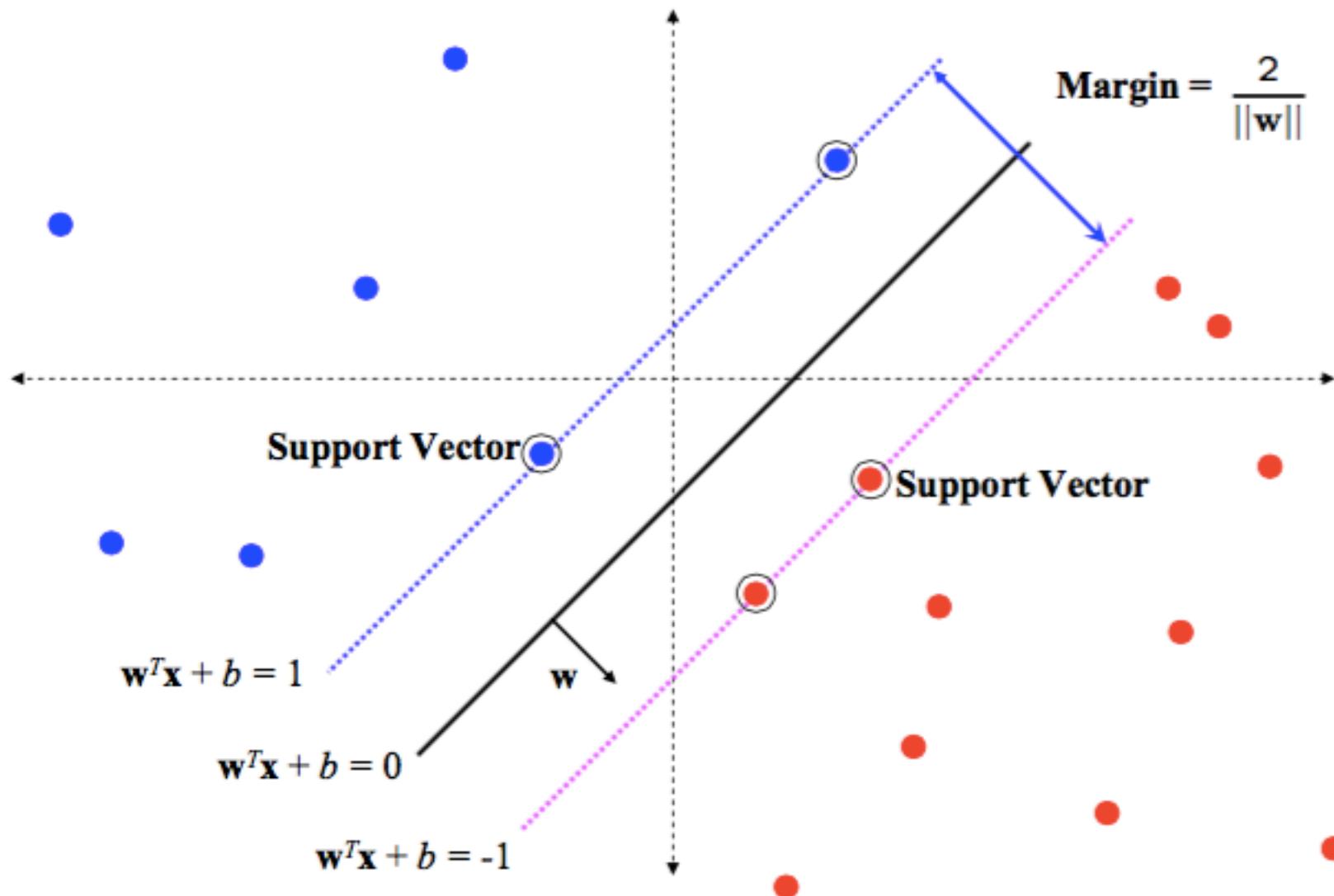
- Compute the convex hull of the positive points, and the convex hull of the negative points
- For each pair of points, one on positive hull and the other on the negative hull, compute the margin
- Choose the largest margin

# Intuitive justification of theorem



- only need to consider points on hull  
(internal points irrelevant) for separation

# Support Vector Machine (SVM)



# Primal and dual problems

Primal, in terms of  $\mathbf{w}$ :  $\min_{\mathbf{w}} \|\mathbf{w}\|^2$   
 s.t. :  $y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1, \quad \forall i$

**But:**  $\|\mathbf{w}^*\|^2 = \langle \mathbf{w}^*, \mathbf{w}^* \rangle$

$$\begin{aligned} \mathbf{w}^* &= \sum_{i=1}^N \alpha^i (y^i \mathbf{x}^i) \\ &= \left\langle \sum_{i=1}^N \alpha^i y^i \mathbf{x}^i, \sum_{j=1}^N \alpha^j y^j \mathbf{x}^j \right\rangle = \sum_{i=1}^N \sum_{j=1}^N \alpha^i \alpha^j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \end{aligned}$$

Dual, in terms of  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)$ :  $\min_{\boldsymbol{\alpha}} \sum_{i=1}^N \sum_{j=1}^N \alpha^i \alpha^j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$   
 s.t. :  $y^i \left( \sum_{j=1}^N \alpha^j y^j \langle \mathbf{x}^j, \mathbf{x}^i \rangle + b \right) \geq 1, \quad i = 1, \dots, N$

## Primal vs dual

Primal:  $\min_{\mathbf{w}} \|\mathbf{w}\|^2$   $\mathbf{w} \in \mathbb{R}^D \rightarrow O(D^3)$

s.t. :  $y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1, \quad \forall i$

Dual:  $\min_{\boldsymbol{\alpha}} \sum_{i=1}^N \sum_{j=1}^N \alpha^i \alpha^j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$   $\boldsymbol{\alpha} \in \mathbb{R}^N \rightarrow O(N^3)$

s.t. :  $y^i \left( \sum_{j=1}^N \alpha^j y^j \langle \mathbf{x}^j, \mathbf{x}^i \rangle + b \right) \geq 1, \quad \forall i$

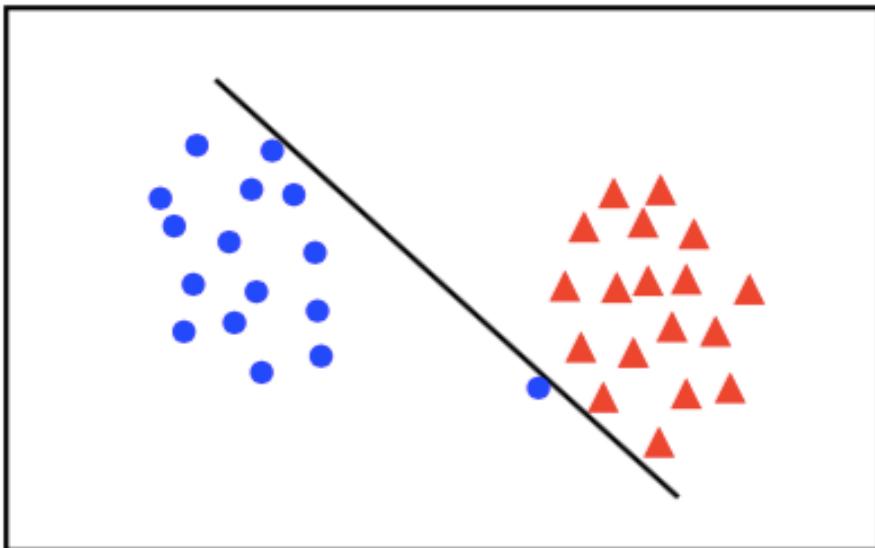
**Dual can be faster if N<D!**

Primal and dual classifier forms:

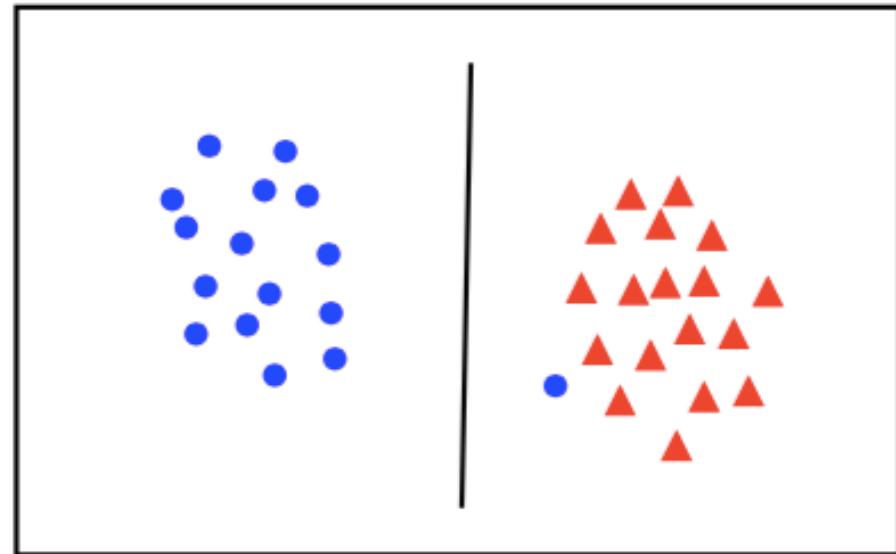
$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i=1}^N \alpha^i y^i \langle \mathbf{x}^i, \mathbf{x} \rangle + b$$

**Dual form involves only inner products of features ( $\Rightarrow$  kernel trick)**

# What is the “best” decision plane?



**All points on the  
correct side!**



**But this looks  
better overall!**

Best: understood at test time

Maybe we could sacrifice classifying some training points correctly

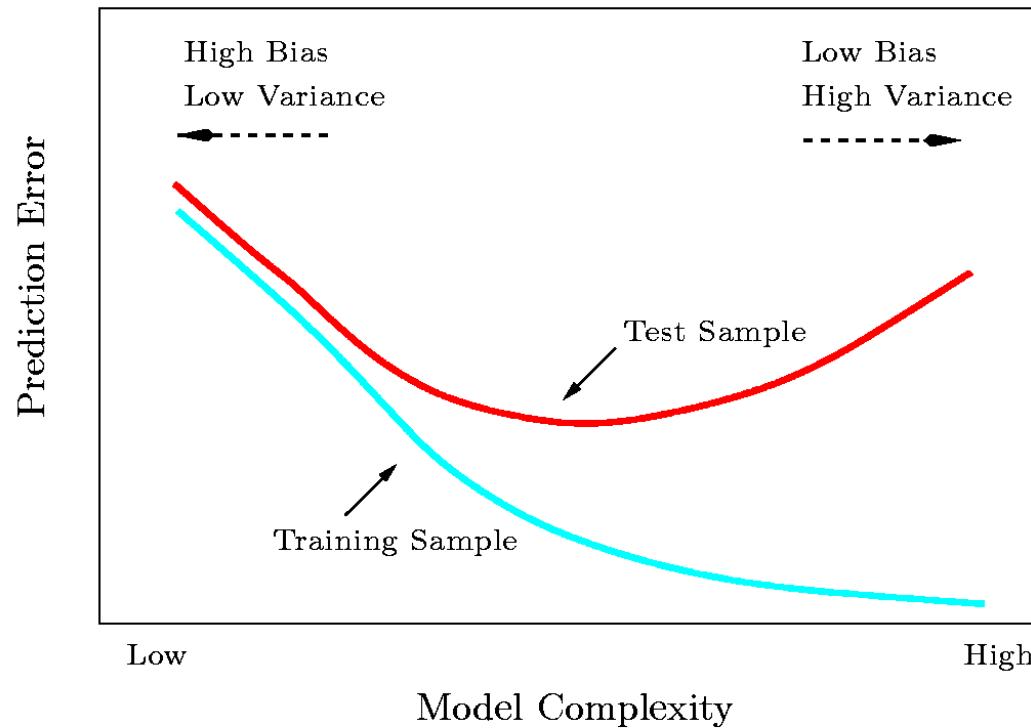
# Tuning the model's complexity

A flexible model approximates the target function well in the training set

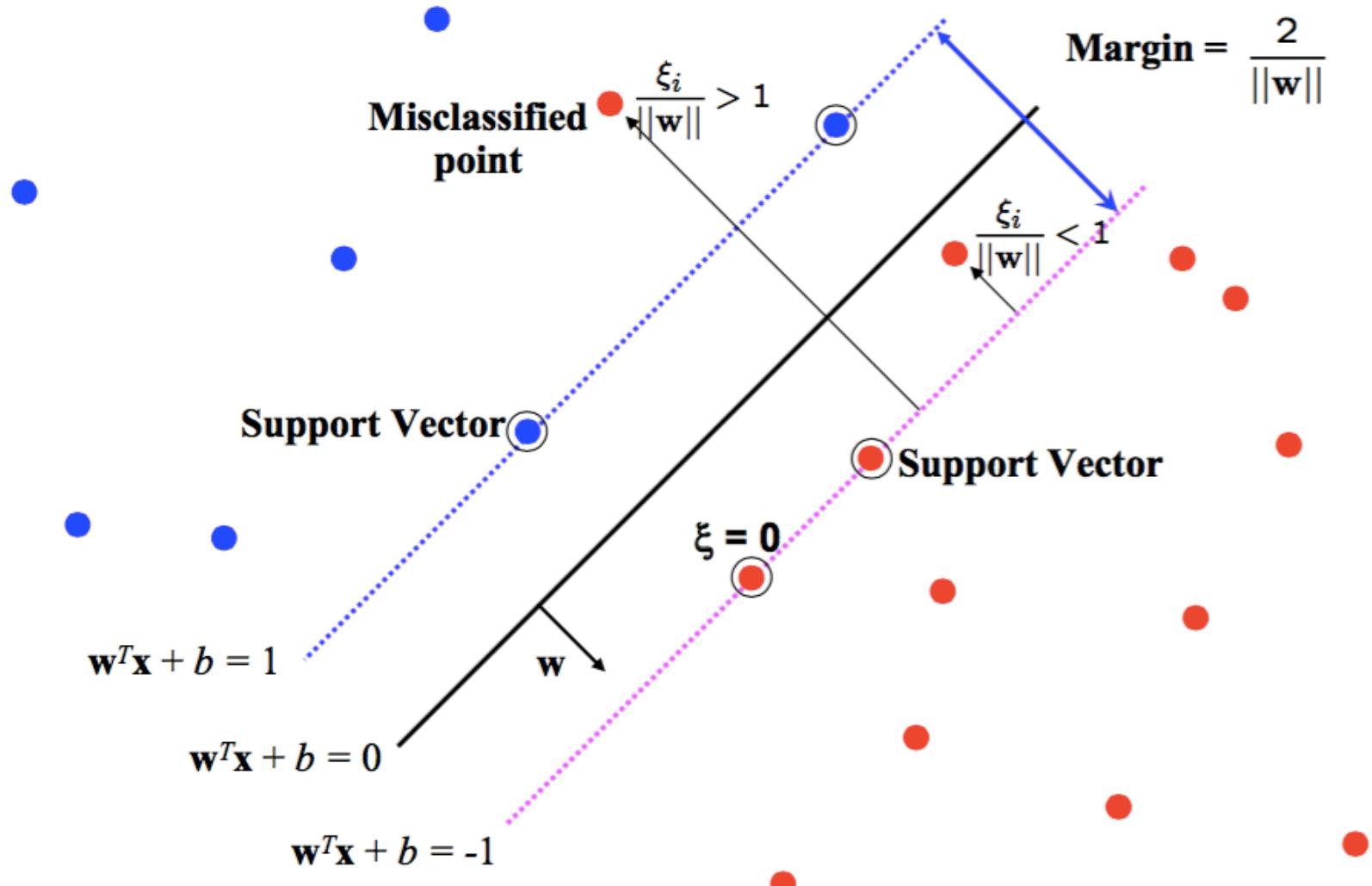
*but can “overtrain” and have poor performance on the test set (“variance”)*

A rigid model’s performance is more predictable in the test set

*but the model may not be good even on the training set (“bias”)*



# Slack variables: let us make (but also pay) some errors



# Objective for non-separable data

$$\min_{\mathbf{w}, \xi} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi^i$$

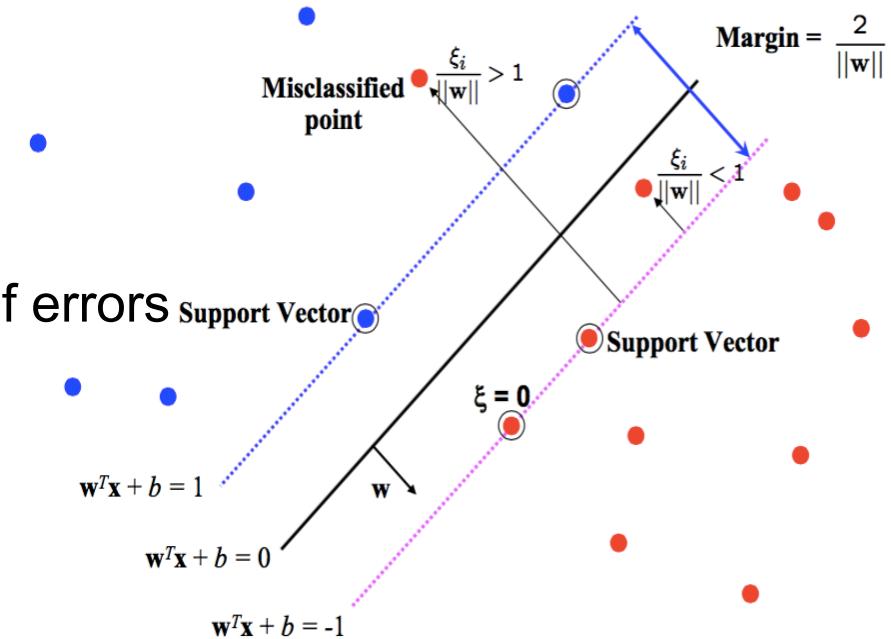
s.t. :  $y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi^i, \quad \forall i$

$$\xi^i \geq 0, \quad \forall i$$

misclassification when  $\xi > 1$

$\sum_i \xi^i$  : upper bound on number of errors

C: hyperparameter  
(cross-validation!)



Primal problem: for  $\mathbf{w} \in \mathbb{R}^d$

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \|\mathbf{w}\|^2 + C \sum_i^N \xi_i \text{ subject to } y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \text{ for } i = 1 \dots N$$

The constraint  $y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i$ , can be written more concisely as

$$y_i f(\mathbf{x}_i) \geq 1 - \xi_i$$

which is equivalent to

$$\xi_i = [1 - y_i f(\mathbf{x}_i)]_+$$

where  $[.]_+$  indicates the positive part. Hence the optimization problem is equivalent to

$$\min_{\mathbf{w} \in \mathbb{R}^d} \underbrace{\|\mathbf{w}\|^2}_{\text{regularization}} + C \underbrace{\sum_i^N [1 - y_i f(\mathbf{x}_i)]_+}_{\text{loss function}}$$

# Loss function

Optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi^i \\ \text{s.t.} \quad & y^i (\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi^i \\ & \xi^i \geq 0 \end{aligned}$$

Rewrite constraint:  $y^i h_{\mathbf{w}, b}(\mathbf{x}) \geq 1 - \xi^i$

Compact form:  $\xi^i = [1 - y^i h_{\mathbf{w}, b}(\mathbf{x})]_+$

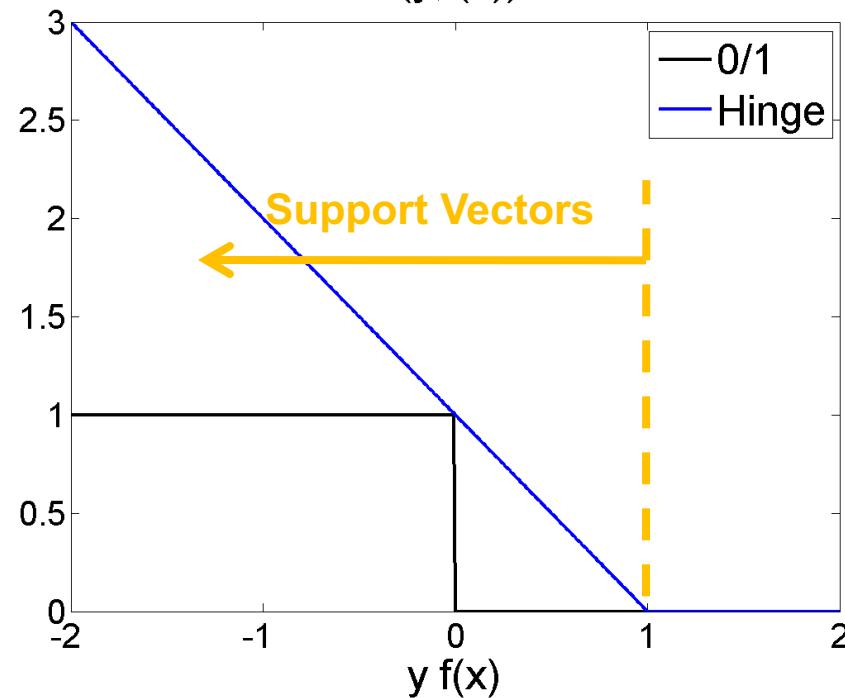
$$= \max(1 - y^i h_{\mathbf{w}, b}(\mathbf{x}), 0)$$

What if we plug that in the optimization objective?

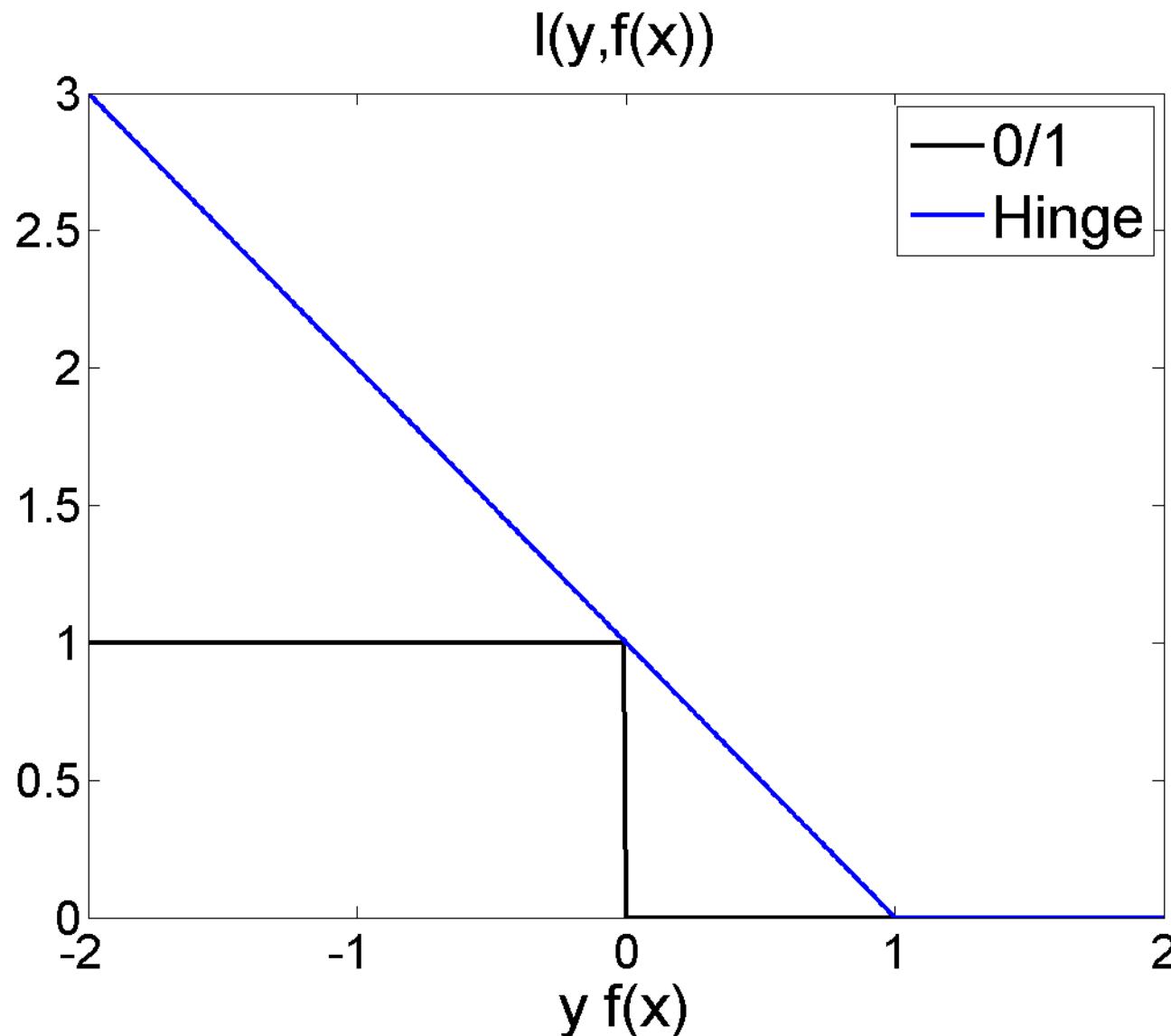
# Loss function

$$\begin{aligned}
 \text{Optimization problem: } L(\mathbf{w}) &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y^i h_{\mathbf{w}, b}(x^i)) \\
 &\propto \lambda \|\mathbf{w}\|^2 + \underbrace{\sum_{i=1}^N \max(0, 1 - y^i h_{\mathbf{w}, b}(x^i))}_{l(y^i, x^i)} \\
 &\quad \text{regularizer} \qquad \qquad \qquad \text{additive loss}
 \end{aligned}$$

Hinge loss:

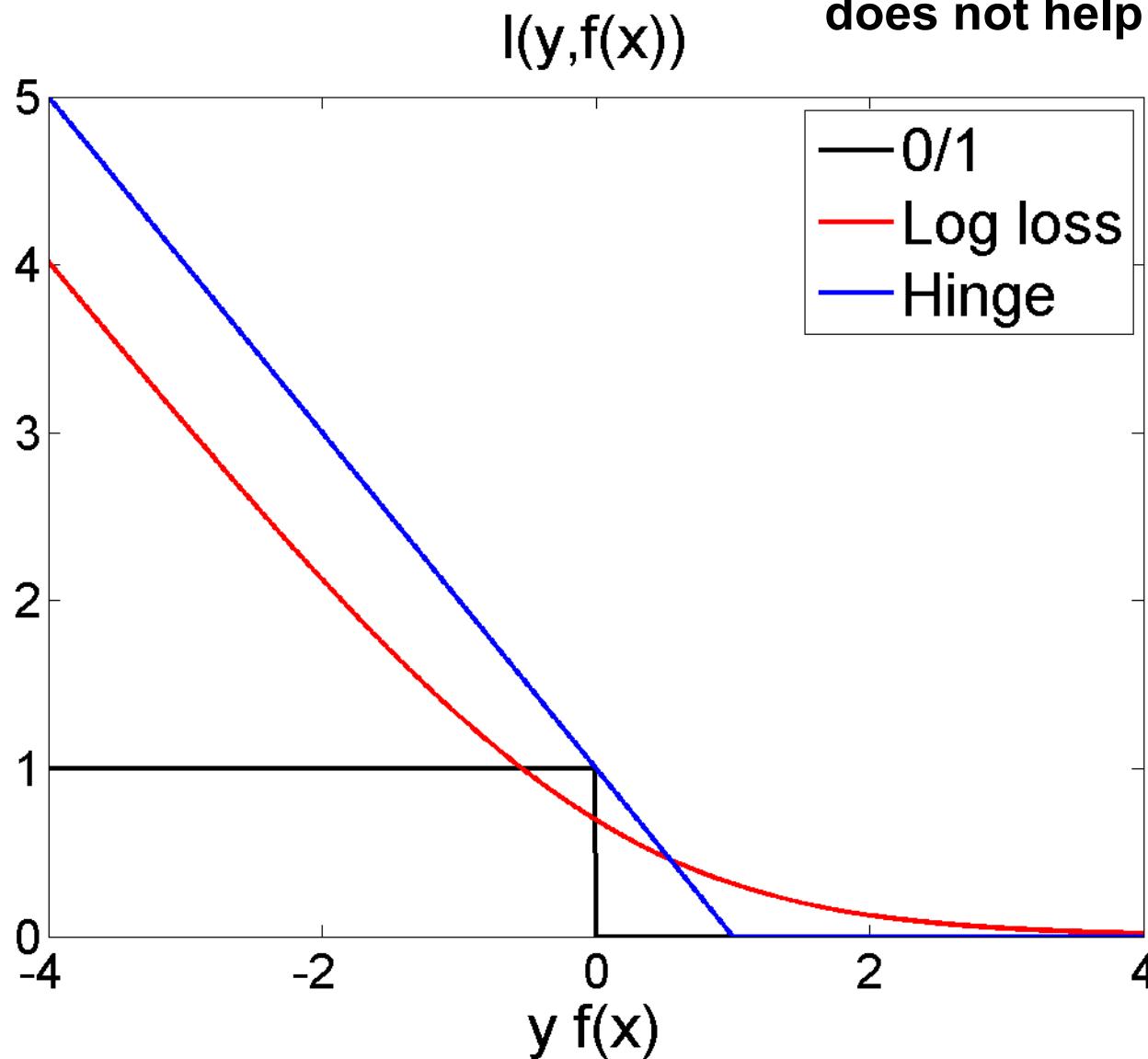


# Hinge loss

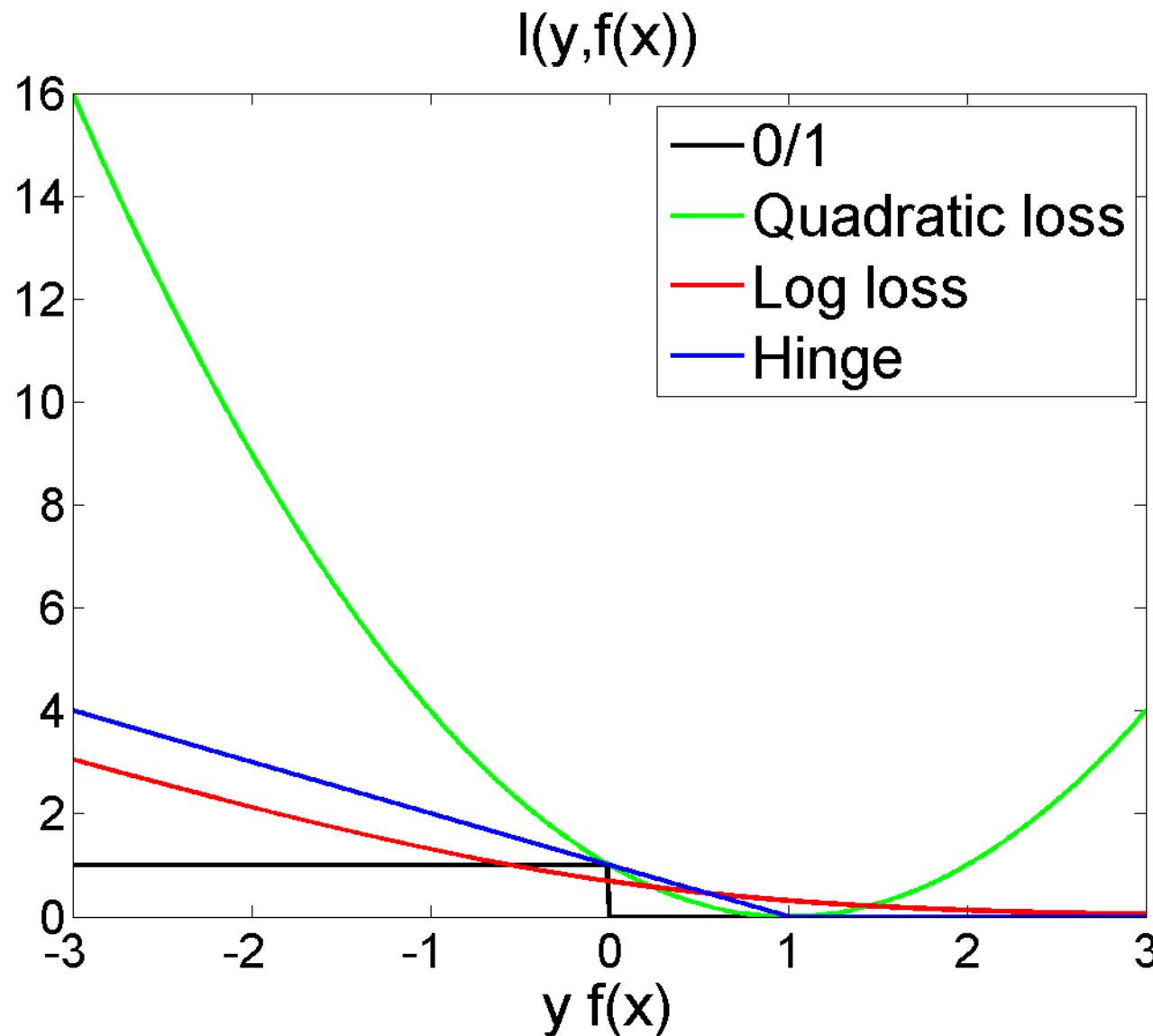


# Hinge loss vs log-loss

getting larger than 1:  
does not harm, but also  
does not help



# Hinge loss vs log-loss vs quadratic





# Lecture outline

Recap

Large margins and generalization

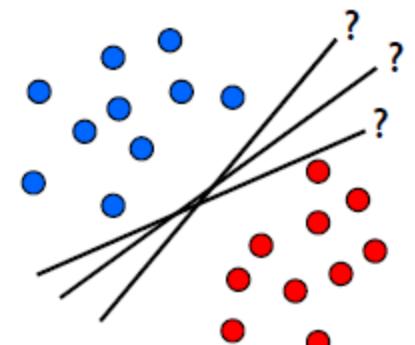
Optimization

Kernels

Applications to vision

# Generalization Error

- What is model complexity?
  - Number of parameters, magnitude of discriminant  $w$ ?
  - Analyze complexity of hypothesis class
  
- Linear classifiers:
  - Different decision boundaries
    - Different generalization performance
  - Test error > training error
  - Which line gives smallest test error?



# Learning Theory

- V. Vapnik, 1968
  - Mainstream Statistics: Large-sample analysis ('in the limit')
  - Pattern Recognition: Small sample properties
- Distribution-free bounds on worst performance

# Empirical and Actual risk

- Empirical risk
  - Measured on the training/validation set

$$R_{emp}(\alpha) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i; \alpha))$$

- Actual risk (= Expected risk)
  - Expectation of the error on *all* data.

$$R(\alpha) = \int L(y_i, f(\mathbf{x}; \alpha)) dP_{X,Y}(\mathbf{x}, y)$$

- $P_{X,Y}(\mathbf{x}, y)$  is the probability distribution of  $(\mathbf{x}, y)$ .  
It is fixed, but typically unknown.

# Actual and Empirical Risk

- Idea

- Compute an upper bound on the actual risk based on the empirical risk

$$R(\alpha) \leq R_{emp}(\alpha) + \epsilon(N, p^*, h)$$

- where

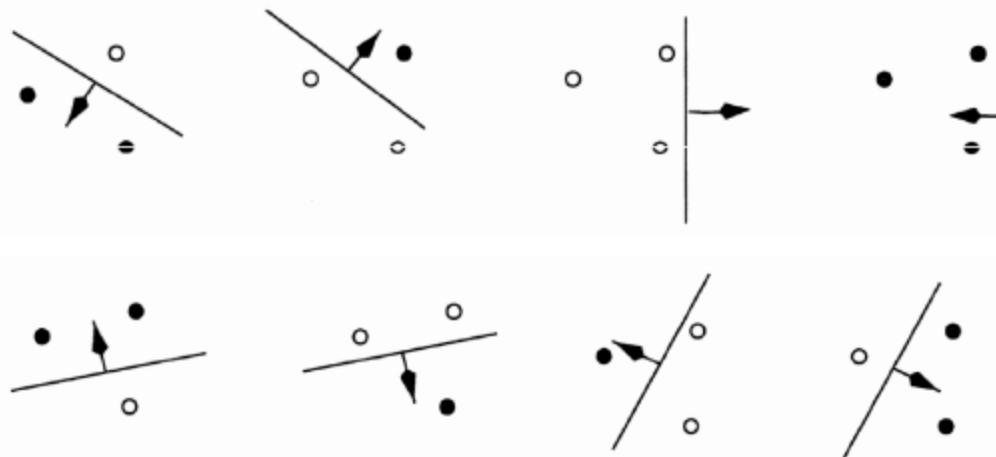
$N$ : number of training examples

$p^*$ : probability that the bound is correct

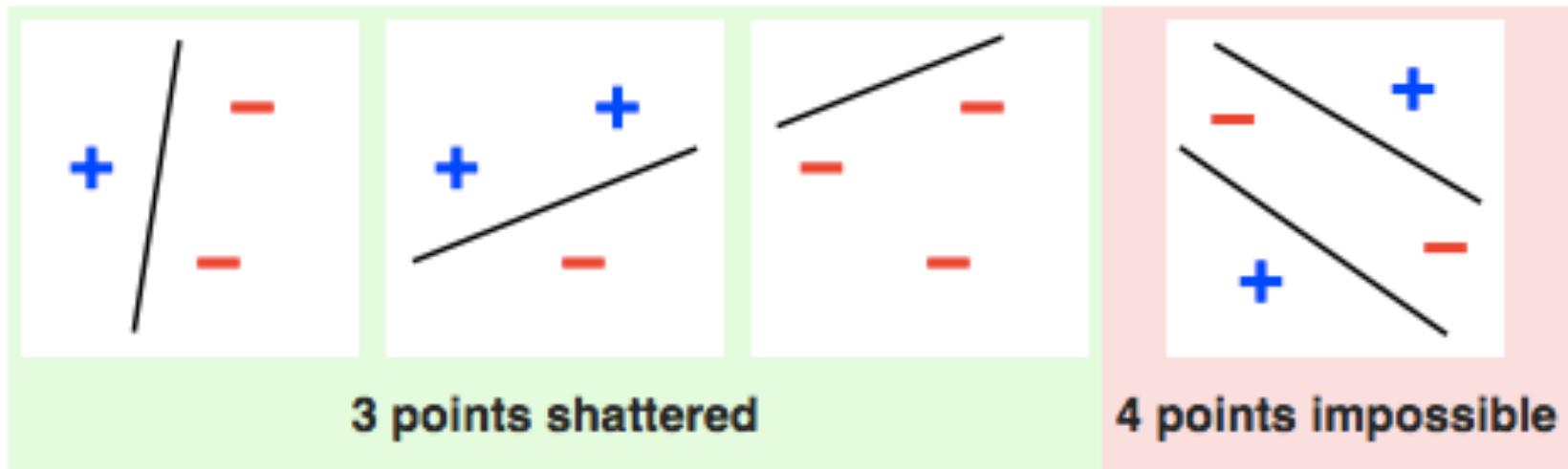
$h$ : capacity of the learning machine (“VC-dimension”)

# Vapnik Chervonenkis (VC) Dimension

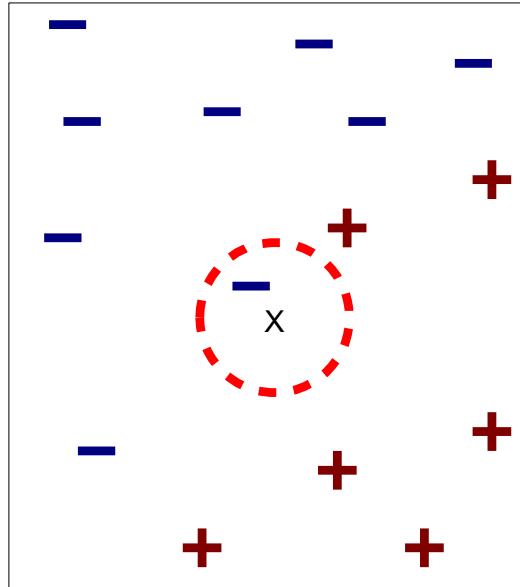
- Shattering: *If a given set of  $\ell$  points can be labeled in all possible  $2^\ell$  ways, and for each labeling, a member of the set  $\{f(\alpha)\}$  can be found which correctly assigns those labels, we say that the set of points is shattered by the set of functions.*
- VC dimension *The VC dimension for the set of functions  $\{f(\alpha)\}$  is defined as the maximum number of training points that can be shattered by  $\{f(\alpha)\}$ .*
- Example



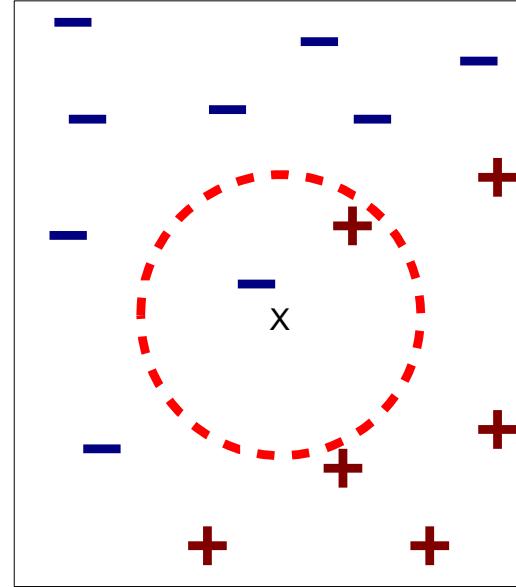
# Arbitrary linear classifier in N-dimensions: VC-dim= N+1



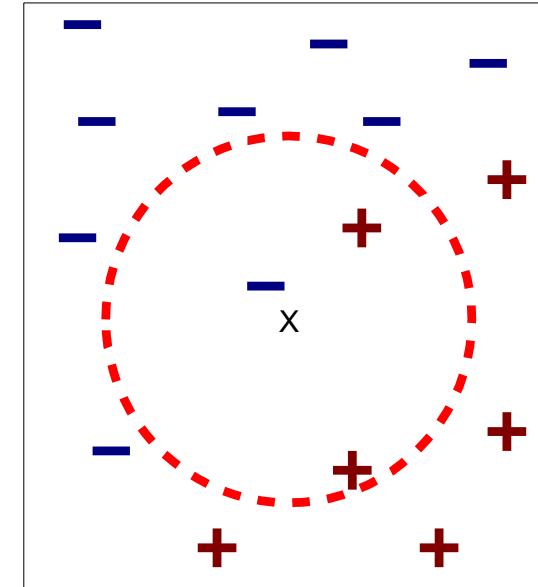
# Reminder: K-nearest neighbor classifier



(a) 1-nearest neighbor



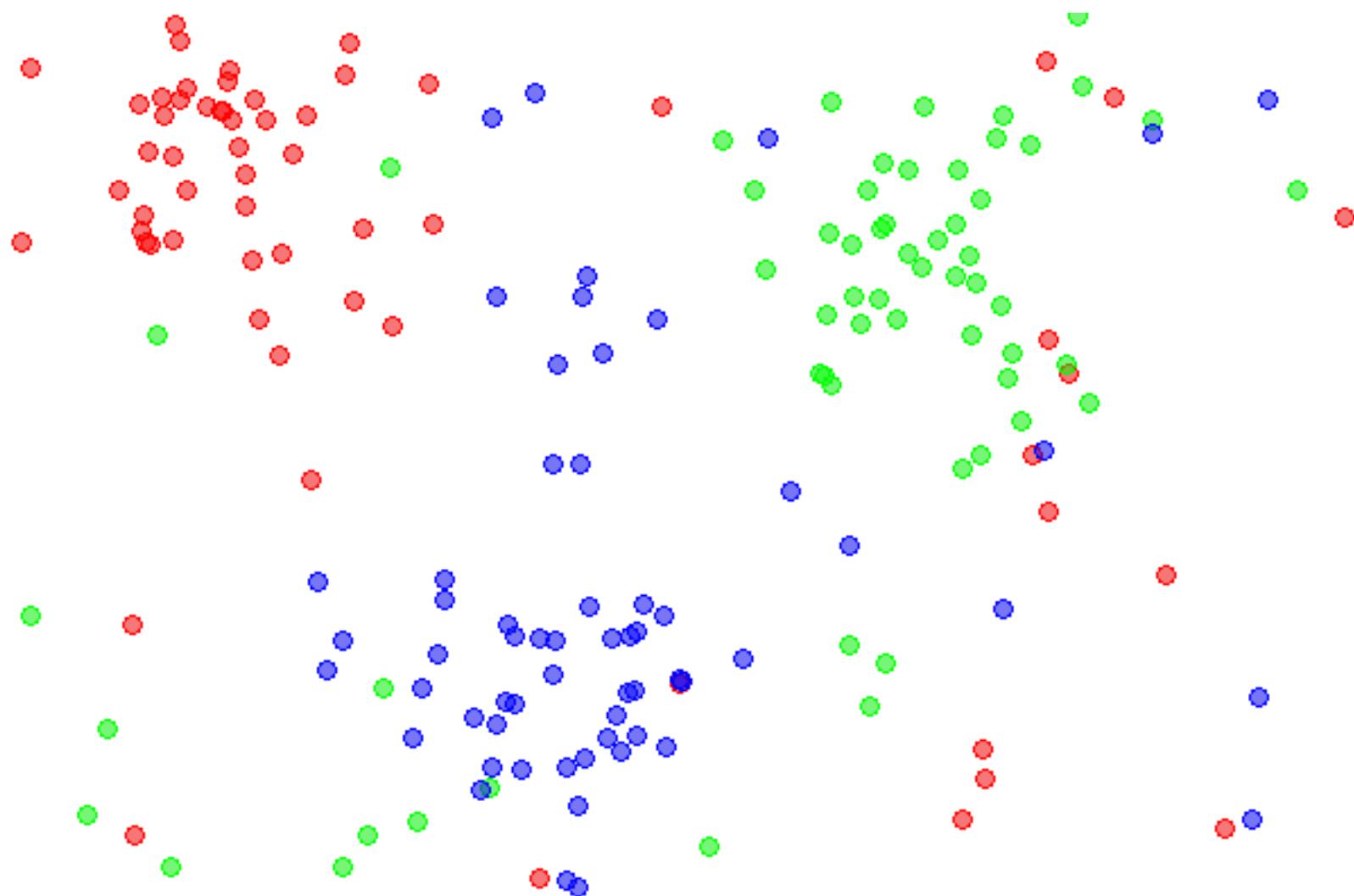
(b) 2-nearest neighbor



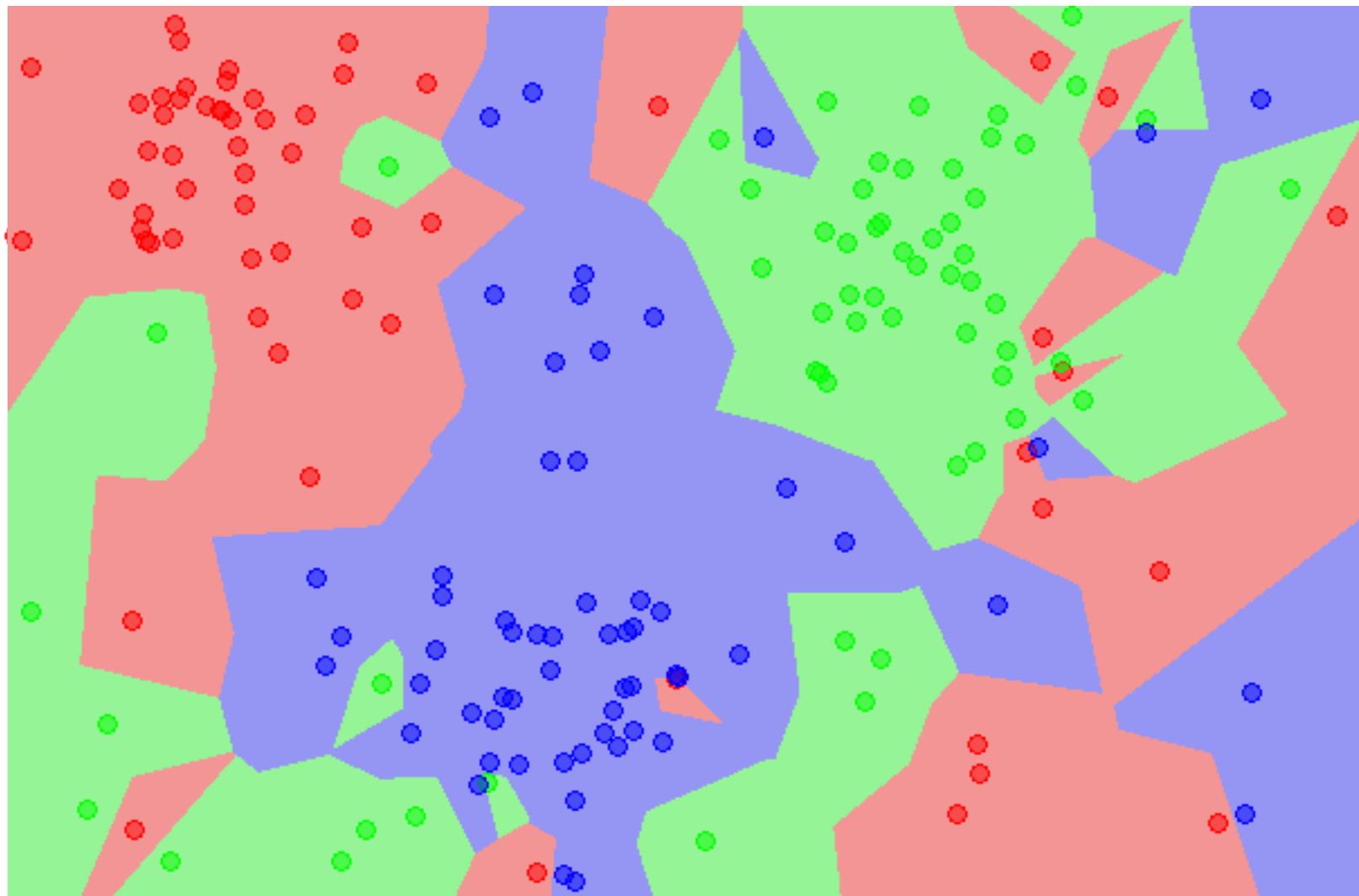
(c) 3-nearest neighbor

- Compute distance to other training records
- Identify  $K$  nearest neighbors
- Take majority vote

# Training data for NN classifier (in $\mathbb{R}^2$ )

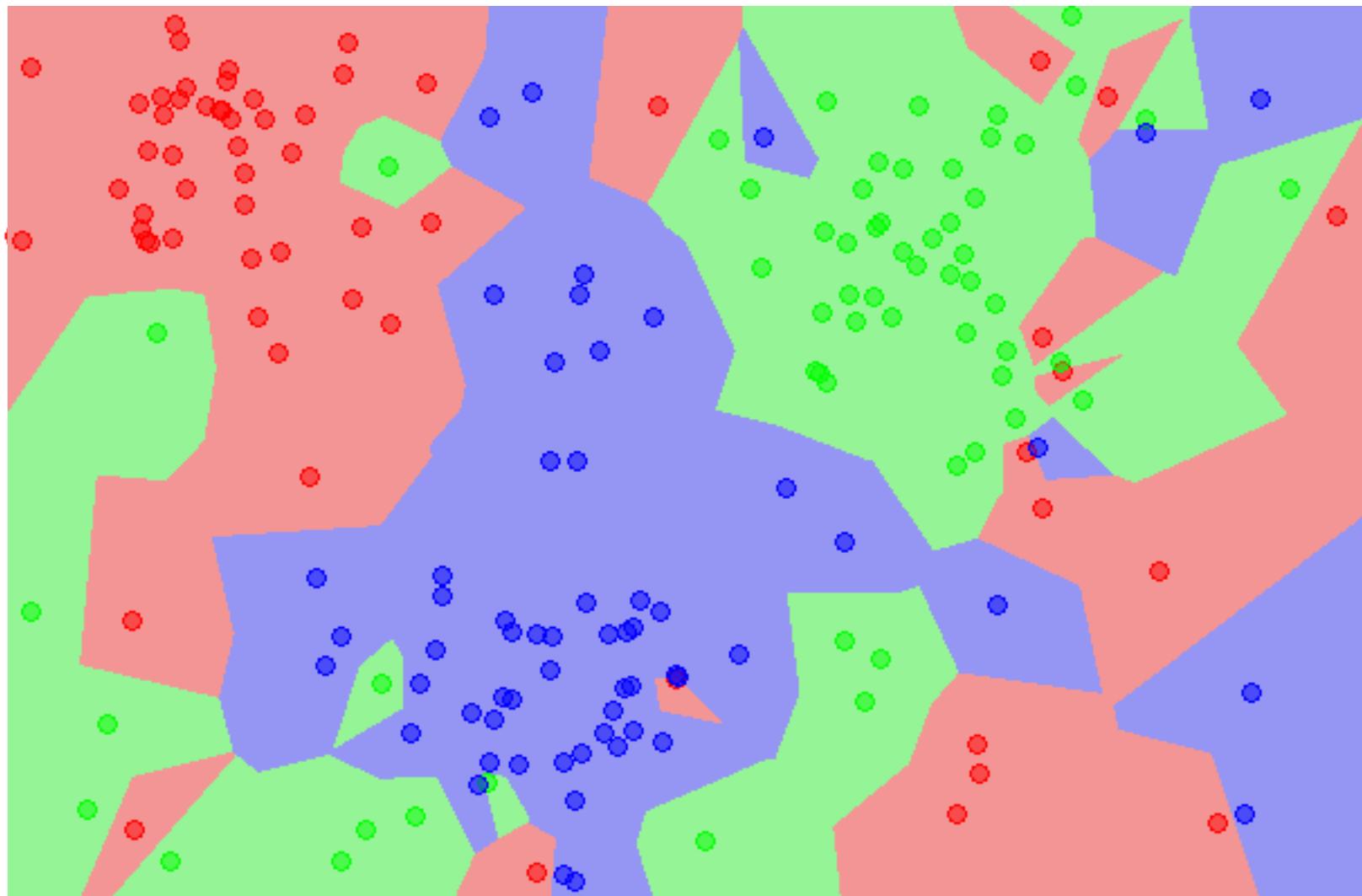


# 1-nn classifier prediction (in $\mathbb{R}^2$ )



What is the VC dimension of this classifier?

# 1-nn classifier prediction (in $\mathbb{R}^2$ )



What is the VC dimension of this classifier?

# Large Margins & VC Dimension

- Vapnik: *The class of optimal linear separators has VC dimension  $h$  bounded from above as*

$$h \leq \min \left\{ \left\lceil \frac{D^2}{\rho^2} \right\rceil, m_0 \right\} + 1$$

*where  $\rho$  is the margin,  $D$  is the diameter of the smallest sphere that can enclose all of the training examples, and  $m_0$  is the dimensionality.*

- If we maximize the margins, feature dimensionality does not matter

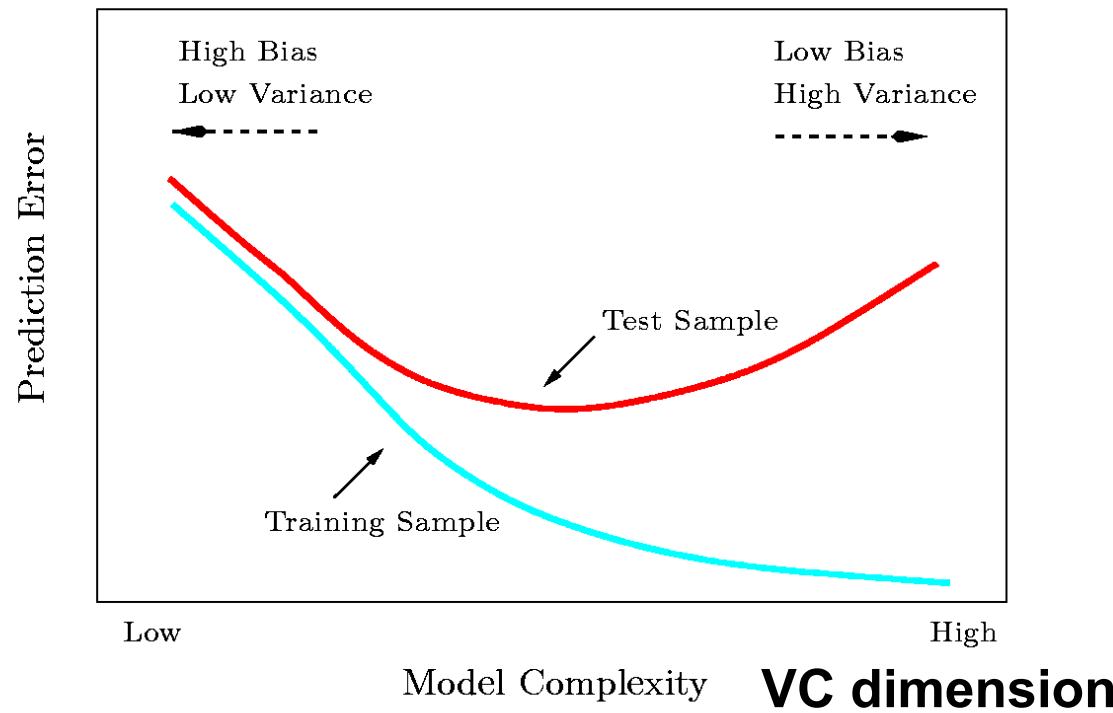
# Tuning the model's complexity

A flexible model approximates the target function well in the training set

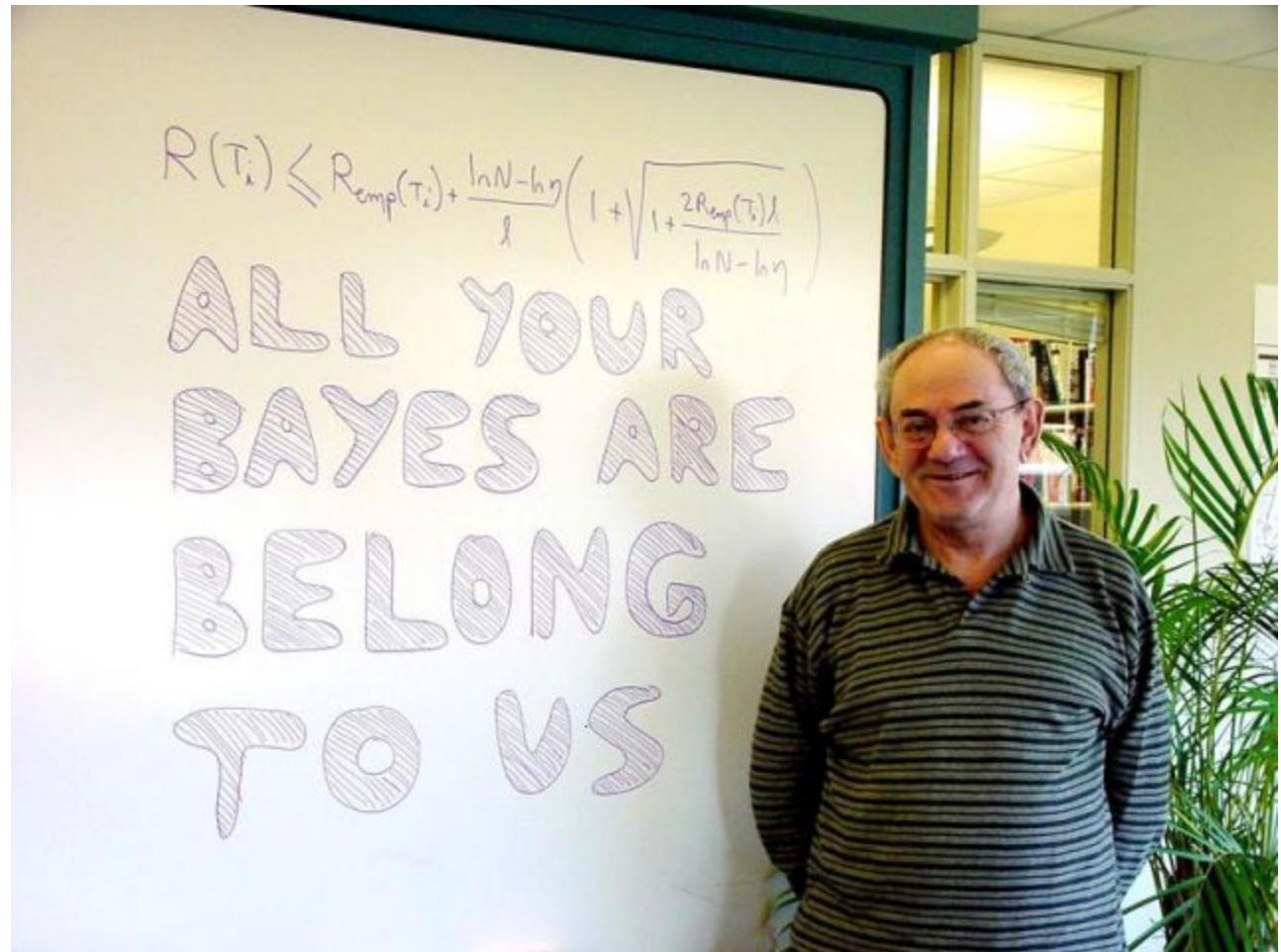
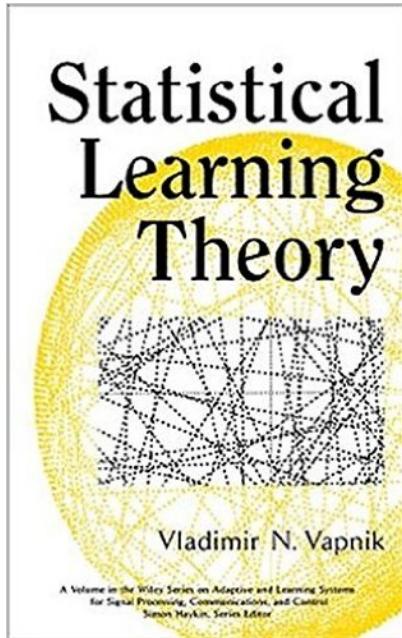
A rigid model's performance is more predictable in the test set

- With probability  $(1-\eta)$ , the following bound holds

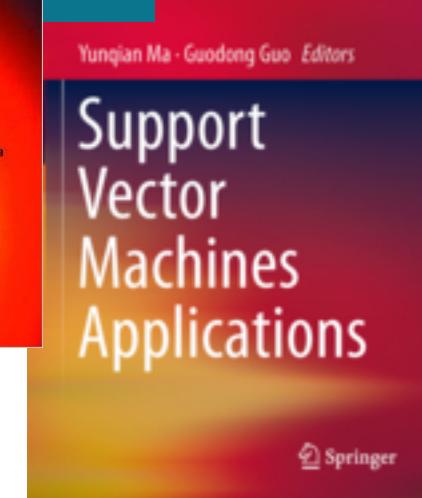
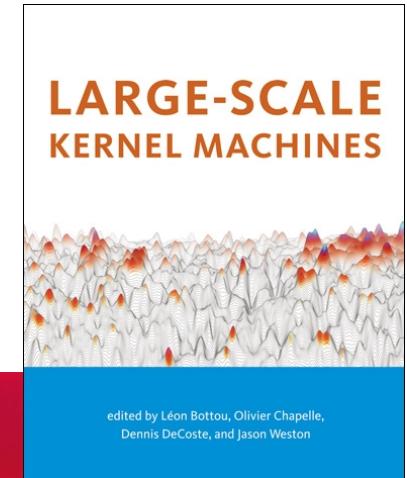
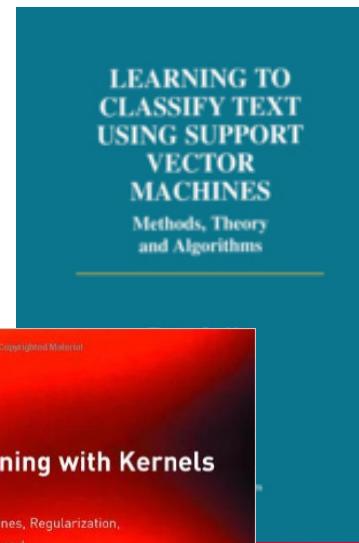
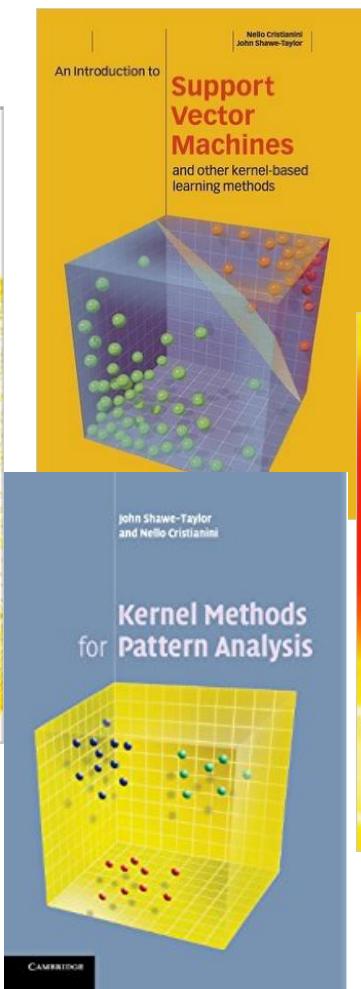
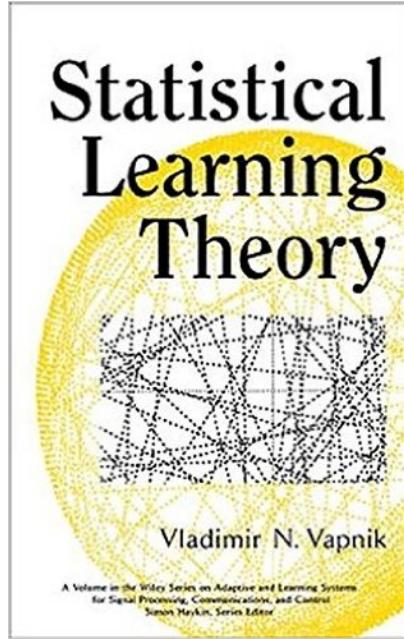
$$R(\alpha) \leq R_{emp}(\alpha) + \underbrace{\sqrt{\frac{h(\log(2N/h) + 1) - \log(\eta/4)}{N}}}_{\text{"VC confidence"}}$$



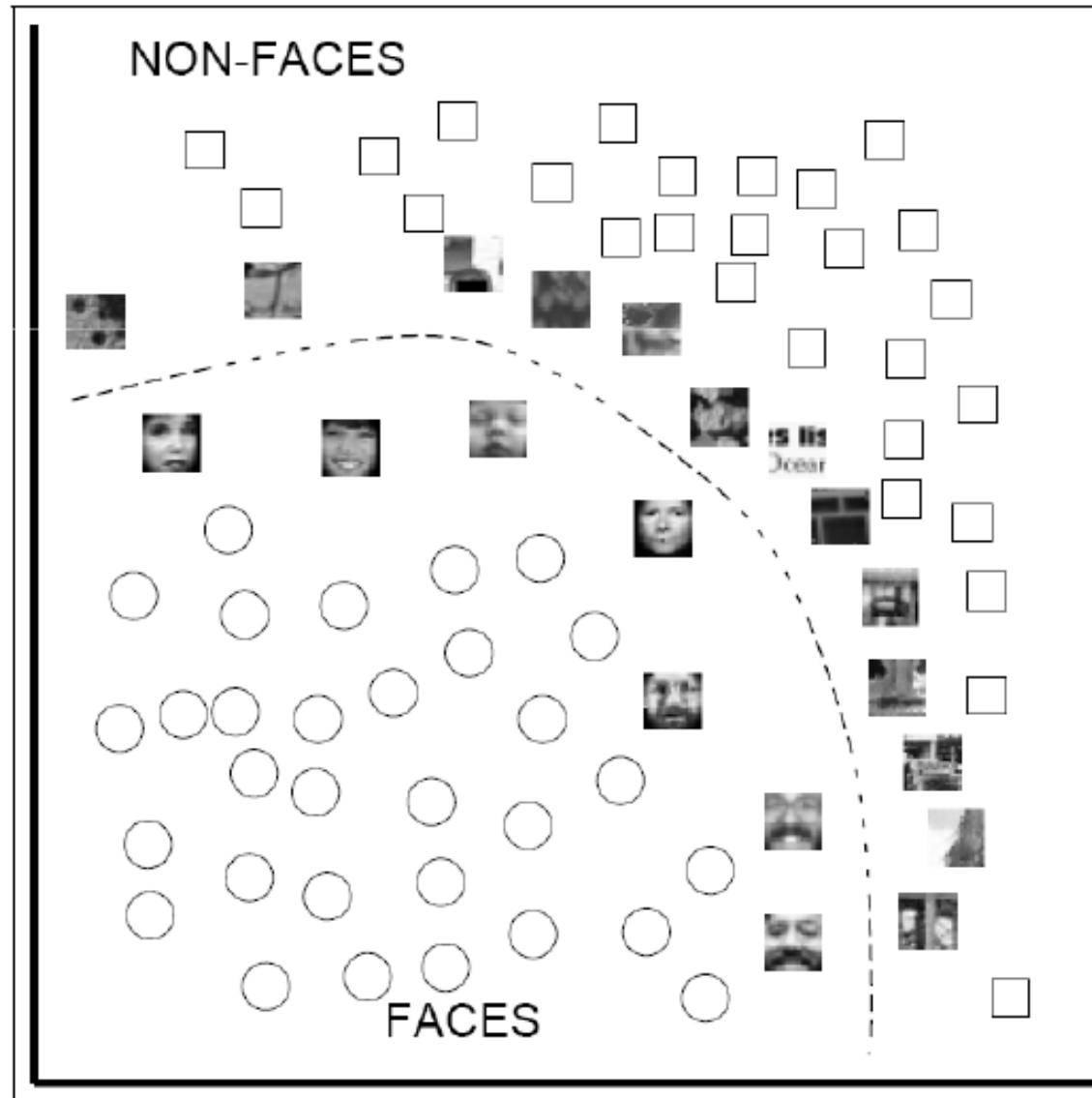
“There’s nothing more practical than a good theory”



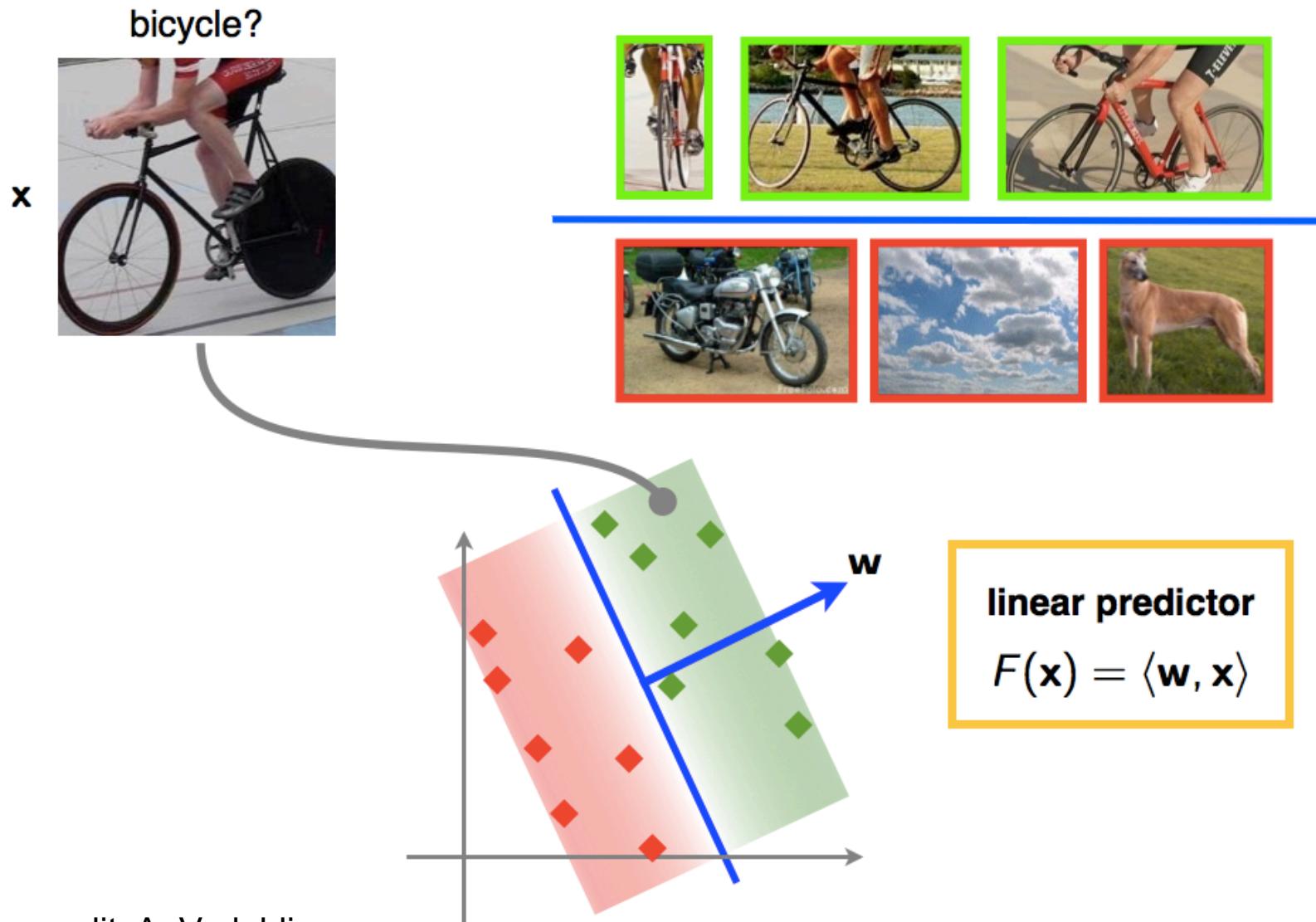
# “There’s nothing more practical than a good theory”



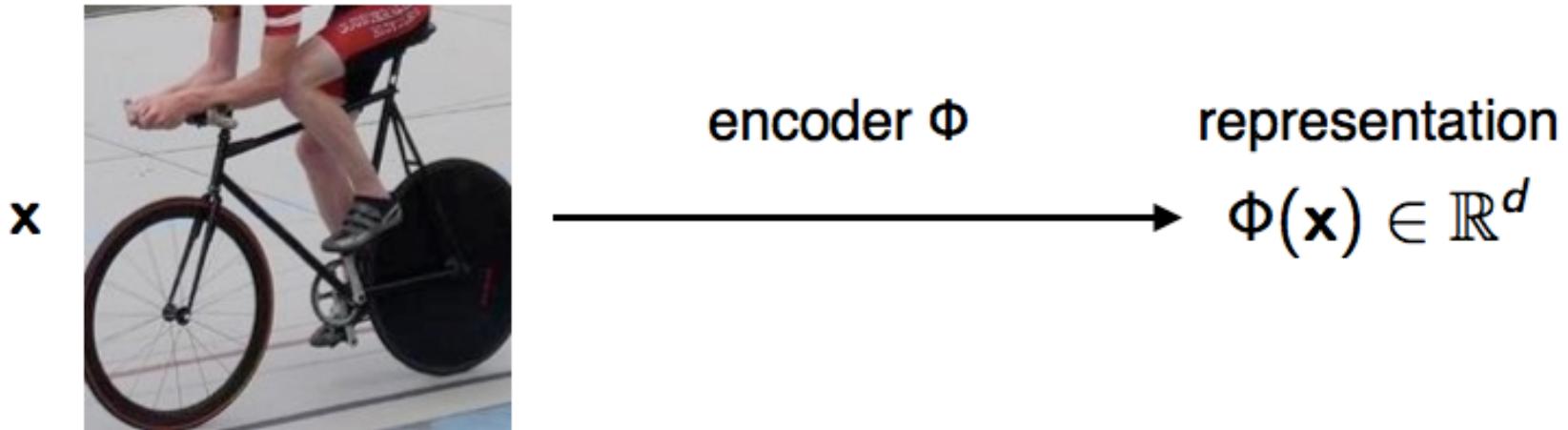
# Support vectors for Faces (P&P 98)



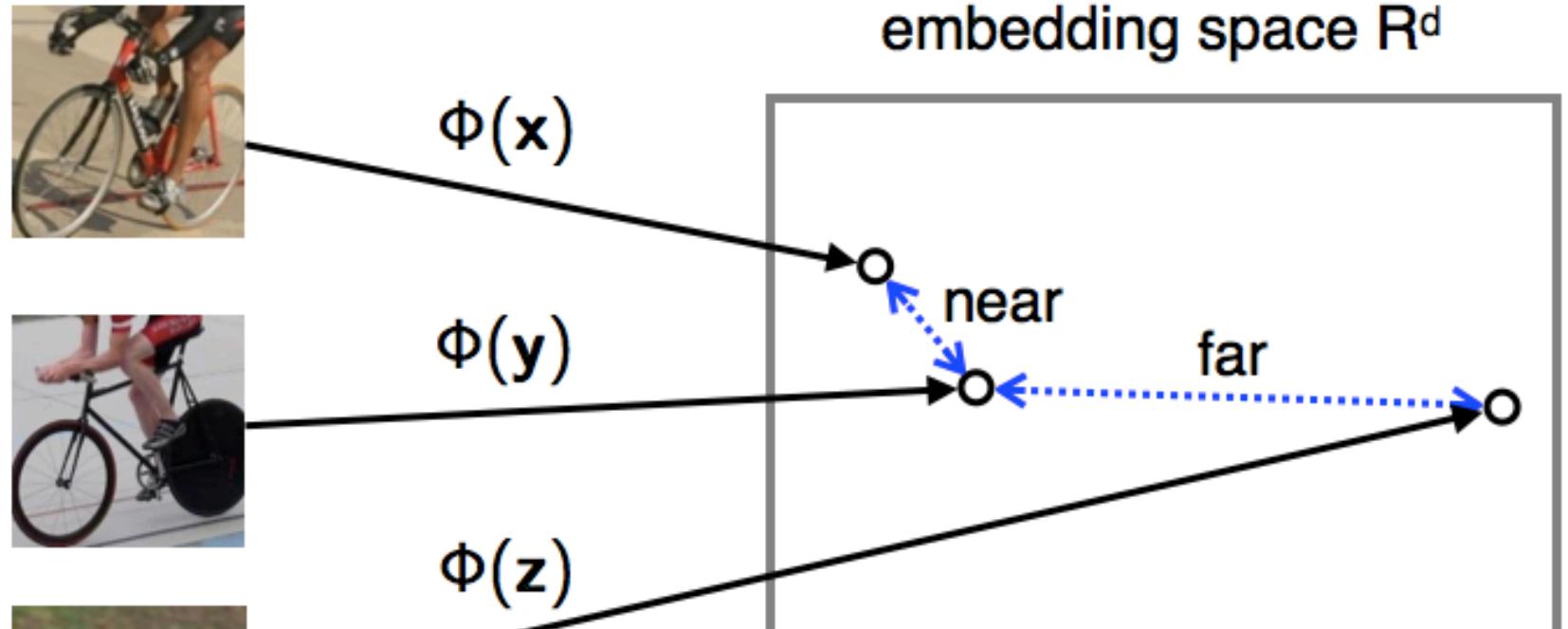
# SVMs in computer vision



# Image features

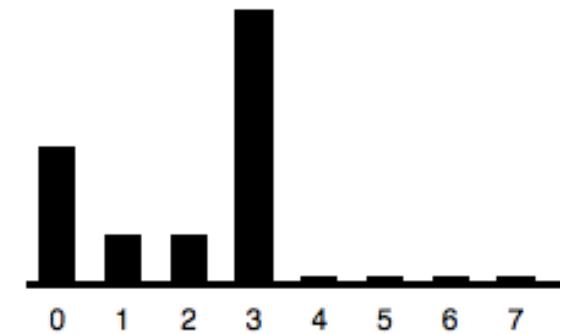
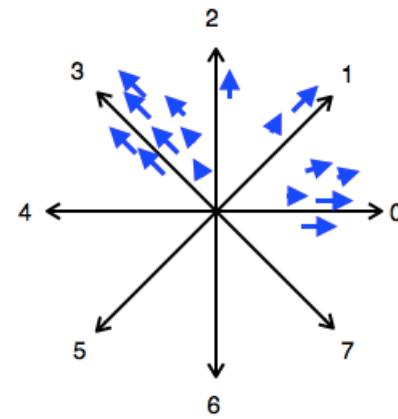
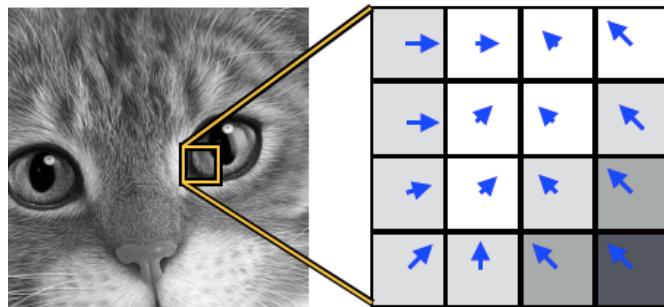
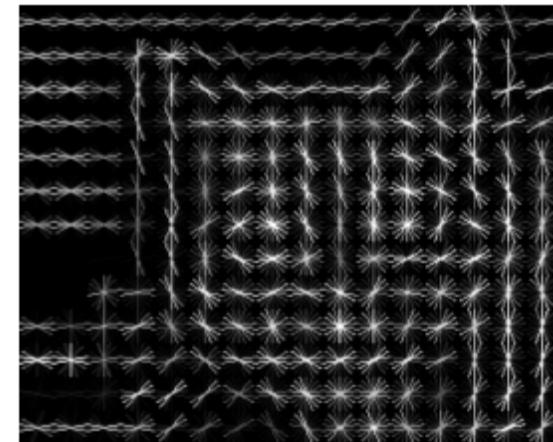


# Desirable feature properties



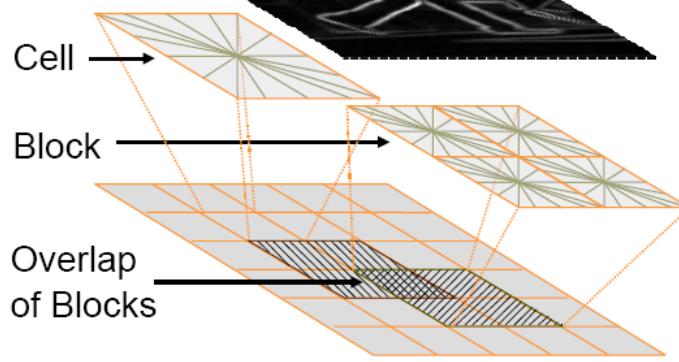
$\Phi$  is **invariant** to nuisance factors, **sensitive** to semantic variations

# Histogram of Gradient (HOG)/SIFT Features



# Dalal and Triggs, ICCV 2005

- Histogram of Oriented Gradient (HOG) features
- Highly accurate detection using linear SVM



Feature vector  $f = [ \dots, \dots, \dots ]$



# HOG features for pedestrians

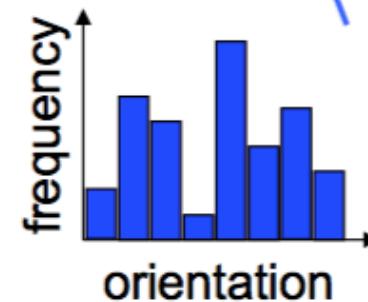
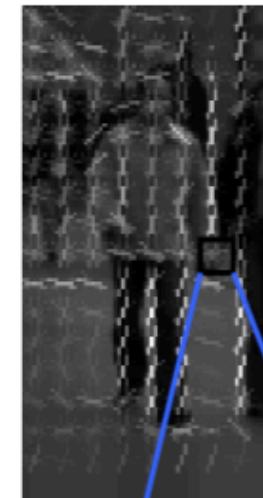
image



dominant direction



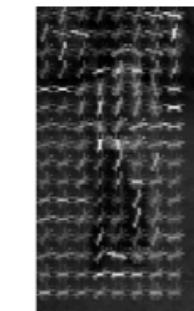
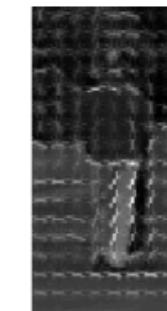
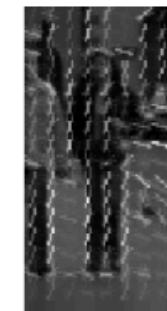
HOG



- tile window into  $8 \times 8$  pixel cells
- each cell represented by HOG

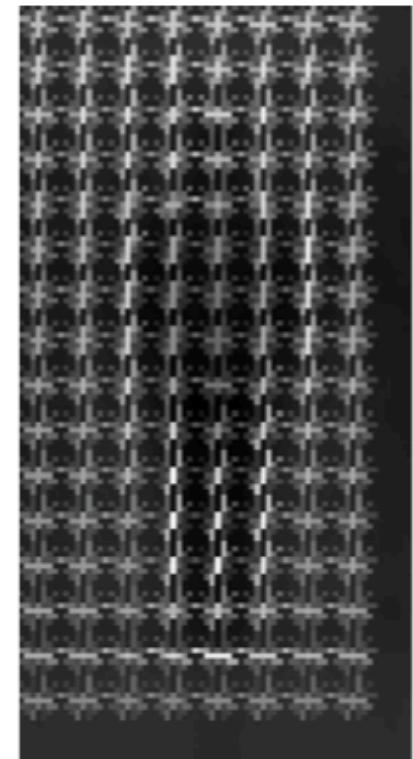
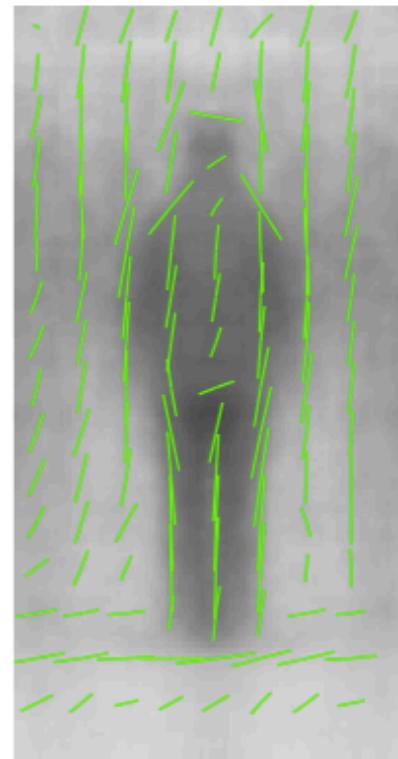
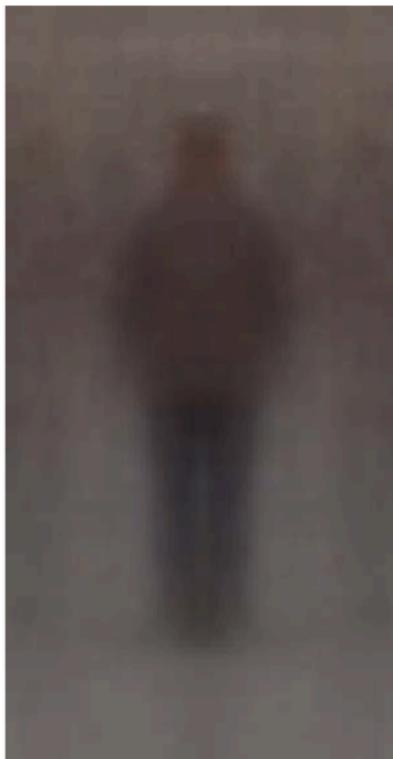
Feature vector dimension =  $16 \times 8$  (for tiling)  $\times 8$  (orientations) = 1024

# SVMs and Pedestrians



# SVMs and Pedestrians

## Averaged examples

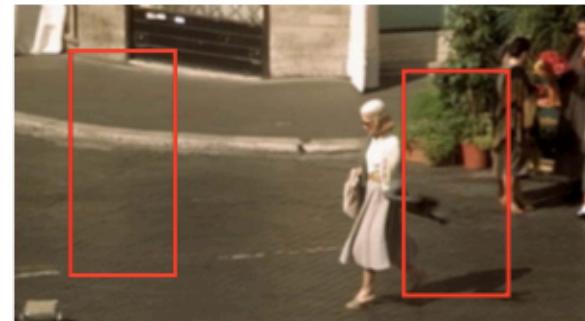


# SVMs and Pedestrians

- Positive data – 1208 positive window examples

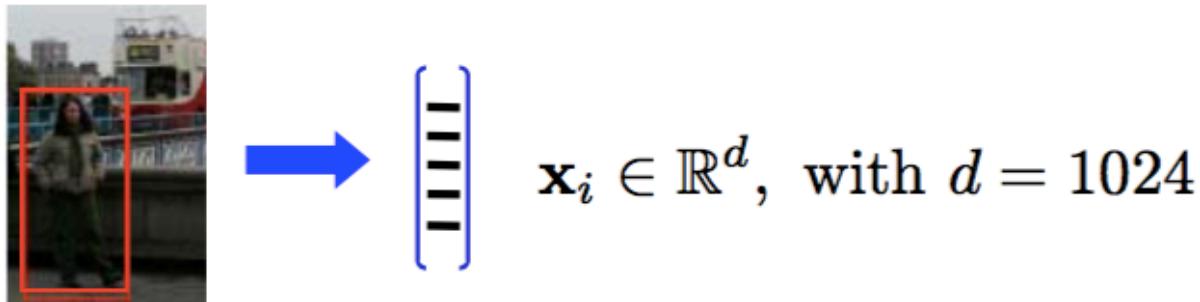


- Negative data – 1218 negative window examples (initially)



## Training (Learning)

- Represent each example window by a HOG feature vector



- Train a SVM classifier

## Testing (Detection)

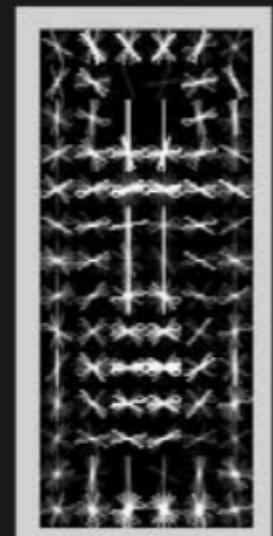
- Sliding window classifier

$$f(x) = \sum_i \alpha_i y_i (\mathbf{x}_i^\top \mathbf{x}) + b$$

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$



positive  
weights

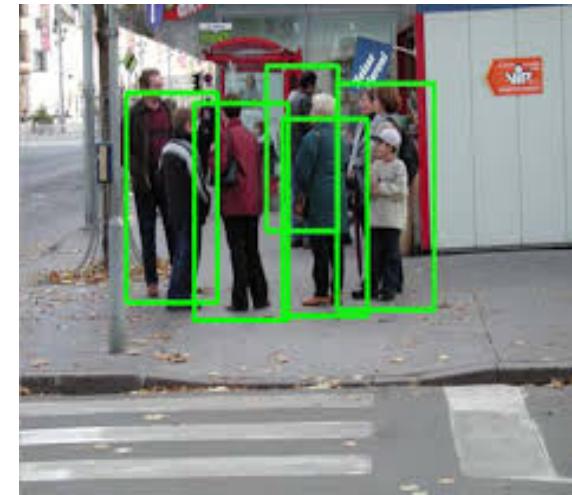
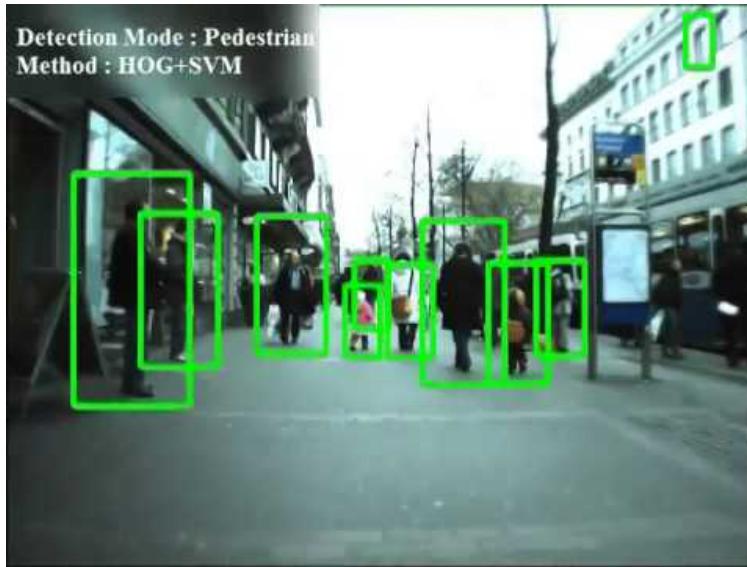
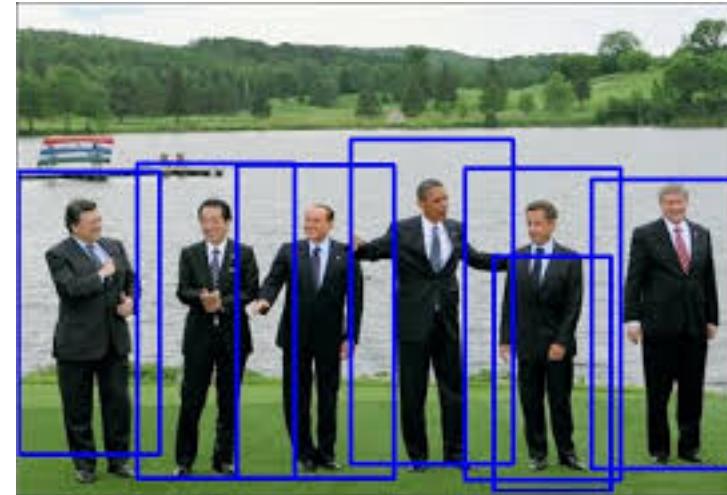
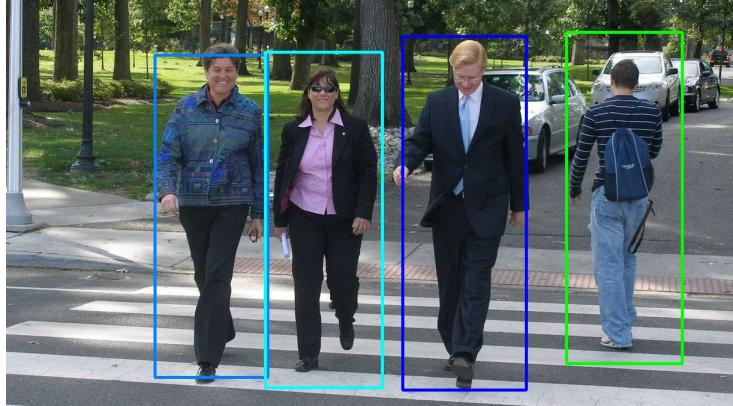


negative  
weights

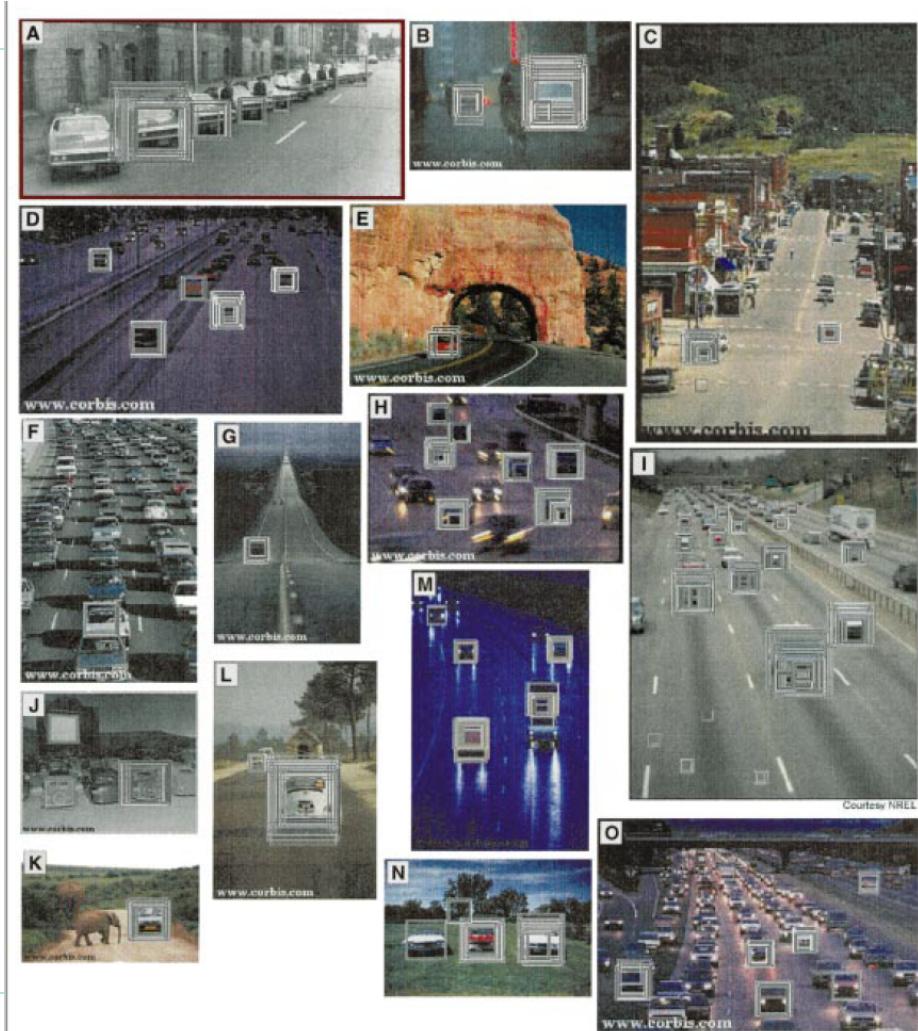
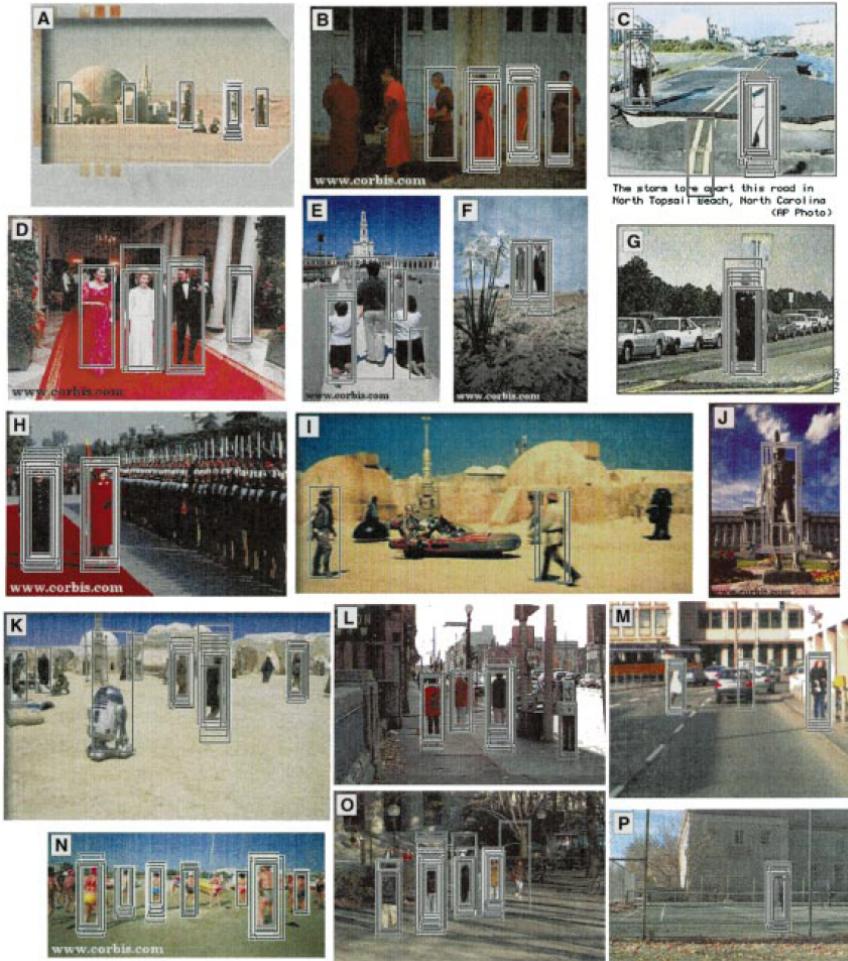


Dalal and Triggs, CVPR 2005

# Pedestrian detection: almost done in 2005



# Papageorgiou & Poggio (1998)





# Lecture outline

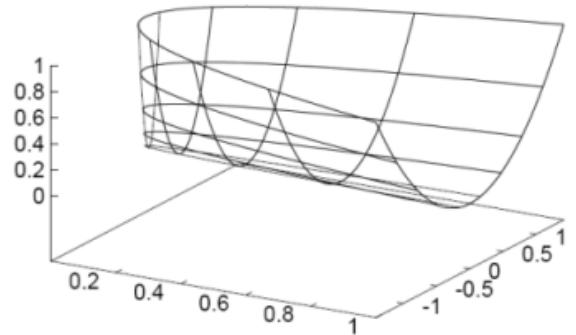
Recap

Large margins and generalization

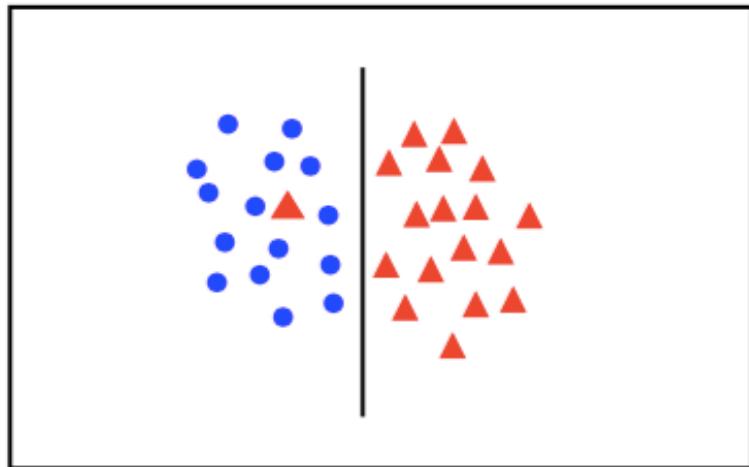
Optimization

Kernels

Applications to vision



# Non-separable data

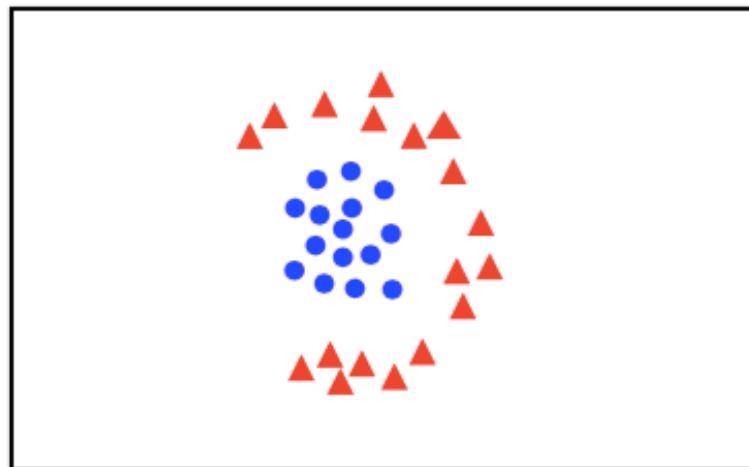


- introduce slack variables

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \|\mathbf{w}\|^2 + C \sum_i^N \xi_i$$

subject to

$$y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \text{ for } i = 1 \dots N$$

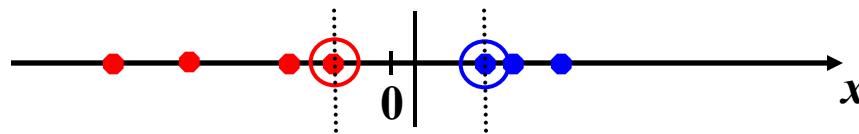


- linear classifier not appropriate

??

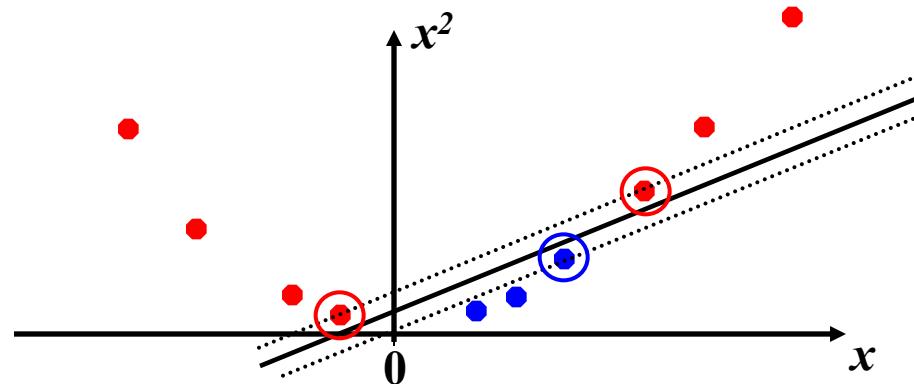
# Non-linear SVMs

- Datasets that are linearly separable (with some noise) work out great:



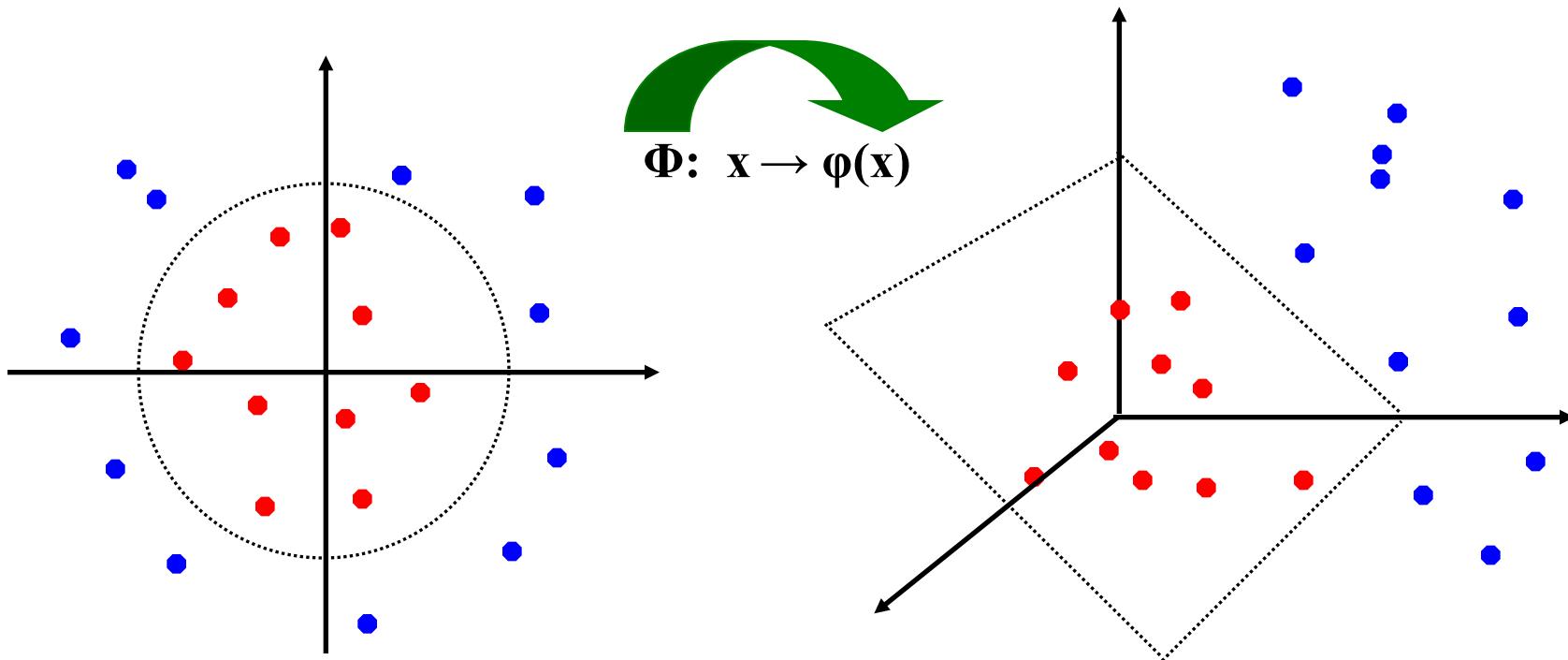
- But what are we going to do if the dataset is just too hard?

- How about ... mapping data to a higher-dimensional space:

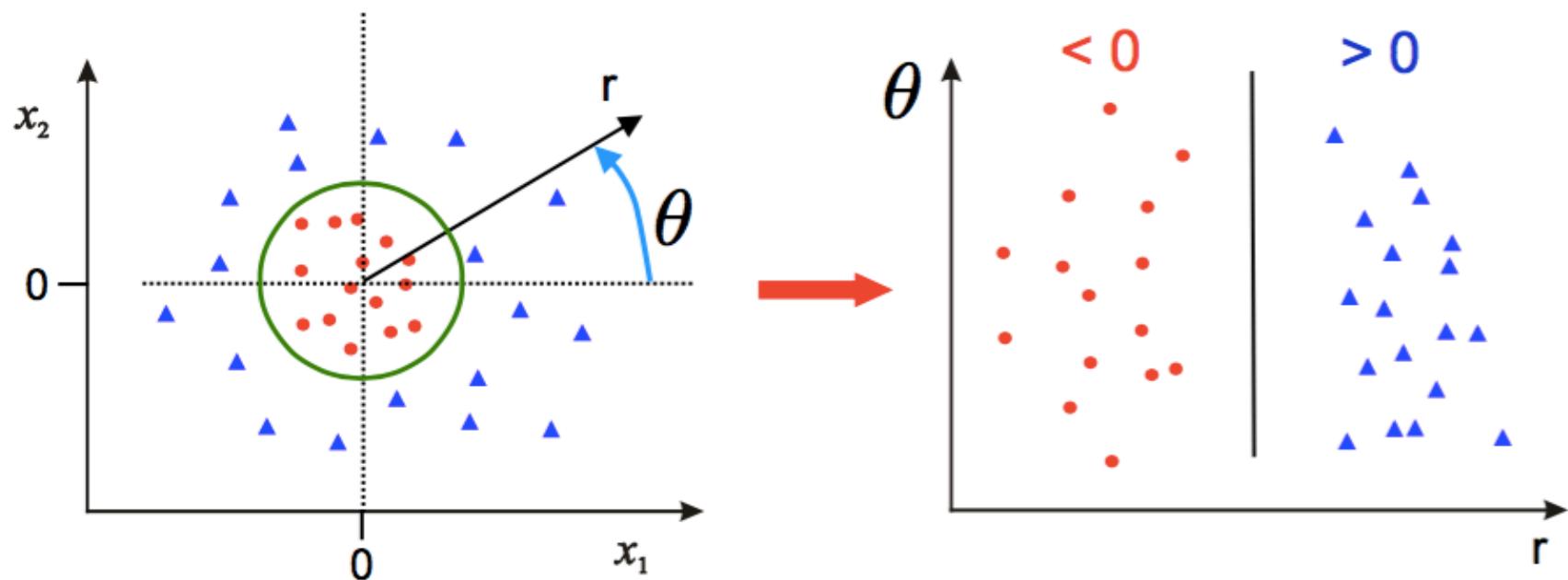


# Non-linear SVMs: Feature spaces

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



## Solution by inspection: hand-crafted features



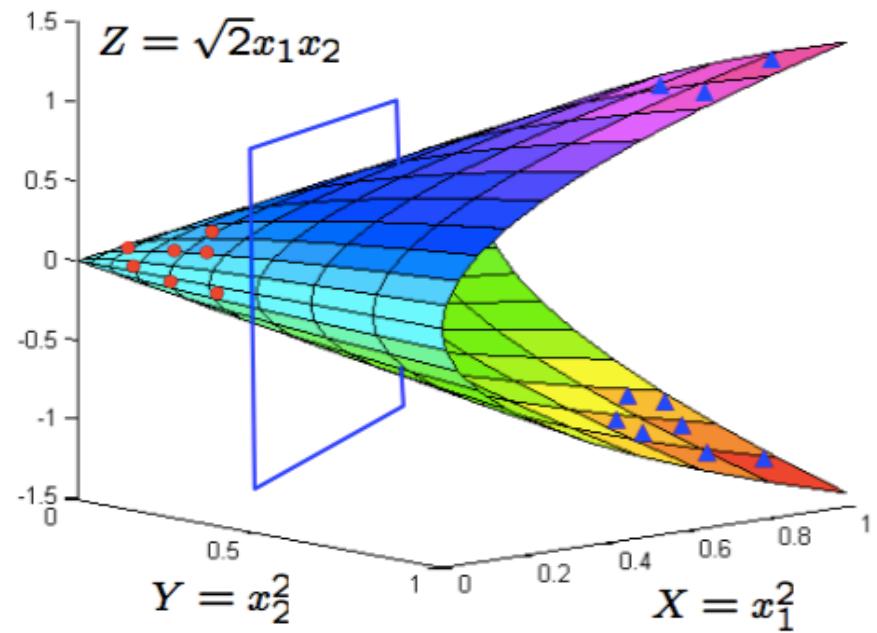
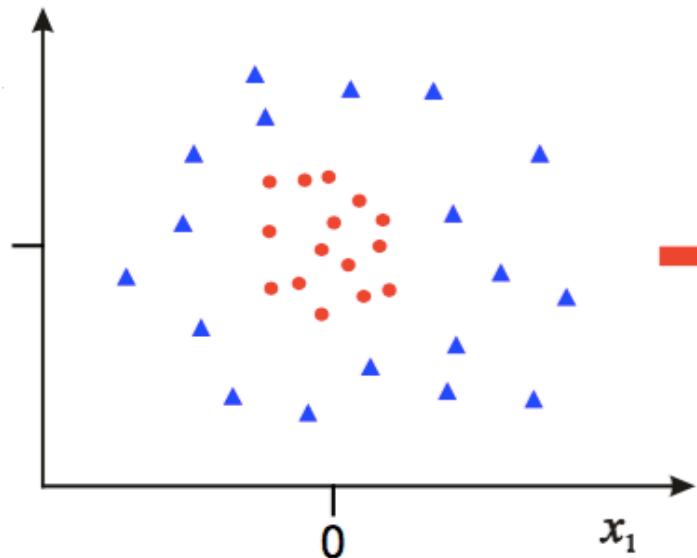
- Data **is** linearly separable in polar coordinates
- Acts non-linearly in original space

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} r \\ \theta \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

## More general method

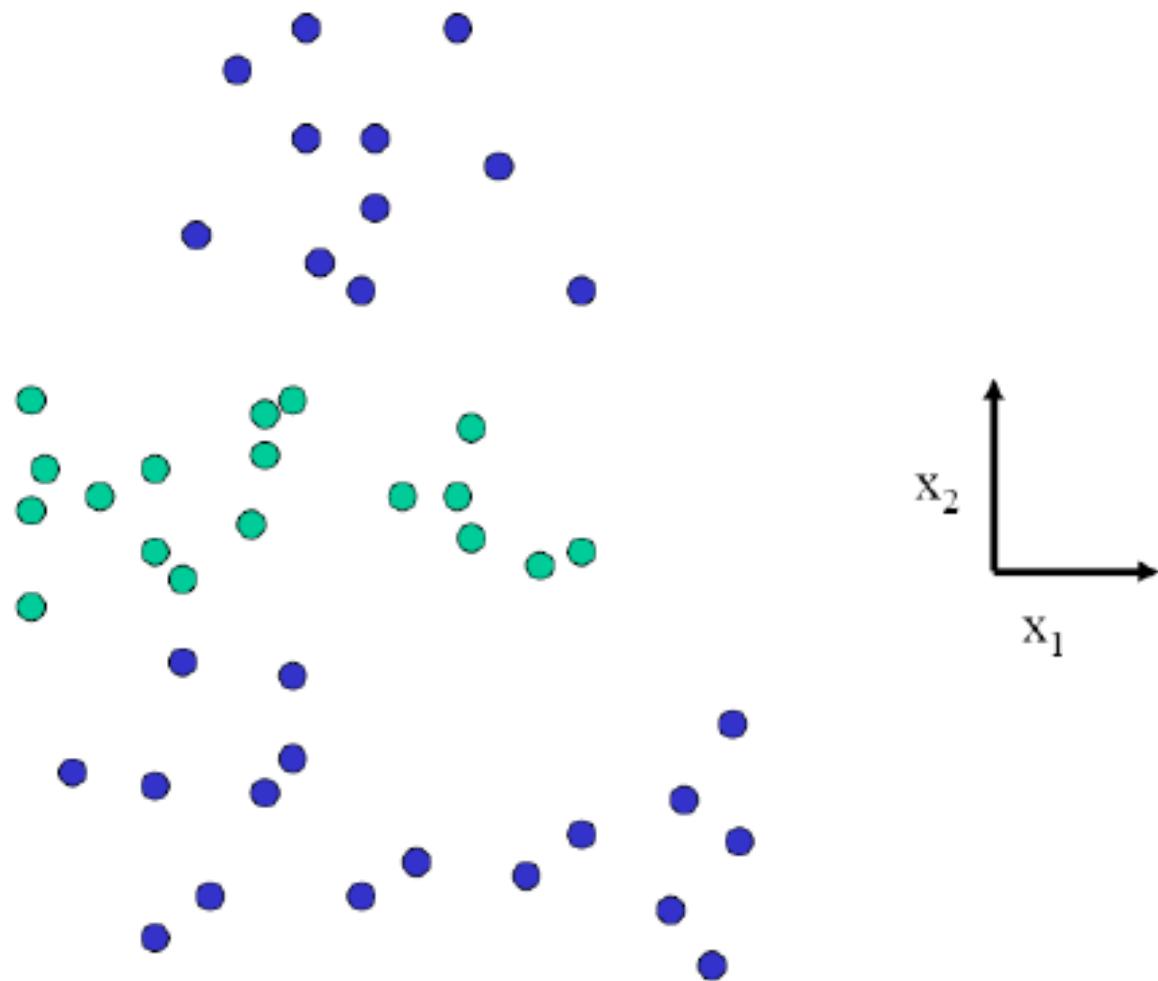
---

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

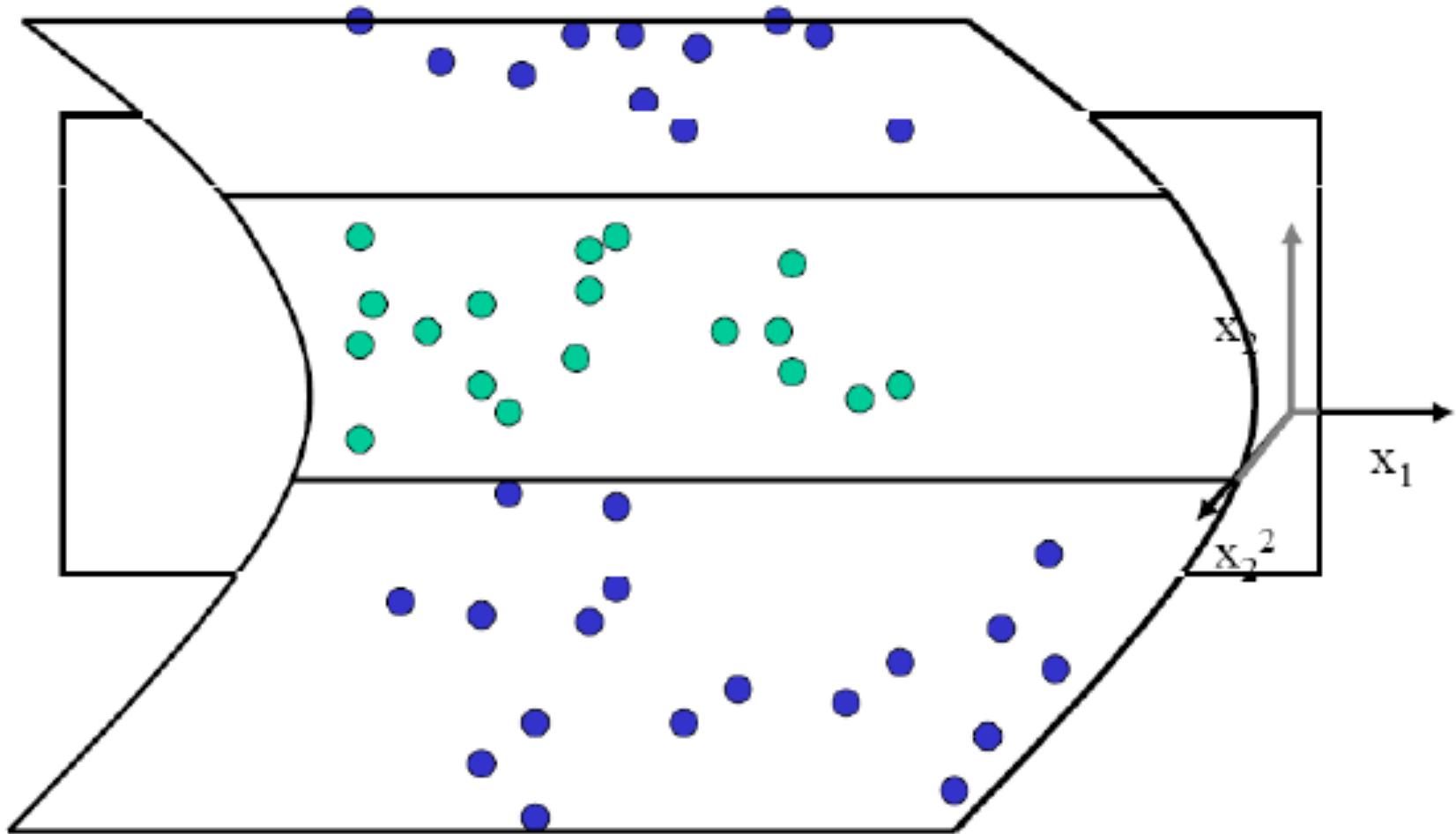


- Data is linearly separable in 3D
- This means that the problem can still be solved by a linear classifier

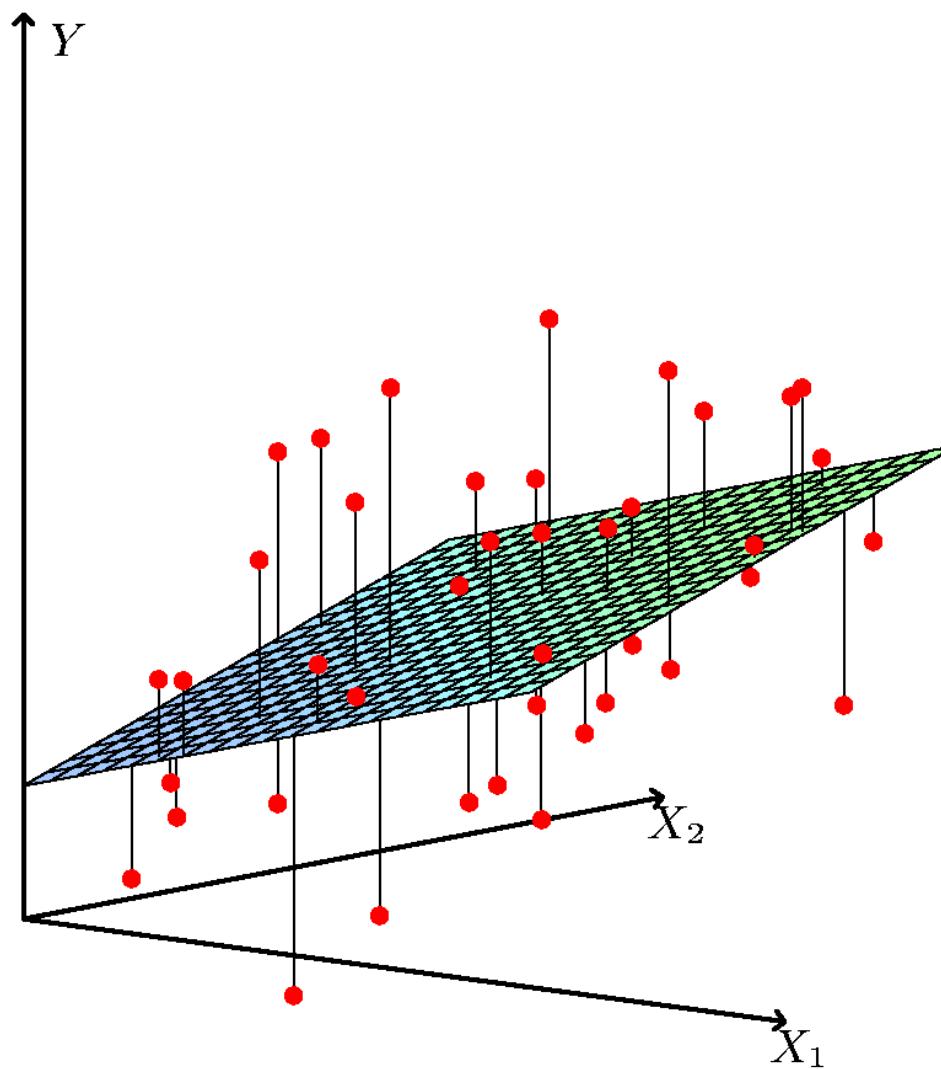
# Nonseparable in 2D



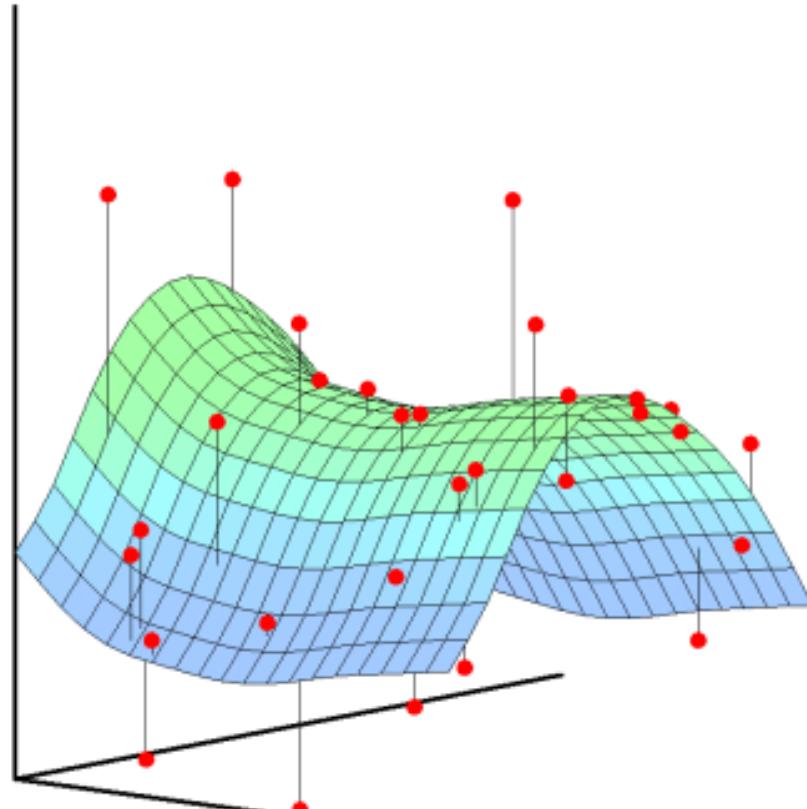
# Separable in 3D



# Linear regression



# Nonlinear regression



$$\mathbf{x} \rightarrow \phi(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \vdots \\ \phi_M(\mathbf{x}) \end{bmatrix}$$

## Example: second-order polynomials

$$\mathbf{x} = (x_1, x_2)$$

$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ (x_1)^2 \\ (x_2)^2 \\ x_1 x_2 \end{bmatrix}$$

$$\langle \mathbf{w}, \phi(\mathbf{x}) \rangle = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$$

# Non-linear Classifiers

So far, decision is based on the sign of  $y = \mathbf{w}^T \mathbf{x}$

Use non-linear transformation,  $\phi(\mathbf{x})$  of our data,  $\mathbf{x}$

$$\text{e.g. } \mathbf{x} = (x_1, x_2) \quad \phi(\mathbf{x}) =$$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ (x_1)^2 \\ (x_2)^2 \\ x_1 x_2 \end{bmatrix}$$

Discriminant:

$$\langle \mathbf{w}, \phi(\mathbf{x}) \rangle = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$$

Non-linear in  $\mathbf{x}$ , linear in  $\phi(\mathbf{x})$

# Dual form of SVM & kernel trick

Optimization:

$$\min_{\alpha} \sum_{i=1}^N \sum_{j=1}^N \alpha^i \alpha^j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

s.t. :  $y^i \left( \sum_{j=1}^N \alpha^j y^j \langle \mathbf{x}^j, \mathbf{x}^i \rangle + b \right) \geq 1, \quad \forall i$

$$\boldsymbol{\alpha} \in \mathbb{R}^N \rightarrow O(N^3)$$

Primal and dual classifier forms:

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i=1}^N \alpha^i y^i \langle \mathbf{x}^i, \mathbf{x} \rangle + b$$

What if we replace  $\mathbf{x}$  with  $\phi(\mathbf{x})$ ?

Everything involves only inner products!

Rewrite everything in terms of Kernel

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

## Dual form of SVM & kernel trick

Optimization:

$$\min_{\alpha} \sum_{i=1}^N \sum_{j=1}^N \alpha^i \alpha^j y^i y^j K(\mathbf{x}^i, \mathbf{x}^j)$$

s.t. :  $y^i \left( \sum_{j=1}^N \alpha^j y^j K(\mathbf{x}^j, \mathbf{x}^i) + b \right) \geq 1, \quad i = 1, \dots, N$

Dual classifier form:

$$\begin{aligned} f(\mathbf{x}) &= \sum_{i=1}^N \alpha^i y^i K(\mathbf{x}^i, \mathbf{x}) + b \\ &= \sum_{\{i: \alpha^i \neq 0\}} w^i K(\mathbf{x}^i, \mathbf{x}) + b, \quad w^i = y^i \alpha^i \end{aligned}$$

Compare with general nonlinear form:

$$f(\mathbf{x}) = \sum_k w_k \phi_k(\mathbf{x})$$

N nonlinear functions – smart choice of sparse coefficients

## 'Kernel trick'

Consider:  $\phi(\mathbf{x}) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ 1 \end{bmatrix}$

We then have:  $\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle =$

$$\begin{aligned} &= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2 + 2x_1 y_1 + 2x_2 y_2 + 1 \\ &= (x_1 y_1 + x_2 y_2 + 1)^2 \\ &= (\mathbf{x}^T \mathbf{y} + 1)^2 \doteq K(\mathbf{x}, \mathbf{y}) \end{aligned}$$

Polynomial Kernel  $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^p$

Kernel: linear complexity in D (dimensions of  $\mathbf{x}, \mathbf{y}$ ), constant in p

Feature space complexity: much higher

## Condition for kernel trick: ‘Mercer’ kernel

- Given some arbitrary function  $k(\mathbf{x}_i, \mathbf{x}_j)$ , how do we know if it corresponds to a scalar product  $\Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$  in some space?
  - Mercer kernels: if  $k(, )$  satisfies:
    - Symmetric  $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$
    - Positive definite,  $\boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \geq 0$  for all  $\boldsymbol{\alpha} \in \mathbb{R}^N$ , where  $\mathbf{K}$  is the  $N \times N$  Gram matrix with entries  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ .
- then  $k(, )$  is a valid kernel.

# Mercer Kernel Examples

Linear kernel

$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$$

Polynomial kernel

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^p$$

Radial Basis Function (a.k.a. Gaussian) kernel

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right)$$

Underlying feature dimension: **Infinite**

# RBF kernel SVM

$N = \text{size of training data}$

$$f(\mathbf{x}) = \sum_i^N \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$



weight (may be zero)      support vector

$$\text{Gaussian kernel } k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma^2)$$

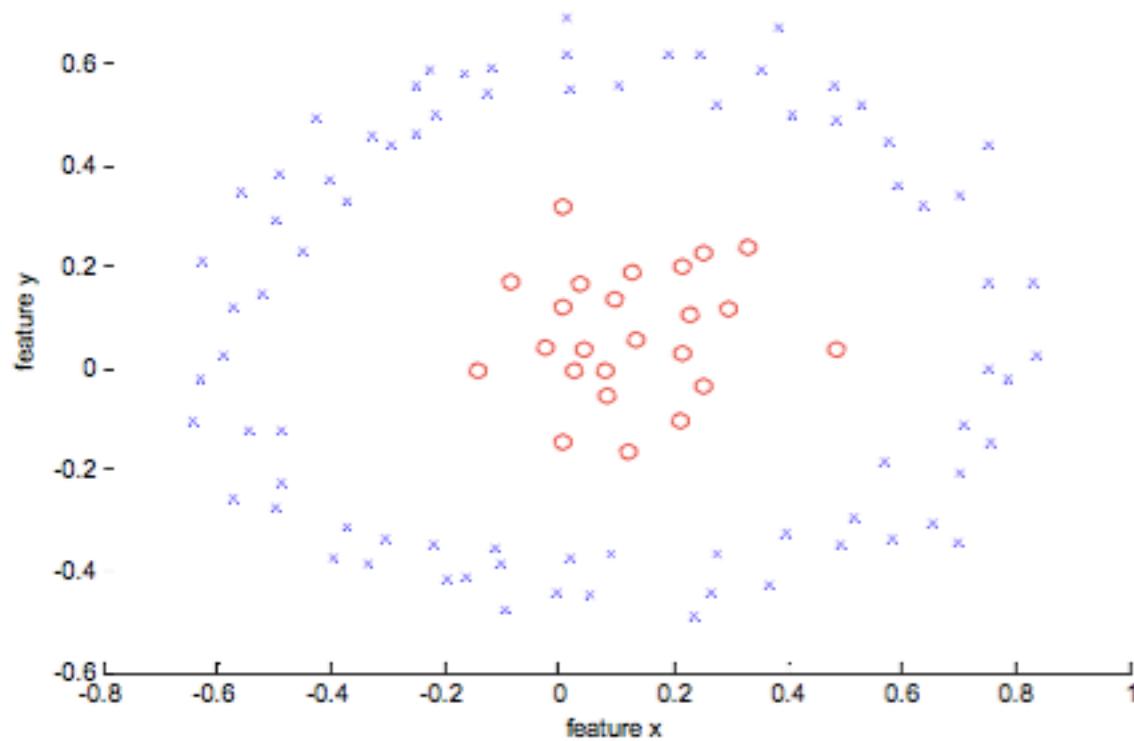
Radial Basis Function (RBF) SVM

# RBF kernel SVM (next week's assignment)

$$\begin{aligned}
 f(\mathbf{x}) &= \sum_{i=1}^N \alpha^i y^i K(\mathbf{x}^i, \mathbf{x}) + b \\
 &= \sum_{\{i: \alpha^i \neq 0\}} \alpha^i y^i K(\mathbf{x}^i, \mathbf{x}) + b \\
 &= \sum_{\{i: \alpha^i \neq 0\}} w^i \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mathbf{x}\|_2^2\right) + b
 \end{aligned}$$

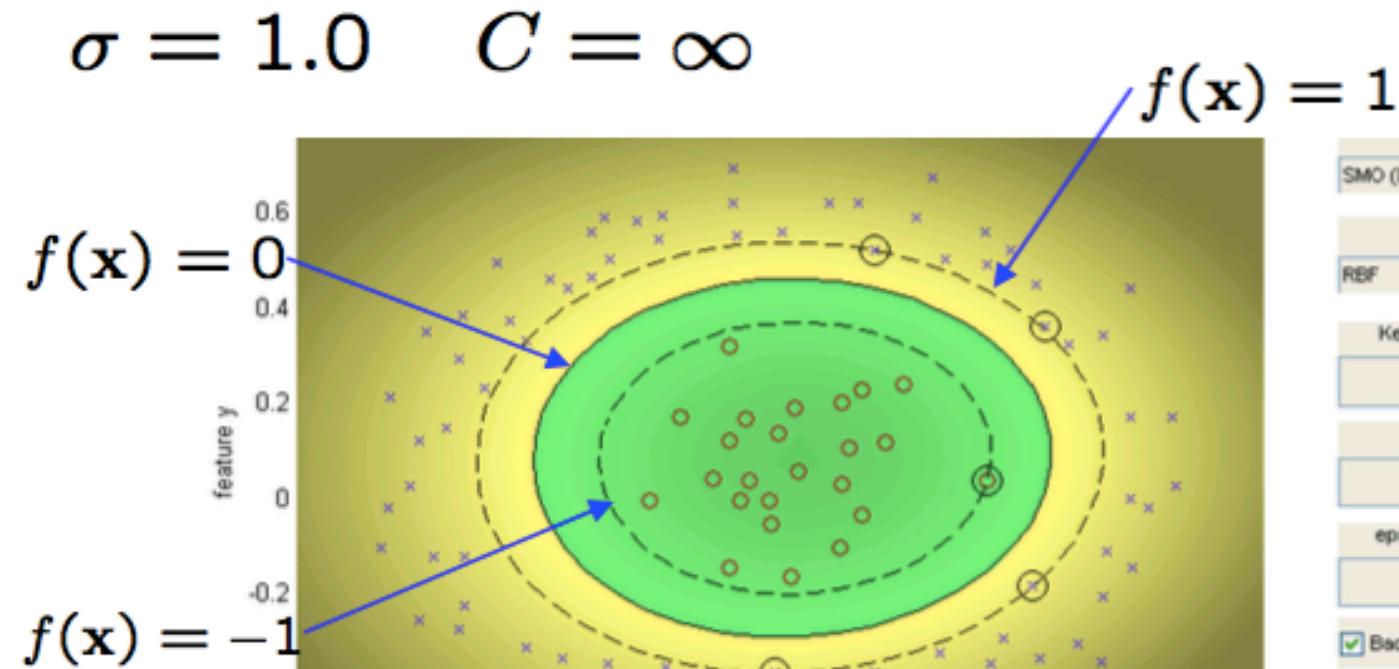
Discriminant form: sum of bumps centered on training points

# RBF-SVM example



- data is not linearly separable in original feature space

# RBF-SVM example



SMO (L1)

Kernel

RBF

Kernel argument

1

C-constant

Inf

epsilon,tolerance

1e-3,1e-3

Background

Load data

Create data

Reset

Train SVM

Info

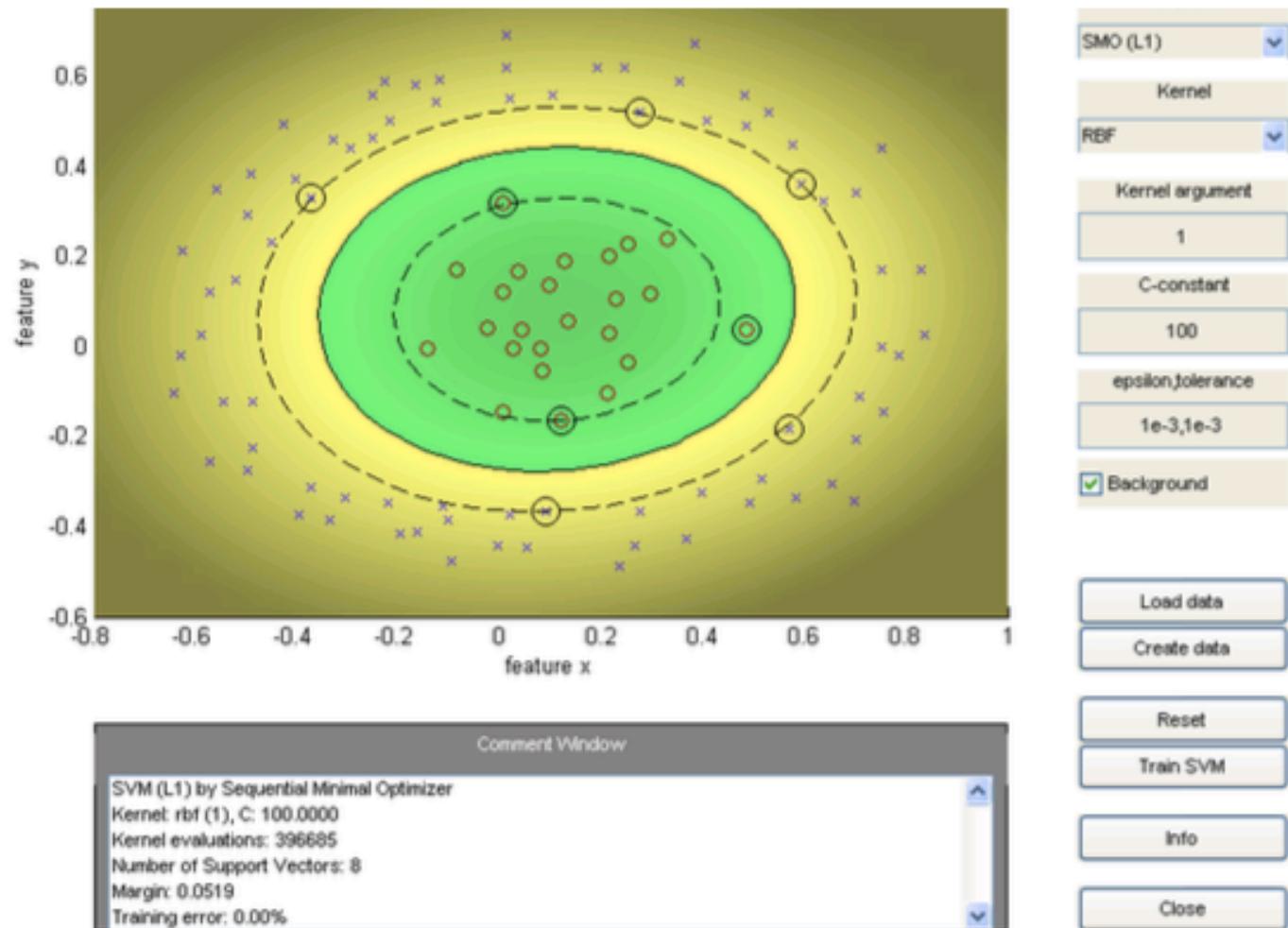
  

Close



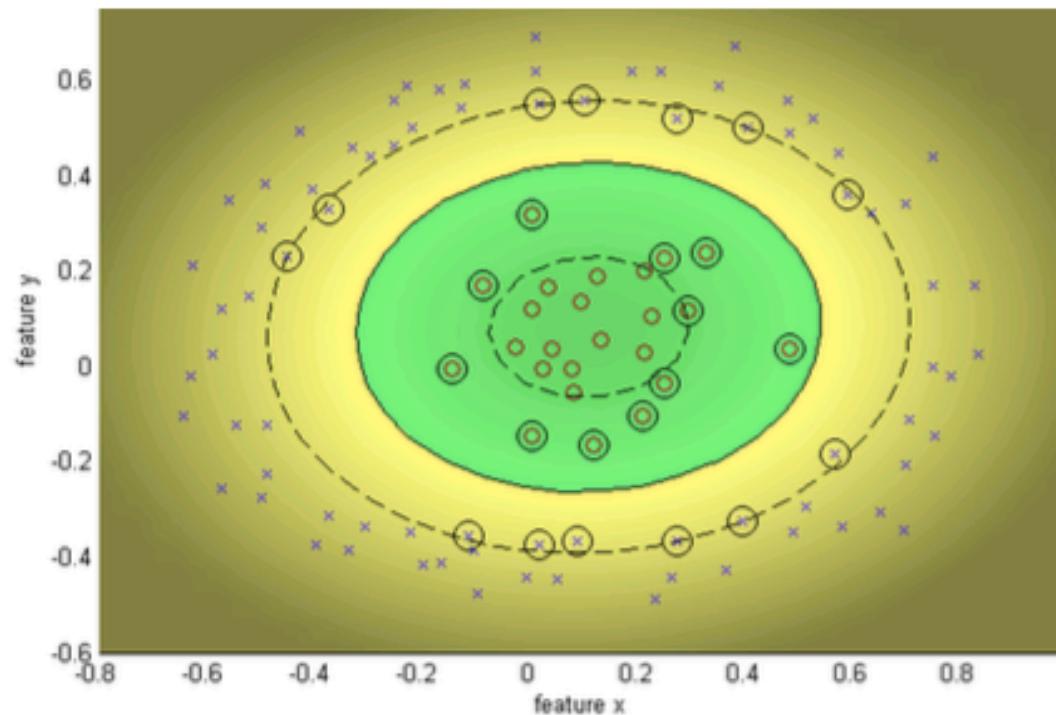
# RBF-SVM example

$$\sigma = 1.0 \quad C = 100$$



# RBF-SVM example

$$\sigma = 1.0 \quad C = 10$$



Comment Window

SVM (L1) by Sequential Minimal Optimizer  
Kernel: rbf (1), C: 10.0000  
Kernel evaluations: 46158  
Number of Support Vectors: 24  
Margin: 0.0755  
Training error: 0.00%

SMO (L1)

Kernel

RBF

Kernel argument

1

C-constant

10

epsilon/tolerance

1e-3,1e-3

Background

Load data

Create data

Reset

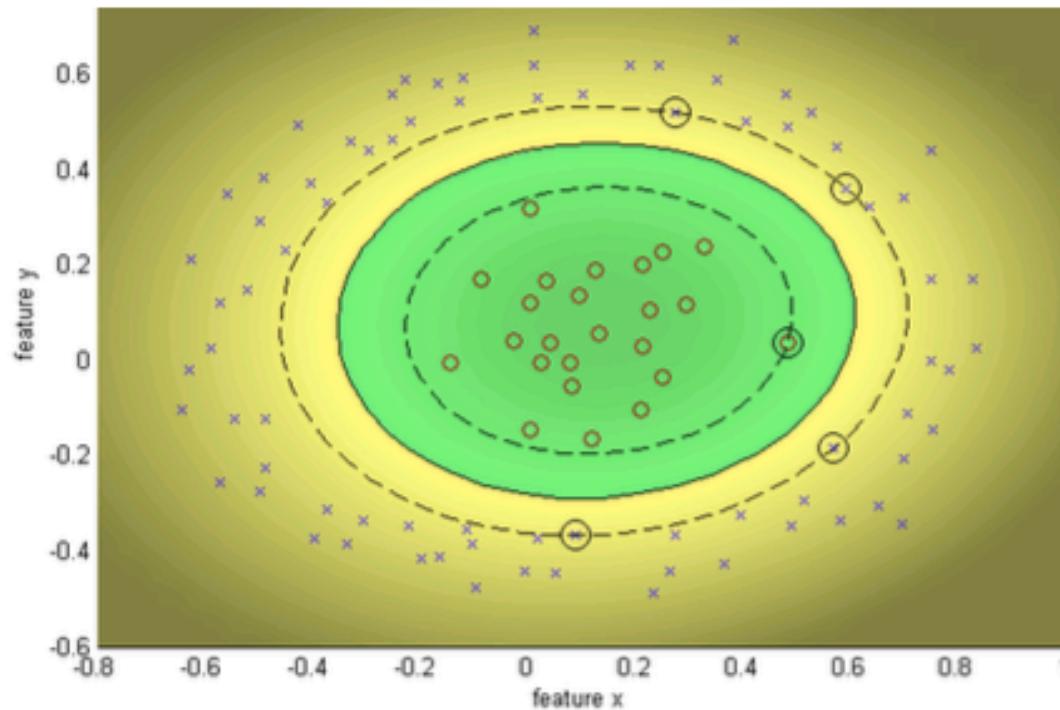
Train SVM

Info

Close

# RBF-SVM example

$$\sigma = 1.0 \quad C = \infty$$



SMO (L1) ▾  
Kernel  
RBF ▾  
Kernel argument  
1  
C-constant  
Inf  
epsilon/tolerance  
1e-3,1e-3  
Background

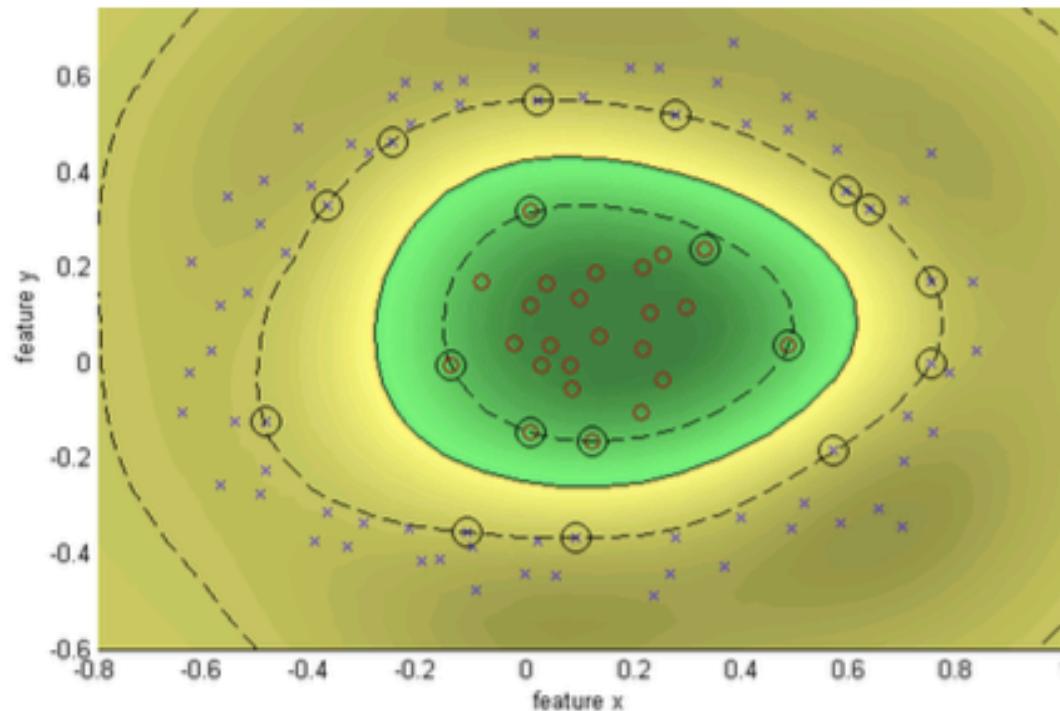
Load data  
Create data  
Reset  
Train SVM  
Info  
Close

Comment Window

SVM (L1) by Sequential Minimal Optimizer  
Kernel: rbf (1), C: Inf  
Kernel evaluations: 62739  
Number of Support Vectors: 5  
Margin: 0.0445  
Training error: 0.00%

# RBF-SVM example

$$\sigma = 0.25 \quad C = \infty$$



SMO (L1)	<input type="button" value="▼"/>
Kernel	
RBF	<input type="button" value="▼"/>
Kernel argument	0.25
C-constant	Inf
epsilon,tolerance	1e-3,1e-3
<input checked="" type="checkbox"/> Background	

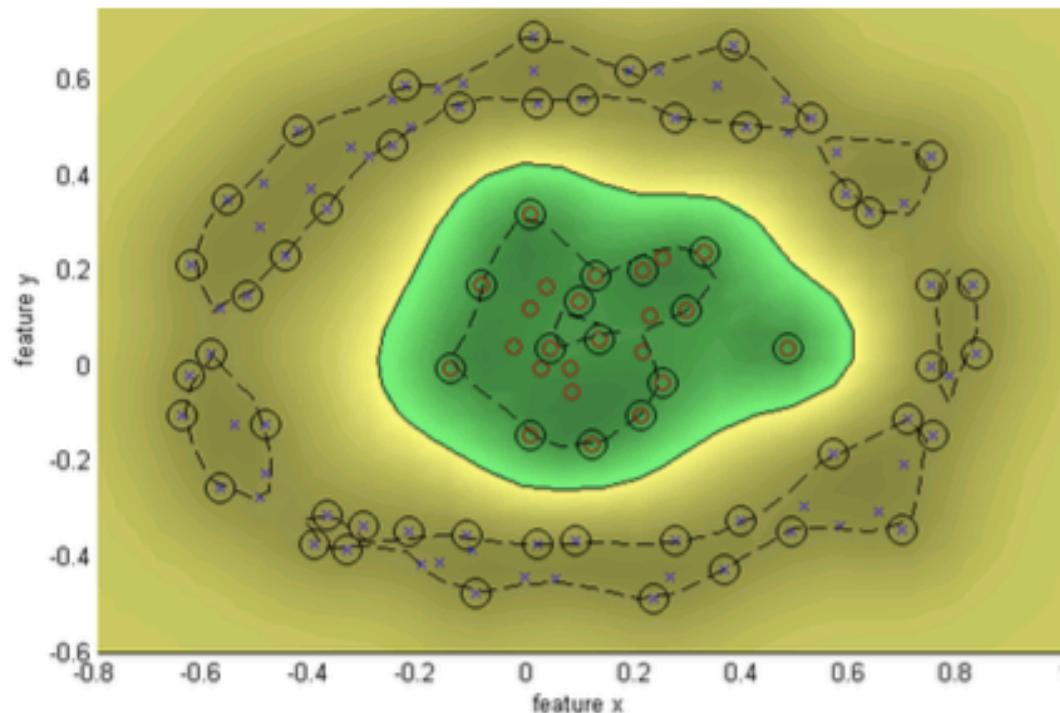
<input type="button" value="Load data"/>
<input type="button" value="Create data"/>
<input type="button" value="Reset"/>
<input type="button" value="Train SVM"/>
<input type="button" value="Info"/>
<input type="button" value="Close"/>

Comment Window

SVM (L1) by Sequential Minimal Optimizer  
 Kernel: rbf (0.25), C: Inf  
 Kernel evaluations: 42795  
 Number of Support Vectors: 18  
 Margin: 0.2358  
 Training error: 0.00%

# RBF-SVM example

$$\sigma = 0.1 \quad C = \infty$$

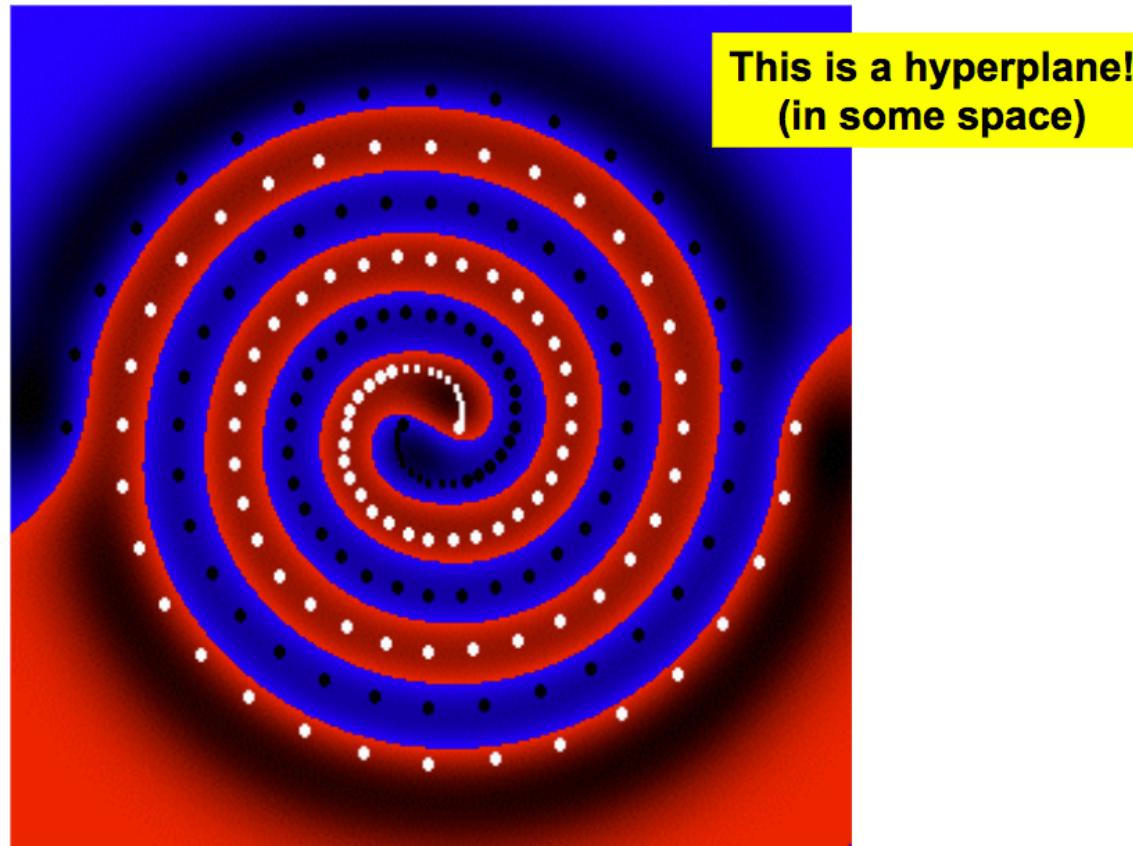


SMO (L1)	<input type="button" value="▼"/>
Kernel	
RBF	<input type="button" value="▼"/>
Kernel argument	<input type="text" value="0.1"/>
C-constant	<input type="text" value="Inf"/>
epsilon,tolerance	<input type="text" value="1e-3,1e-3"/>
<input checked="" type="checkbox"/> Background	
<input type="button" value="Load data"/>	
<input type="button" value="Create data"/>	
<input type="button" value="Reset"/>	
<input type="button" value="Train SVM"/>	
<input type="button" value="Info"/>	
<input type="button" value="Close"/>	

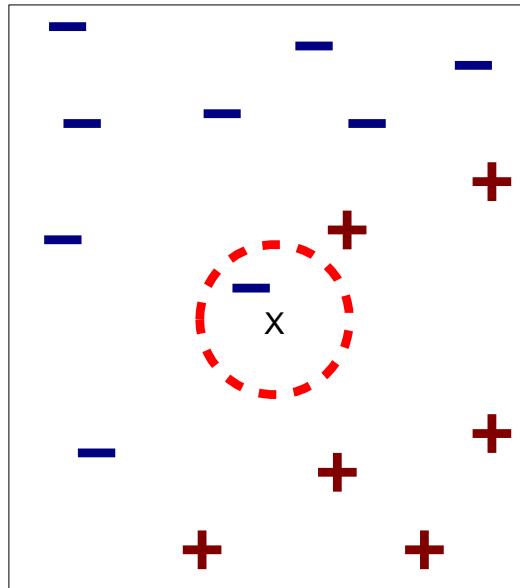
Comment Window

```
SVM (L1) by Sequential Minimal Optimizer
Kernel: rbf (0.1), C: Inf
Kernel evaluations: 173935
Number of Support Vectors: 62
Margin: 0.2196
Training error: 0.00%
```

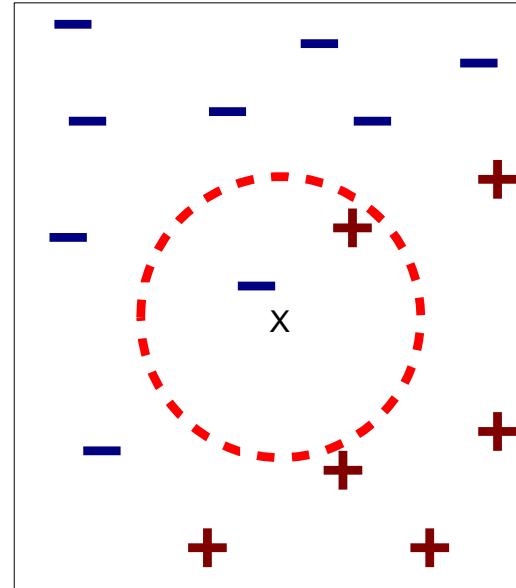
# All of the flexibility you may need is there



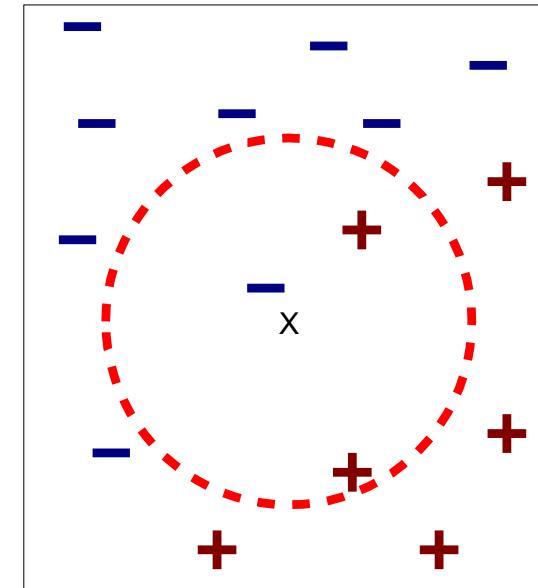
# Reminder: K-nearest neighbor classifier



(a) 1-nearest neighbor



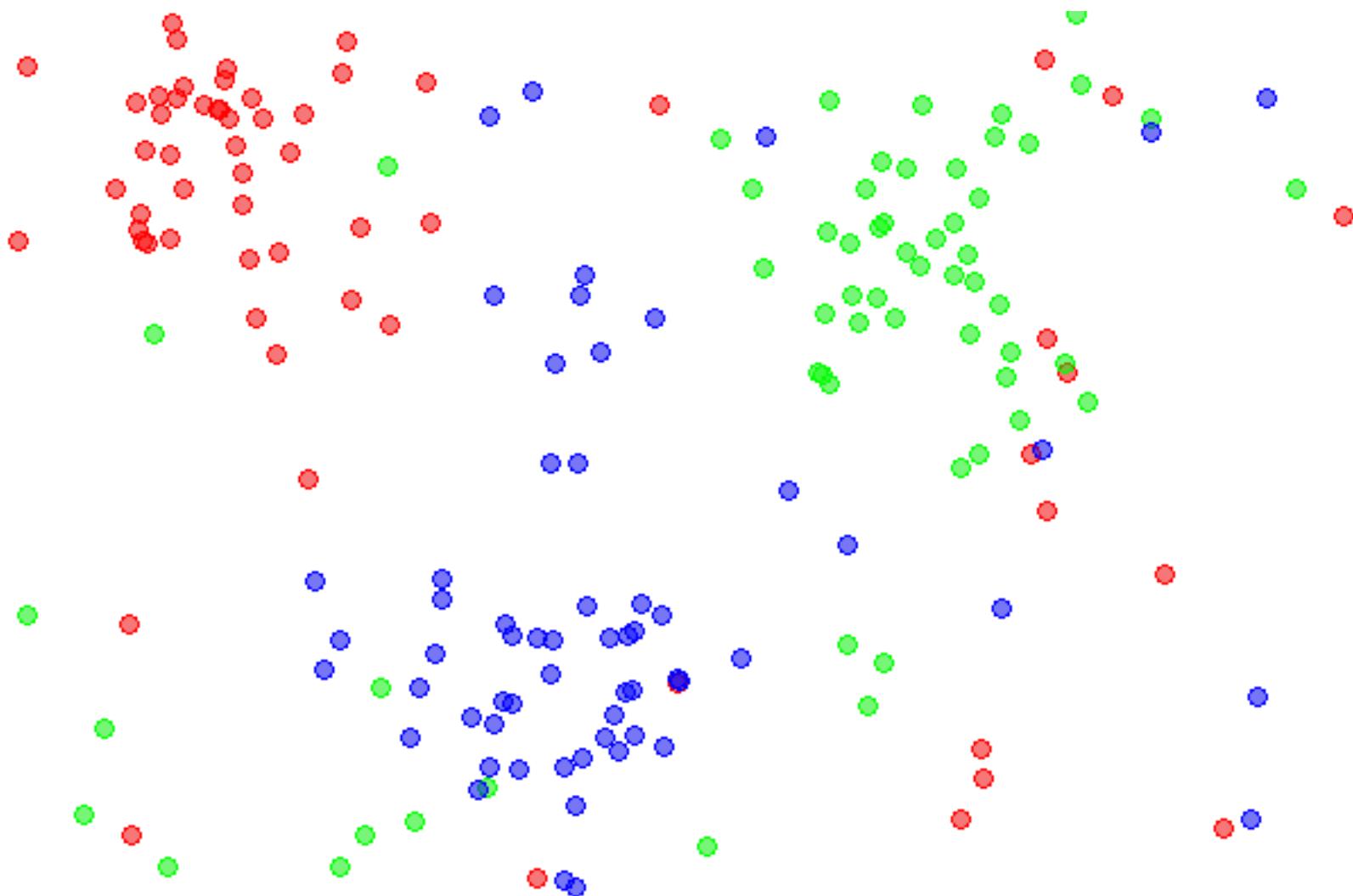
(b) 2-nearest neighbor



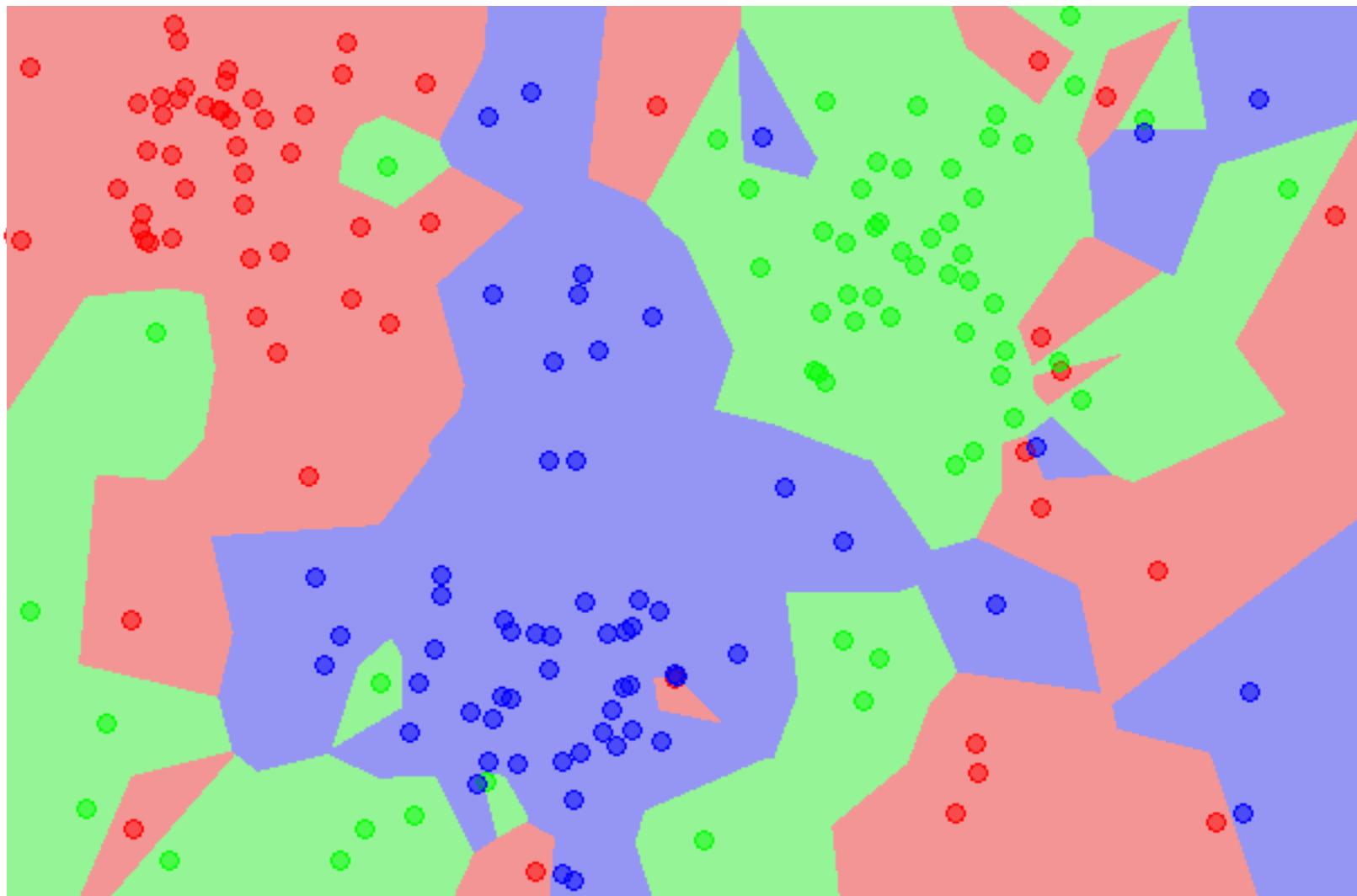
(c) 3-nearest neighbor

- Compute distance to other training records
- Identify  $K$  nearest neighbors
- Take majority vote

# Training data for NN classifier (in $\mathbb{R}^2$ )



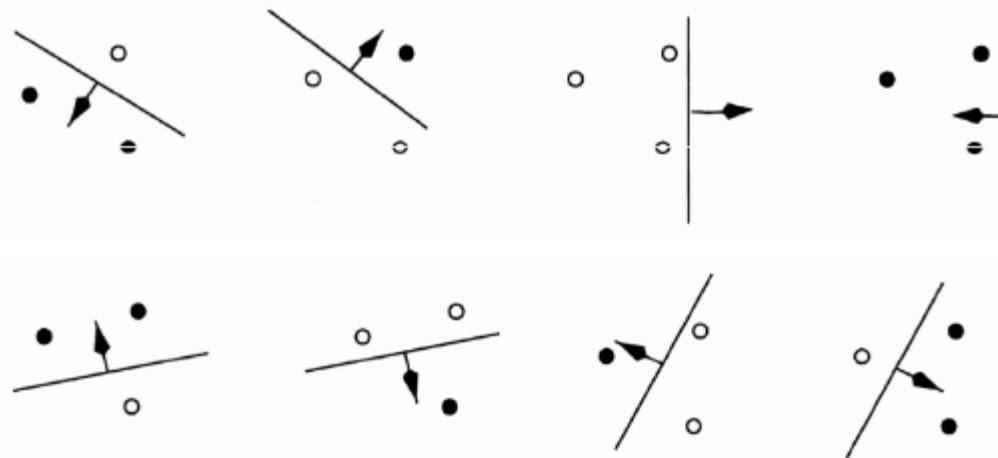
# 1-nn classifier prediction (in $R^2$ )



# VC dimension of 1-nearest neighbor classifier?

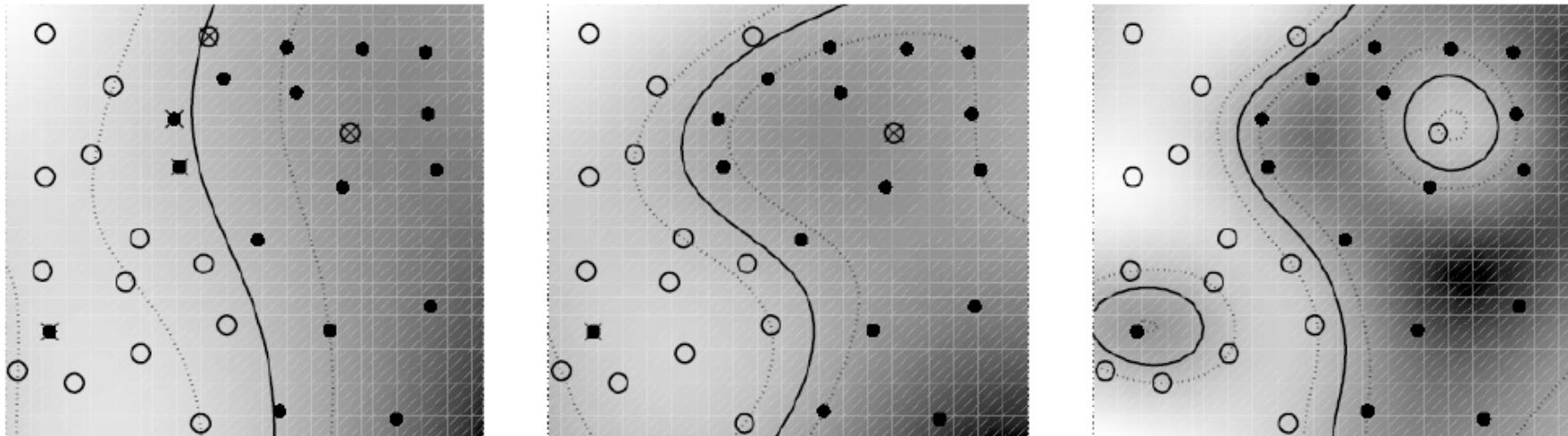
The VC dimension for the set of functions  $\{f(\alpha)\}$  is defined as the maximum number of training points that can be shattered by  $\{f(\alpha)\}$ .

- VC dimension of N-dimensional linear classifier: N+1



- VC dimension of 1-NN: infinite

# Large margins for nonlinear classifiers



RBF Kernel width ( $\sigma$ )

**Margin size:** determined by both  $\sigma$  and regularizer

We can slide between a linear and a Nearest-Neighbor classifier!

# Guyon & Vapnik, 1995

- Handwritten digit recognition
  - US Postal Service Database
  - Standard benchmark task for many learning algorithms

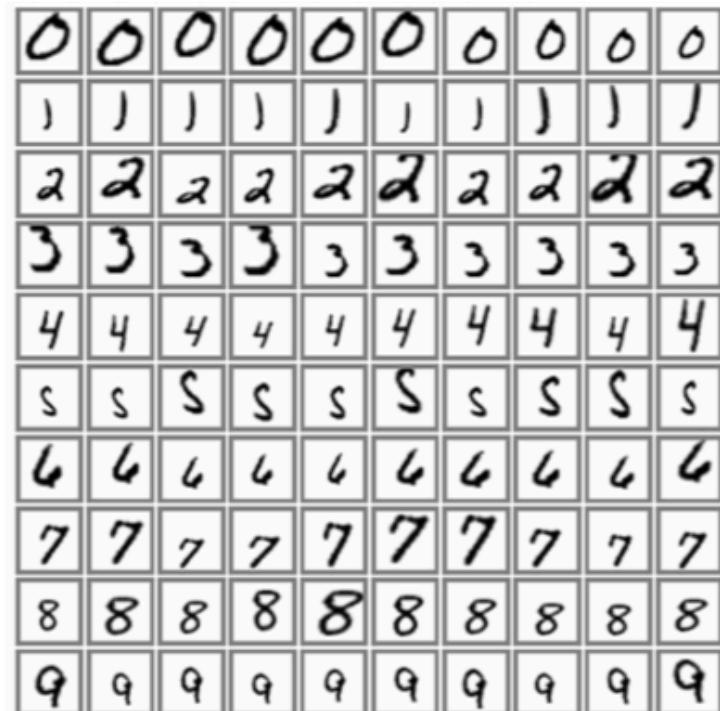
A 20x20 grid of handwritten digits, likely from the USPS dataset. The digits are rendered in black on a white background. They vary in size and orientation, representing a diverse set of handwritten samples. The digits include numbers from 0 to 9, with some appearing more frequently than others.

2 6 0 1 4 4 6 7 5 7 1 4 6 3 7 1 0 3 7 1 1 4 4 9 7  
1 1 0 5 7 1 1 1 6 9 9 8 1 1 9 2 1 6 0 0 2 8 8 2 0  
3 3 0 1 0 3 3 0 1 0 2 7 9 6 0 2 8 1 0 0 2 9 0 1 2  
1 4 0 5 2 8 0 6 7 2 9 5 0 1 3 1 5 3 0 2 9 9 5 5  
5 1 0 1 2 9 2 0 1 3 0 3 2 7 7 0 1 3 9 4 3 4 8 6 4  
1 1 4 1 1 7 6 0 5 7 1 8 8 4 0 0 1 5 8 7 0 1 8 2 2  
1 1 5 7 5 5 7 2 1 2 5 7 0 6 8 8 2 2 1 4 9 9 8 1 6  
9 9 5 0 5 7 2 0 0 1 5 3 6 2 7 2 2 0 3 3 2 1 3 7 2  
3 3 5 7 2 2 1 2 7 2 3 1 5 3 9 5 0 5 3 8 8 0 3 1 1  
1 3 7 1 9 1 4 1 1 9 1 2 9 1 2 8 3 1 9 1 7 0 1 4  
1 0 1 1 9 1 2 1 3 5 7 3 6 8 9 3 2 2 6 4 1 5 1 8 6  
6 3 5 9 7 2 0 2 9 9 2 9 1 7 2 2 5 1 0 0 4 6 7 0 1  
3 0 9 4 1 1 1 5 9 1 0 1 0 6 1 5 4 0 6 1 0 3 6 3 1  
1 0 4 4 1 1 1 0 3 0 4 3 5 3 6 2 0 0 1 7 7 9 9 6 6  
8 9 1 8 0 5 4 7 0 8 5 5 7 1 2 1 4 2 2 1 5 5 4 6 0  
1 0 1 1 2 3 0 1 0 7 1 1 3 9 9 1 0 8 9 9 2 0 9 8 4  
0 1 0 9 7 0 7 5 9 1 3 3 1 9 7 3 0 1 2 5 1 2 0 5 6  
1 0 7 4 3 1 8 2 5 5 1 8 2 8 1 4 3 5 8 0 9 0 9 4 3  
1 2 8 7 5 4 1 6 3 5 4 6 0 5 5 4 6 0 3 5 4 6 0 5 5  
1 8 2 5 5 1 0 8 5 0 3 0 8 2 5 2 0 1 3 9 4 0 1

# Application: Handwritten digit recognition

- Feature vectors: each image is  $28 \times 28$  pixels. Rearrange as a 784-vector  $\mathbf{x}$
- Training: learn  $k=10$  two-class 1 vs the rest SVM classifiers  $f_k(\mathbf{x})$
- Classification: choose class with most positive score

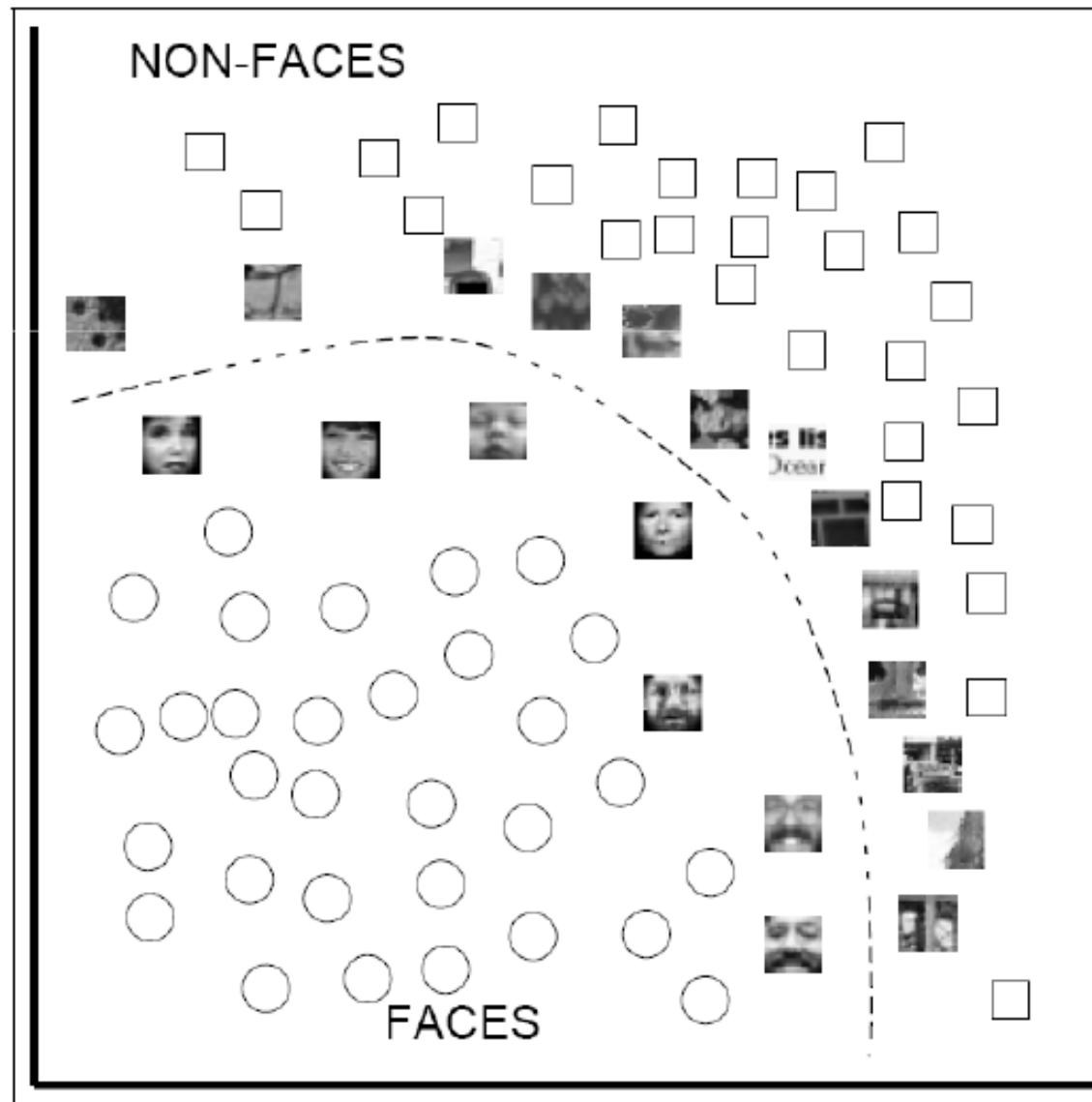
$$f(\mathbf{x}) = \max_k f_k(\mathbf{x})$$



# Guyon & Vapnik 1995

- **USPS benchmark**
  - 2.5% error: human performance
- **Different learning algorithms**
  - 16.2% error: Decision tree (C4.5)
  - 5.9% error: (best) 2-layer Neural Network
  - 5.1% error: LeNet 1 - (massively hand-tuned) 5-layer network
- **Different SVMs**
  - 4.0% error: Polynomial kernel ( $p=3$ , 274 support vectors)
  - 4.1% error: Gaussian kernel ( $\sigma=0.3$ , 291 support vectors)

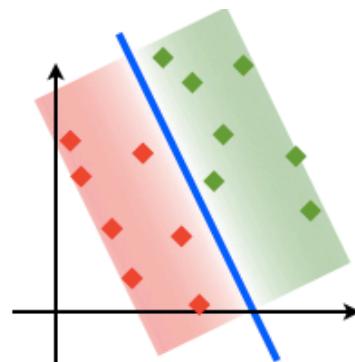
# Support vectors for Faces (P&P 98)



# Linear vs. Nonlinear

## Linear SVM

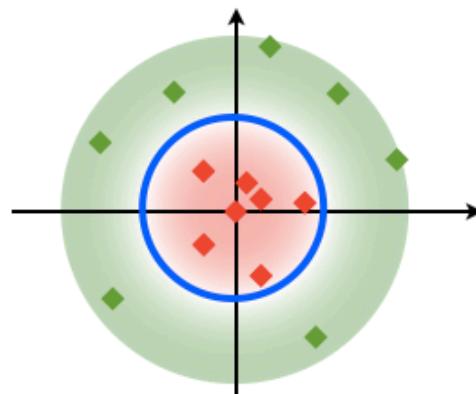
- ✓ fast
- ✗ restrictive



$$F(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$$

## Non-linear SVM

- ✗ much slower
- ✓ powerful



$$F(\mathbf{x}) = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

# Other kernels

- From <http://www.kernel-methods.net/kernels.html>

Kernel Functions Described in the Book:

Definition 9.1 Polynomial kernel 286  
 Computation 9.6 All-subsets kernel 289  
 Computation 9.8 Gaussian kernel 290  
 Computation 9.12 ANOVA kernel 293  
 Computation 9.18 Alternative recursion for ANOVA kernel 296  
 Computation 9.24 General graph kernels 301  
 Definition 9.33 Exponential diffusion kernel 307  
 Definition 9.34 von Neumann diffusion kernel 307  
 Computation 9.35 Evaluating diffusion kernels 308  
 Computation 9.46 Evaluating randomised kernels 315  
 Definition 9.37 Intersection kernel 309  
 Definition 9.38 Union-complement kernel 310  
 Remark 9.40 Agreement kernel 310  
 Section 9.6 Kernels on real numbers 311  
 Remark 9.42 Spline kernels 313  
 Definition 9.43 Derived subsets kernel 313  
 Definition 10.5 Vector space kernel 325  
 Computation 10.8 Latent semantic kernels 332  
 Definition 11.7 The p-spectrum kernel 342  
 Computation 11.10 The p-spectrum recursion 343  
 Remark 11.13 Blended spectrum kernel 344  
 Computation 11.17 All subsequences kernel 347  
 Computation 11.24 Fixed length subsequences kernel 352

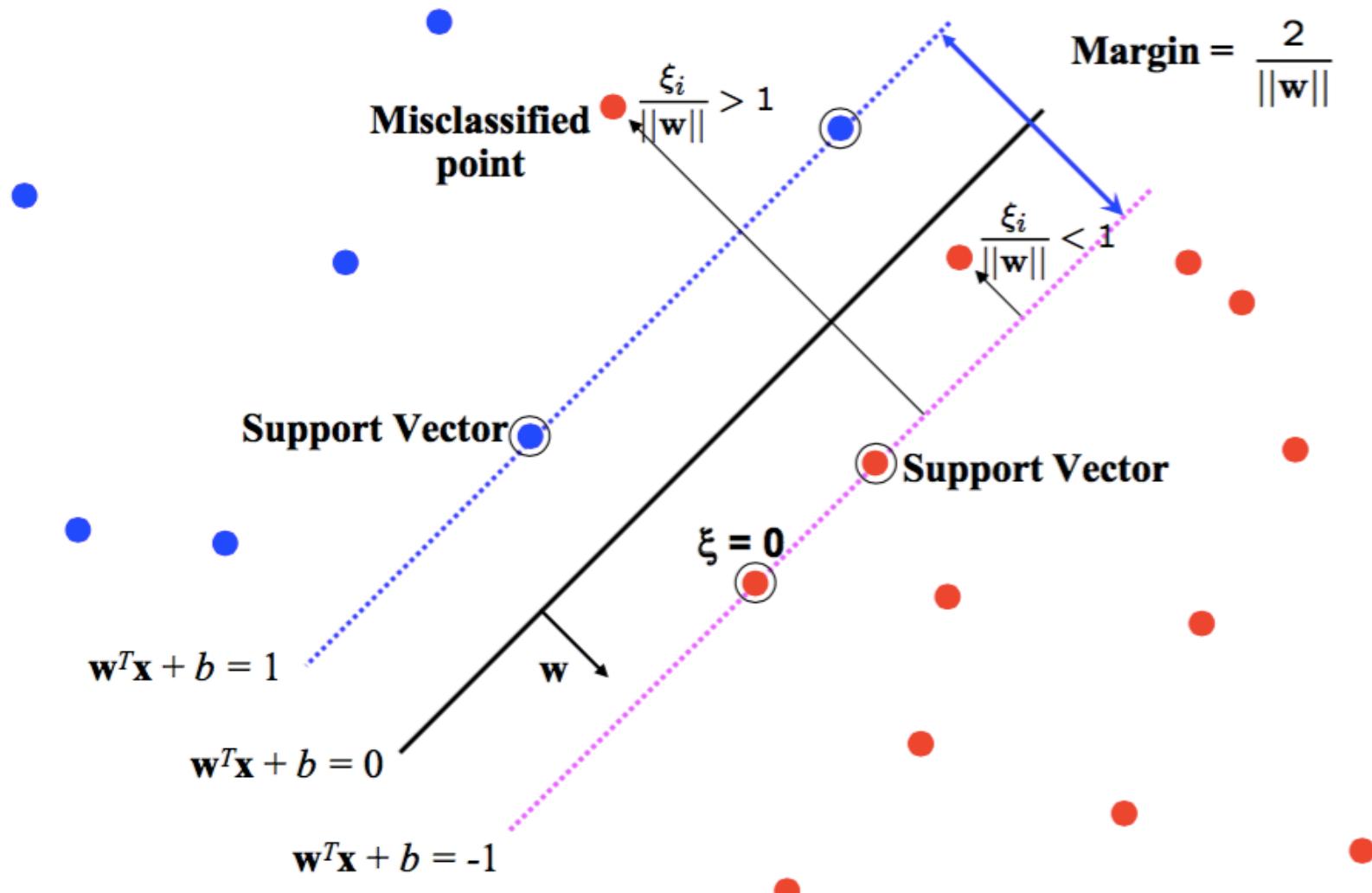
Computation 11.33 Naive recursion for gap-weighted subsequences kernel 358  
 Computation 11.36 Gap-weighted subsequences kernel 360  
 Computation 11.45 Trie-based string kernels 367  
 Algorithm 9.14 ANOVA kernel 294  
 Algorithm 9.25 Simple graph kernels 302  
 Algorithm 11.20 All non-contiguous subsequences kernel 350  
 Algorithm 11.25 Fixed length subsequences kernel 352  
 Algorithm 11.38 Gap-weighted subsequences kernel 361  
 Algorithm 11.40 Character weighting string kernel 364  
 Algorithm 11.41 Soft matching string kernel 365  
 Algorithm 11.42 Gap number weighting string kernel 366  
 Algorithm 11.46 Trie-based p-spectrum kernel 368  
 Algorithm 11.51 Trie-based mismatch kernel 371  
 Algorithm 11.54 Trie-based restricted gap-weighted kernel 374  
 Algorithm 11.62 Co-rooted subtree kernel 380  
 Algorithm 11.65 All-subtree kernel 383  
 Algorithm 12.8 Fixed length HMM kernel 401  
 Algorithm 12.14 Pair HMM kernel 407  
 Algorithm 12.17 Hidden tree model kernel 411  
 Algorithm 12.34 Fixed length Markov model Fisher kernel 427

# Text Classification: Examples

- Classify news stories as *World, US, Business, SciTech, Sports, Entertainment, Health, Other*
- Add MeSH terms to Medline abstracts
  - e.g. “Conscious Sedation” [E03.250]
- Classify business names by industry.
- Classify student essays as *A,B,C,D*, or *F*.
- Classify email as *Spam, Other*.
- Classify email to tech staff as *Mac, Windows, ..., Other*.
- Classify pdf files as *ResearchPaper, Other*
- Classify documents as *WrittenByReagan, GhostWritten*
- Classify movie reviews as *Favorable, Unfavorable, Neutral*.
- Classify technical papers as *Interesting, Uninteresting*.
- Classify jokes as *Funny, NotFunny*.
- Classify web sites of companies by Standard Industrial Classification (SIC) code.



# Slack variables: let us make (but also pay) some errors



# Objective for non-separable data

$$\min_{\mathbf{w}, \xi} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi^i$$

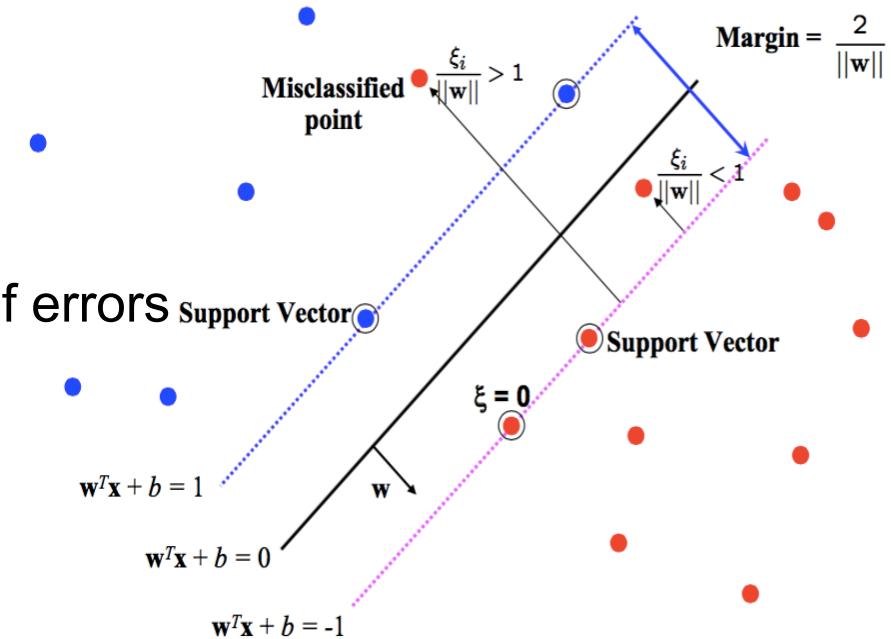
s.t. :  $y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi^i, \quad \forall i$

$$\xi^i \geq 0, \quad \forall i$$

misclassification when  $\xi > 1$

$\sum_i \xi^i$  : upper bound on number of errors

C: hyperparameter  
(cross-validation!)



# Support Vectors

From complementary slackness (KKT)  $\mu_i g_i(x) = 0$

where  $g_i(x) = 1 - y^i(\mathbf{w}^T \mathbf{x}^i + b)$  ( $\leq 0$ )

Therefore:  $\mu_i \neq 0 \rightarrow y^i(\mathbf{w}^T \mathbf{x}^i + b) = 1$

*Interpretation:  $\mu$  is nonzero only for points on the margin (hardest points)*

From minimum w.r.t.  $\mathbf{w}$ :  $\mathbf{w}^* = \sum_{i=1}^M \mu_i y^i \mathbf{x}^i$

*Interpretation: only points on the margin contribute to the solution*

- ‘Support Vectors’

Intuitively ok: we want to maximize the margins of the hardest cases

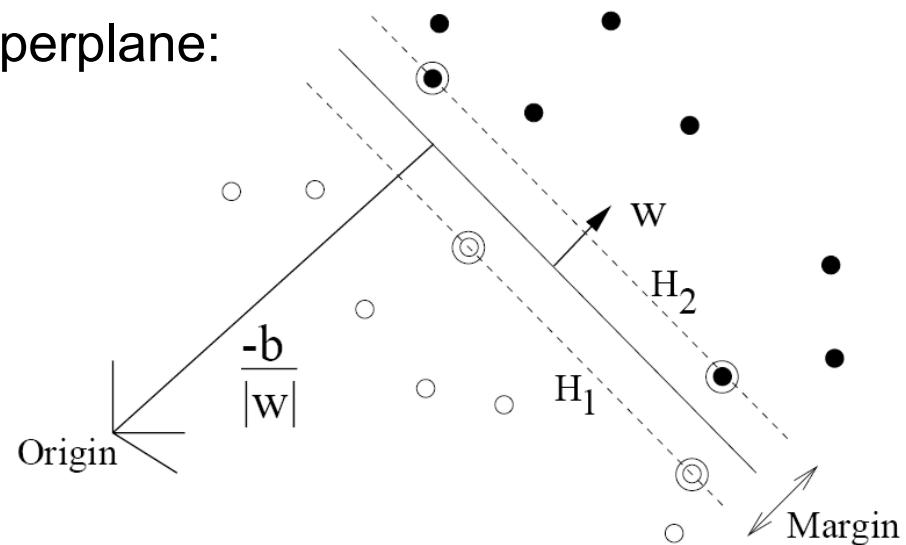
# Decision Hyperplanes & Support Vectors

Use support vectors to determine  $b^*$ :

$$y^i(\mathbf{w}^T \mathbf{x}^i + b) = 1, \quad \forall i \in S$$

$$b^* = \frac{1}{N_S} \sum_{i \in S} (y^i - \mathbf{w}^T \mathbf{x}^i)$$

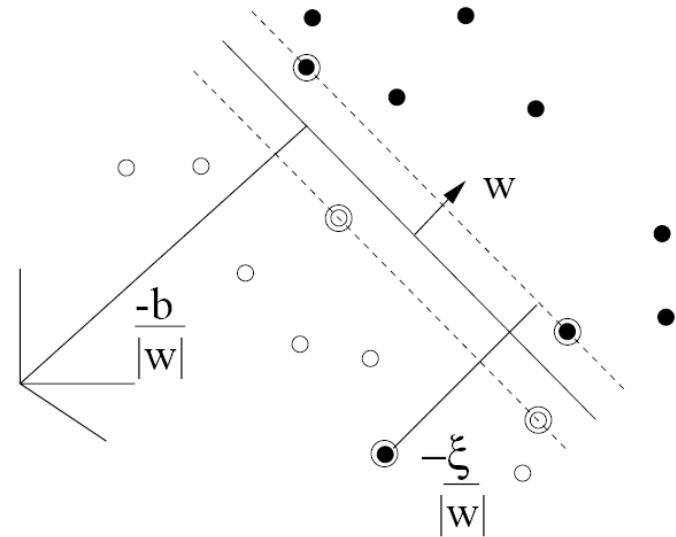
Support Vector Machine decision hyperplane:



# Non-separable data

Primal:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^M \xi^i \\ \text{s.t.} \quad & y^i (\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi^i \\ & \xi^i \geq 0 \end{aligned}$$



Lagrangian:

$$L(\mathbf{w}, b, \xi, \mu, \nu) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^M \xi^i - \sum_{i=1}^M \mu_i [y^i (\mathbf{w}^T \mathbf{x}^i + b) - 1 + \xi] - \sum_{i=1}^M \nu_i \xi_i$$

$$\text{Dual: } \max_{\mu} \quad \sum_{i=1}^M \mu_i - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M y^i y^j \mu_i \mu_j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \quad \begin{array}{l} \mu_i \geq 0, \quad \forall i \\ \nu_i \geq 0, \quad \forall i \end{array}$$

$$\text{s.t.} \quad 0 \leq \mu_i \leq C$$

$$\sum_{i=1}^M \mu_i y^i = 0$$

## KKT conditions – nonseparable case

$$C - \mu^i - \nu^i = 0 \quad (1)$$

$$y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) - 1 + \xi^i \geq 0 \quad (2)$$

$$\mu^i [y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) - 1 + \xi^i] = 0 \quad (3) \quad \text{Complementary slackness}$$

$$\nu^i \xi^i = 0 \quad (4) \quad \text{Complementary slackness}$$

$$\xi^i \geq 0 \quad (5)$$

$$\mu^i \geq 0 \quad (6)$$

$$\nu^i \geq 0 \quad (7)$$

Case analysis:  $y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) > 1 \xrightarrow[3:\xi^i > 0]{ } \mu_i = 0$

$y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) < 1 \xrightarrow[2:\xi^i > 0 \rightarrow 4:\nu^i = 0 \rightarrow 1:]{ } \mu_i = C$

$y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) = 1 \xrightarrow[3:\mu^i \xi^i = 0 \rightarrow 1:]{ } \mu_i \in [0, C]$

*Interpretation: influence,  $\mu$ , of any training point is bounded in  $[0, C]$*

# Hinge Loss

$$C - \mu^i - \nu^i = 0 \quad (1)$$

$$\mu^i [y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b) - 1 + \xi^i] = 0 \quad (3)$$

$$\nu^i \xi^i = 0 \quad (4)$$

$$\xi^i \geq 0 \quad (5)$$

$$\mu^i \neq 0 \xrightarrow{(3)} \xi^i = 1 - y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b)$$

$$\mu^i = 0 \xrightarrow{(1)} \nu^i = C \xrightarrow{(4)} \xi^i = 0$$

$$\xi^i = \max(0, 1 - y^i(\langle \mathbf{x}^i, \mathbf{w} \rangle + b))$$

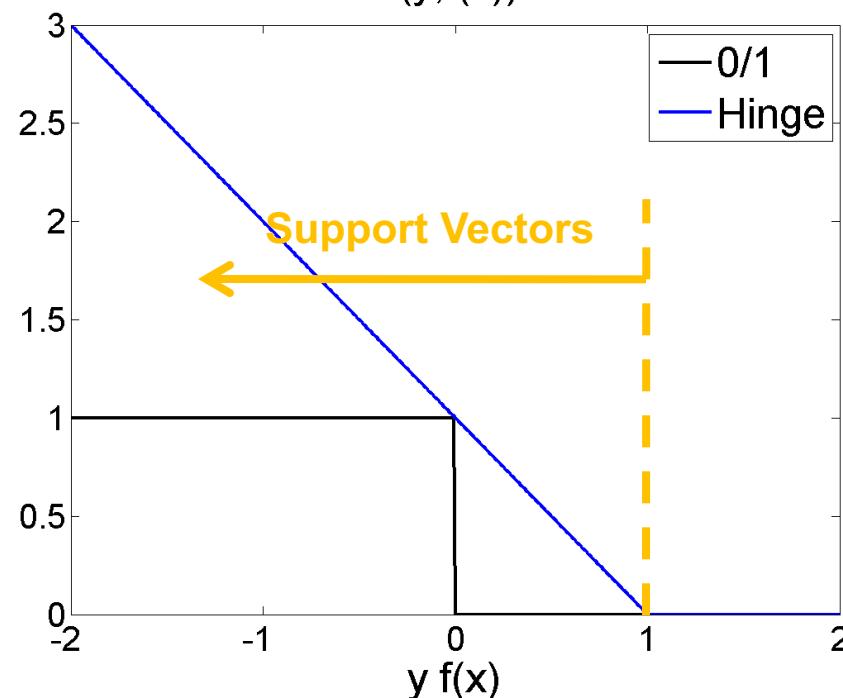
$$\begin{array}{ll} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^M \xi^i \\ s.t. & y^i (\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi^i \\ & \xi^i \geq 0 \end{array} \quad \longleftrightarrow \quad L(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y^i h_{\mathbf{w}, b}(\mathbf{x}^i))$$

# Loss function for SVM training

Optimization problem:

$$\begin{aligned}
 L(\mathbf{w}) &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y^i h_{\mathbf{w}, b}(x^i)) \\
 &\propto \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^N \underbrace{\max(0, 1 - y^i h_{\mathbf{w}, b}(x^i))}_{l(y^i, x^i)}
 \end{aligned}$$

Hinge loss:



# Appendix

- Primal and Dual form of SVMs: the full story

## References:

- S. Boyd and L. Vandenberghe: Convex Optimization (textbook)
- C. Burges: A tutorial on SVMs for pattern recognition

# Duality

- Constrained optimization problem:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & h_i(w) = 0, \quad i = 1 \dots l \\ & g_i(w) \leq 0, \quad i = 1 \dots m \end{aligned}$$

- Equivalent to unconstrained problem:

$$\min_w f_{uc}(w) = f(w) + \sum_{i=1}^l I_0(h_i(w)) + \sum_{i=1}^m I_+(g_i(w))$$

$$I_0(x) = \begin{cases} 0, & x = 0 \\ \infty, & x \neq 0 \end{cases}, \quad I_+(x) = \begin{cases} 0, & x \leq 0 \\ \infty, & x > 0 \end{cases}$$

- Soften constraint terms  $I_0(x_i) \rightarrow \lambda_i x_i$        $I_+(x_i) \rightarrow \mu_i x_i, \mu > 0$

# Lagrangian

- Replace hard constraints with soft ones

$$\min_w \quad f_{uc}(w) = f(w) + \sum_{i=1}^l I_0(h_i(w)) + \sum_{i=1}^m I_+(g_i(w))$$

$$L(w, \lambda, \mu) = f(w) + \sum_{i=1}^l \lambda_i h_i(w) + \sum_{i=1}^m \mu_i g_i(w), \quad \mu_i > 0 \forall i$$

- Observe that

$$f_{uc}(w) = \max_{\lambda, \mu: \mu_i > 0} L(w, \lambda, \mu)$$

- At an optimum:

$$f(w^*) = \min_w \max_{\lambda, \mu: \mu_i > 0} L(w, \lambda, \mu)$$



You do your worst, and we will do our best

# Lagrange Dual Function

- Form  $\theta(\lambda, \mu) = \inf_w L(w, \lambda, \mu)$

- $\theta$  : lower bound on optimal value of the original problem

$$\begin{aligned}
 L(w^*, \lambda, \mu) &= f(w^*) + \sum_{i=1}^l \lambda_i h_i(w^*) + \sum_{i=1}^m \mu_i g_i(w^*) = \\
 &\stackrel{w^*: \text{feasible}}{=} f(w^*) + \sum_{i=1}^l \lambda_i 0 + \sum_{i=1}^m \mu_i \underbrace{g_i(w^*)}_{< 0} = \\
 &\stackrel{\mu_i > 0}{\leq} f(w^*)
 \end{aligned}$$

- Therefore:  $\theta(\lambda, \mu) = \inf_w L(w, \lambda, \mu) \leq L(w^*, \lambda, \mu) \leq f(w^*)$

# Dual Problem

- Maximize the lower bound on the cost of the primal

$$\begin{aligned} \max_{\lambda, \mu} \quad & \theta(\lambda, \mu) \\ \text{s.t.} \quad & \mu_i > 0 \quad \forall i \end{aligned}$$

- In general:

$$\begin{aligned} d^* &= \max_{\lambda, \mu} \theta(\lambda, \mu) \\ &= \max_{\lambda, \mu: \mu_i > 0} \min_w L(w, \lambda, \mu) \\ &\leq \min_w \max_{\lambda, \mu: \mu_i > 0} L(w, \lambda, \mu) \\ &= \min_w f_{uc}(w) = p^* \end{aligned}$$

- For convex cost and convex constraints (SVM case):  $d^* = p^*$

# Complementary Slackness

- Assume  $d^* = p^*$
- There exists a feasible solution  $w^*, \lambda^*, \mu^*$  to the primal and dual problems, such that  $f(w^*) = \theta(\lambda^*, \mu^*)$
- We will have
 
$$\begin{aligned}
 f(w^*) &= \theta(\lambda^*, \mu^*) \\
 &= \inf_w f(w) + \sum_{i=1}^l \lambda_i^* h_i(w) + \sum_{i=1}^m \mu_i^* f_i(w) \\
 &\leq f(w^*) + \sum_{i=1}^M \lambda_i^* h_i(w^*) + \sum_{i=1}^m \mu_i^* f_i(w^*) \\
 &\leq f(w^*)
 \end{aligned}$$
- This means  $\mu_i^* f_i(w^*) = 0, \forall i$

# Karush-Kuhn Tucker (KKT) Conditions

- Solution of the primal problem:
  - minimum of the Lagrangian w.r.t. the primal variables

– therefore  $\nabla f(w^*) + \sum_{i=1}^l \lambda_i \nabla h_i(w^*) + \sum_{i=1}^m \nabla f_i(w^*) = 0$

- Putting all constraints together: KKT conditions

$$h_i(w^*) = 0$$

$$f_i(w^*) \leq 0$$

$$\mu_i f_i(w^*) = 0$$

$$\mu_i \geq 0$$

$$\nabla f(w^*) + \sum_{i=1}^l \lambda_i \nabla h_i(w^*) + \sum_{i=1}^m \nabla f_i(w^*) = 0$$

# Problem Lagrangian

Primal:

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} |\mathbf{w}|^2$$

$$s.t. \quad -y^i(\mathbf{w}^T \mathbf{x}^i + b) + 1 \leq 0, \quad i = 1 \dots M$$

Lagrangian:  $L(\mathbf{w}, b, \mu) = \frac{1}{2} |\mathbf{w}|^2 - \sum_{i=1}^M \mu_i [y^i(\mathbf{w}^T \mathbf{x}^i + b) - 1] \quad \mu_i \geq 0$

Optimum w.r.t.  $\mathbf{w}$ :  $0 = \mathbf{w}^* - \sum_{i=1}^M \mu_i [y^i \mathbf{x}^i] \quad \mathbf{w}^* = \sum_{i=1}^M \mu_i y^i \mathbf{x}^i$

Optimum w.r.t.  $b$ :  $0 = \sum_{i=1}^M \mu_i y^i$

# Dual for Large-Margin Classifier-I

Plug optimal values into Lagrangian:

$$\begin{aligned}
 \theta(\mu) &= L(\mathbf{w}^*, b^*, \mu) \\
 &= \frac{1}{2} |\mathbf{w}^*|^2 - \sum_{i=1}^M \mu_i [y^i (\mathbf{w}^{*T} x^i + b) - 1] \\
 &= \frac{1}{2} \left( \sum_{i=1}^M \mu_i y^i x^i \right)^T \left( \sum_{j=1}^M \mu_j y^j x^j \right) - \sum_{i=1}^M \mu_i [y^i \left( \left( \sum_{j=1}^M \mu_j y^j x^j \right)^T x^i + b \right) - 1] \\
 &= \sum_{i=1}^M \mu_i - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \mu_i \mu_j y^i y^j (x^i)^T (x^j) - b \sum_{i=1}^M \mu_i y^i
 \end{aligned}$$

# Dual for Large-Margin Classifier-II

Equivalent optimization problem:

$$\max_{\mu} \quad \theta(\mu) = \sum_{i=1}^M \mu_i - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \mu_i \mu_j y^i y^j \langle x^i, x^j \rangle$$

$$s.t. \quad \mu_i > 0, \quad \forall i$$

$$\sum_{i=1}^M \mu_i y^i = 0$$