

---

# Support Vector Machine

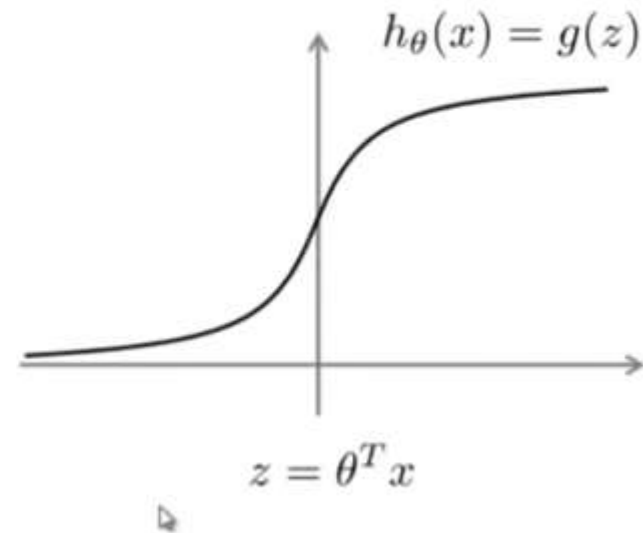
---

Some slides were adapted/taken from various sources, including Prof. Andrew Ng's Coursera Lectures, Stanford University, Prof. Kilian Q. Weinberger's lectures on Machine Learning, Cornell University, Prof. Sudeshna Sarkar's Lecture on Machine Learning, IIT Kharagpur, Prof. Bing Liu's lecture, University of Illinois at Chicago (UIC), CS231n: Convolutional Neural Networks for Visual Recognition lectures, Stanford University, Patrick Winston, MIT OpenCourseWare and many more. We thankfully acknowledge them. Students are requested to use this material for their study only and **NOT** to distribute it.

# Optimization objective

## Alternative view of logistic regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



If  $y = 1$ , we want  $h_{\theta}(x) \approx 1$ ,  $\theta^T x \gg 0$

— If  $y = 0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$  —

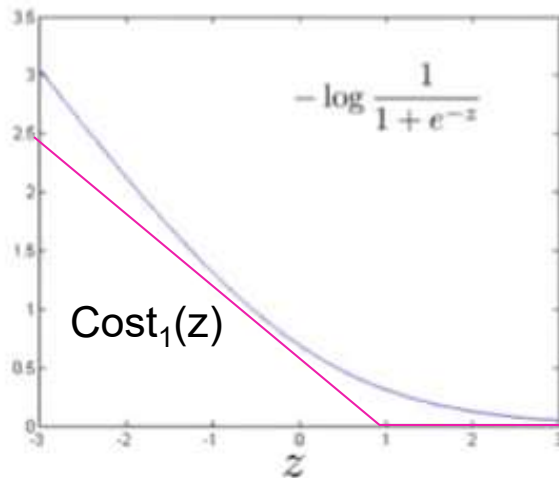
# SVM

## Alternative view of logistic regression

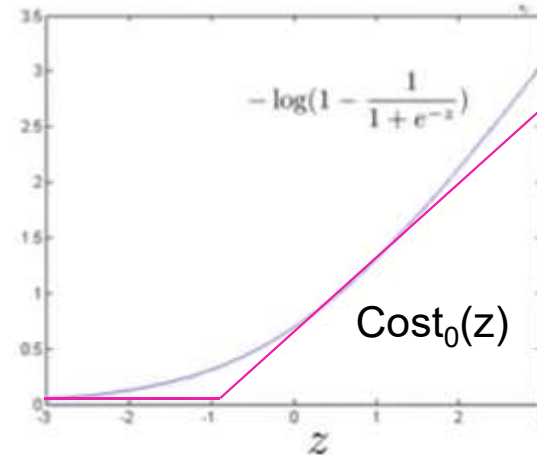
Cost of example:  $-(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log(1 - \frac{1}{1 + e^{-\theta^T x}})$$

If  $y = 1$  (want  $\theta^T x \gg 0$ ):



If  $y = 0$  (want  $\theta^T x \ll 0$ ):



# SVM

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \underbrace{\left( -\log h_{\theta}(x^{(i)}) \right)}_{\text{Cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left( -\log(1 - h_{\theta}(x^{(i)})) \right)}_{\text{Cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Support vector machine:  $\text{Cost}_1(\theta^T x^{(i)})$   $\text{Cost}_0(\theta^T x^{(i)})$

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{Cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{Cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$A + \lambda B = C A + B \quad \text{where } C = 1/\lambda$$

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

# SVM

## SVM hypothesis

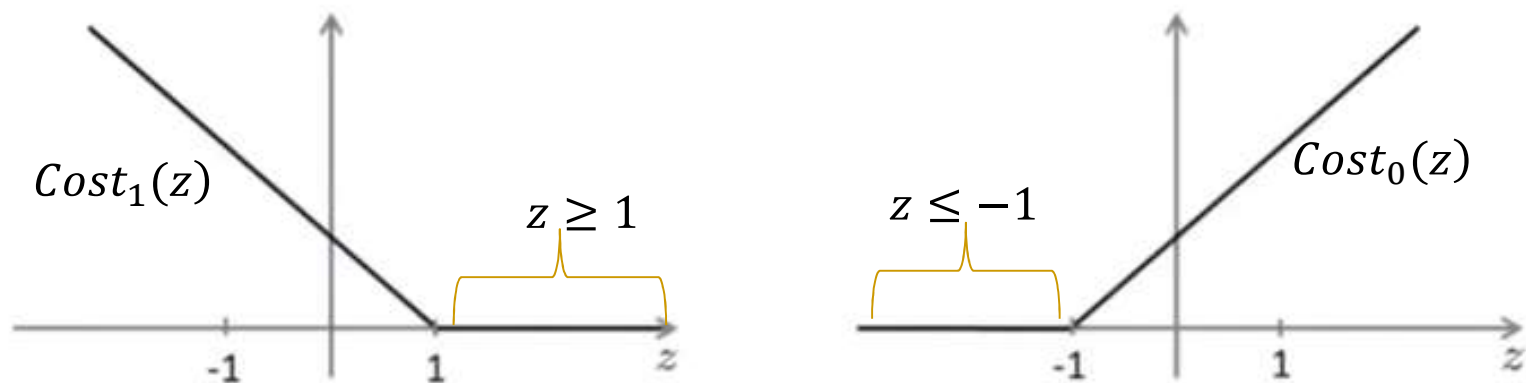
$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis:

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

# SVM: As Large Margin Classifier

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



If  $y = 1$ , we want  $\theta^T x \geq 1$  (not just  $\geq 0$ )

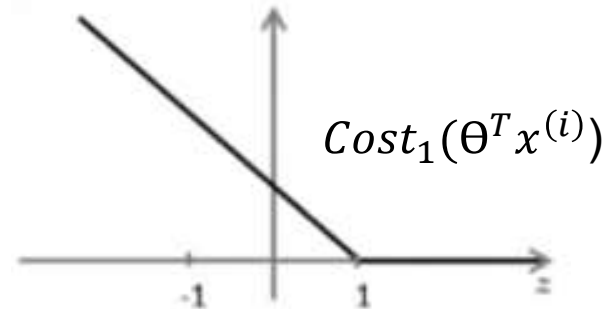
If  $y = 0$ , we want  $\theta^T x \leq -1$  (not just  $< 0$ )

# SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

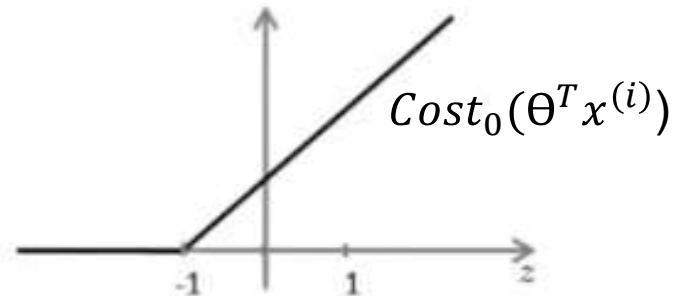
Whenever  $y^{(i)} = 1$ :

$$\theta^T x^{(i)} \geq 1$$



Whenever  $y^{(i)} = 0$ :

$$\theta^T x^{(i)} \leq -1$$



# SVM Decision Boundary

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

**= 0** as C is very big number

Whenever  $y^{(i)} = 1$ :

$$\theta^T x^{(i)} \geq 1 \quad \min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \text{ i.e. } \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

**s.t.**

$$\theta^T x^{(i)} \begin{cases} \geq 1 & \text{if } y^{(i)} = 1 \\ \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$$

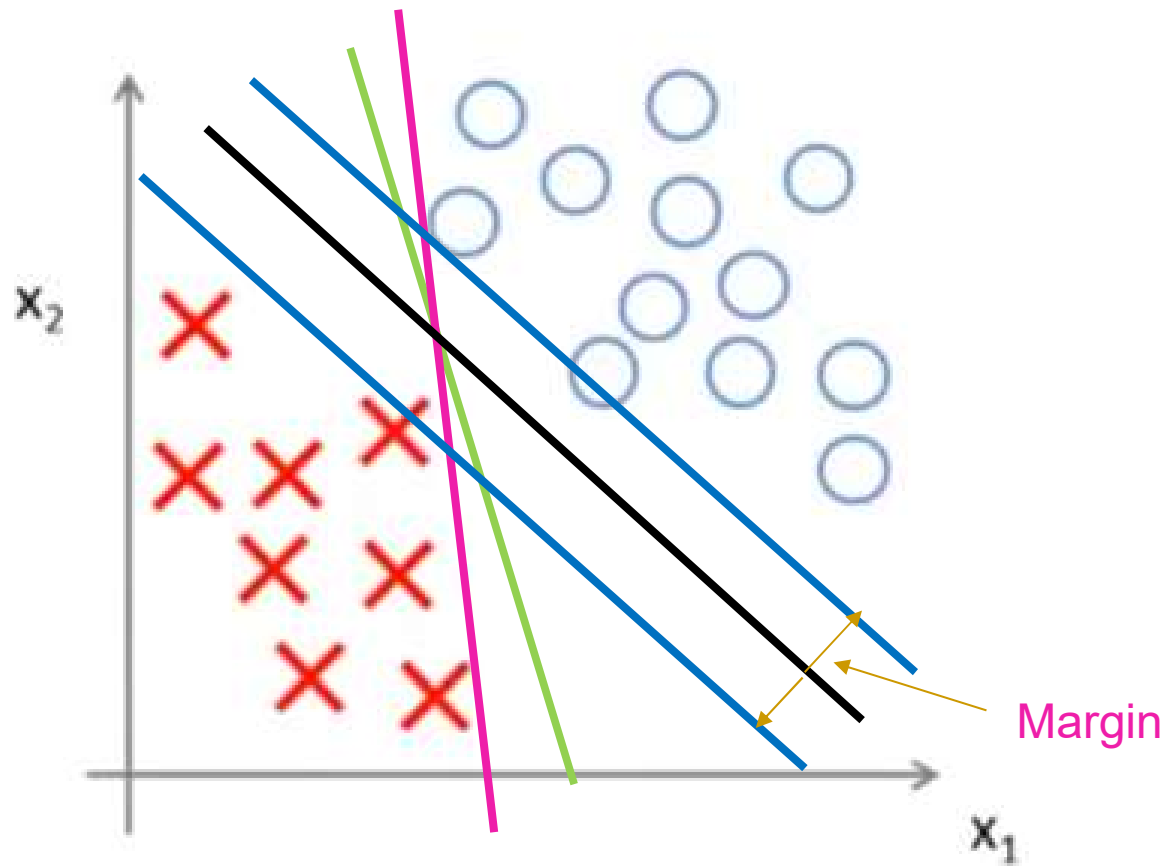
Whenever  $y^{(i)} = 0$ :

$$\theta^T x^{(i)} \leq -1$$



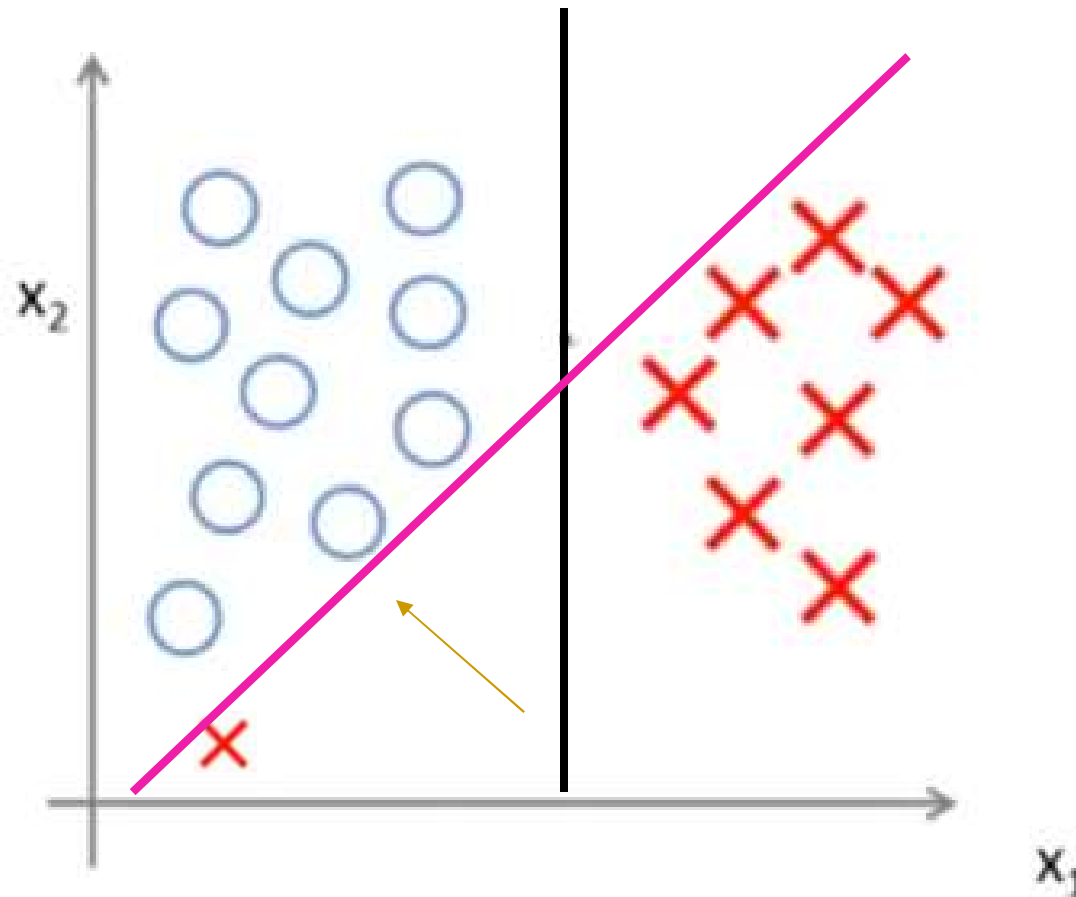
# SVM Decision Boundary

Linearly separable case



# Large Margin Classifier

In case of Outliers



C is very large

Sensitive to  
outliers

---

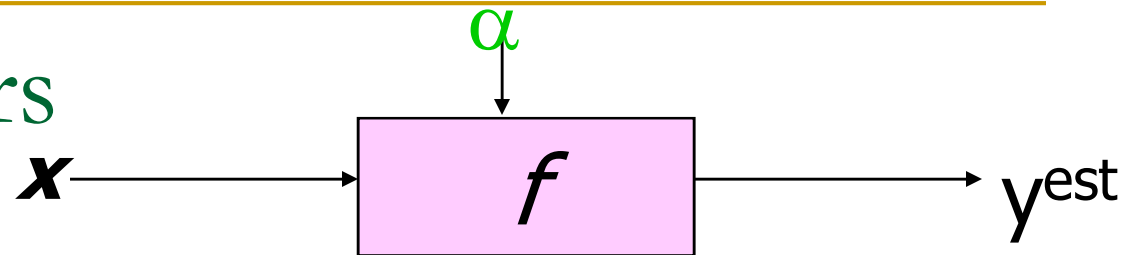
to continue...

---

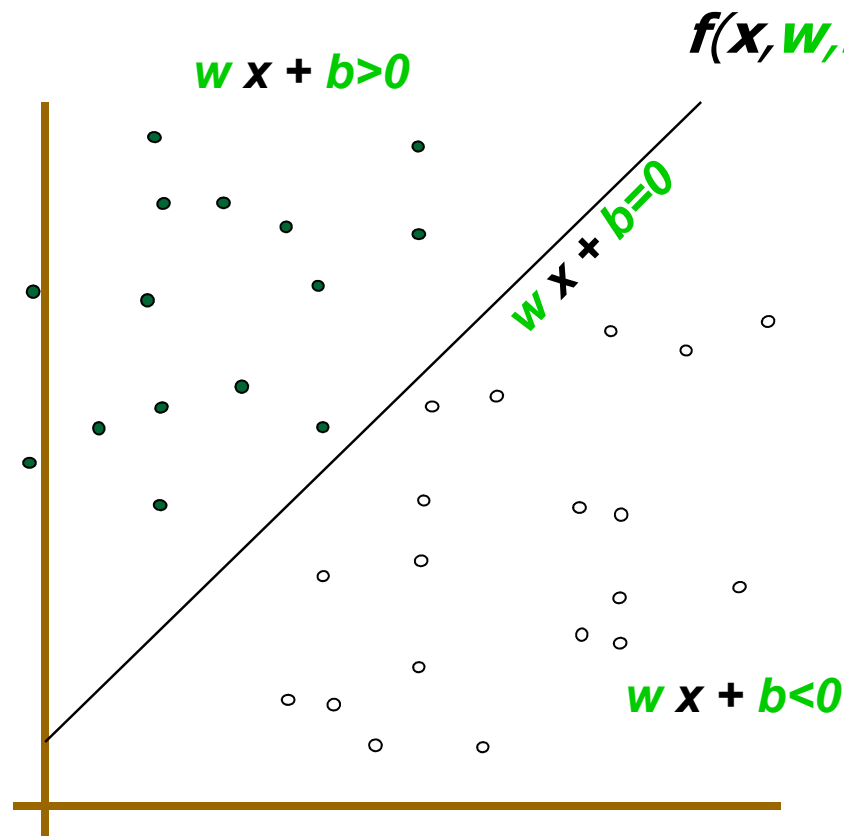
# Support Vector Machine

---

# Linear Classifiers



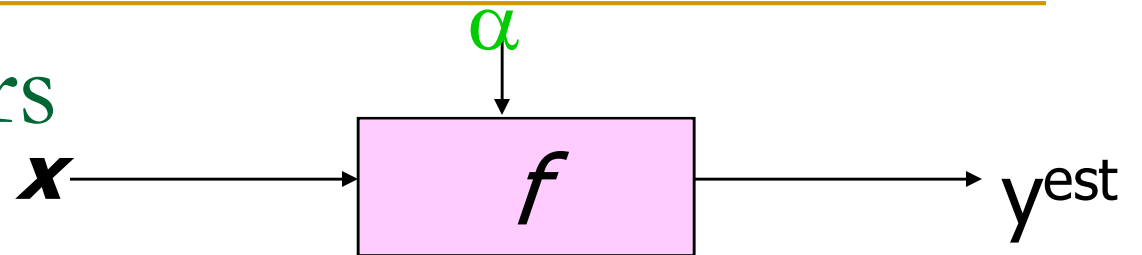
- denotes +1
- denotes -1



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w}\mathbf{x} + b)$$

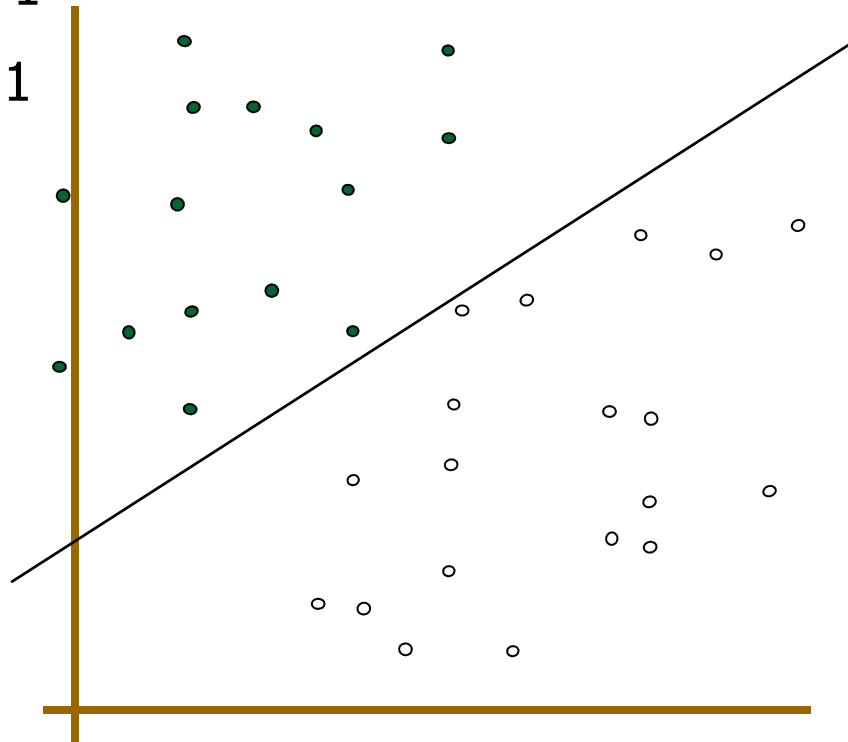
How would you classify this data?

# Linear Classifiers



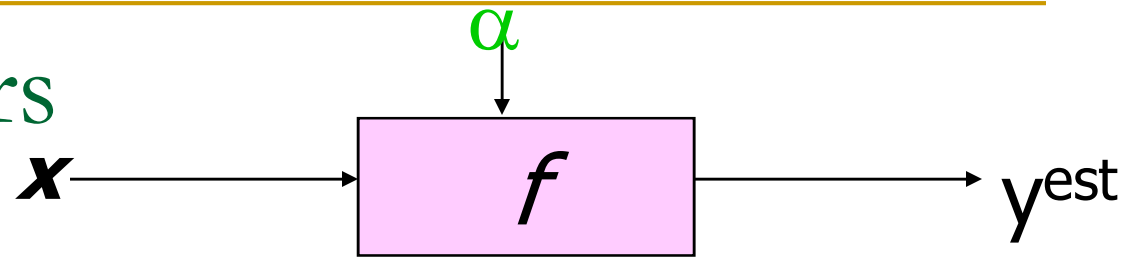
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

- denotes +1
- denotes -1

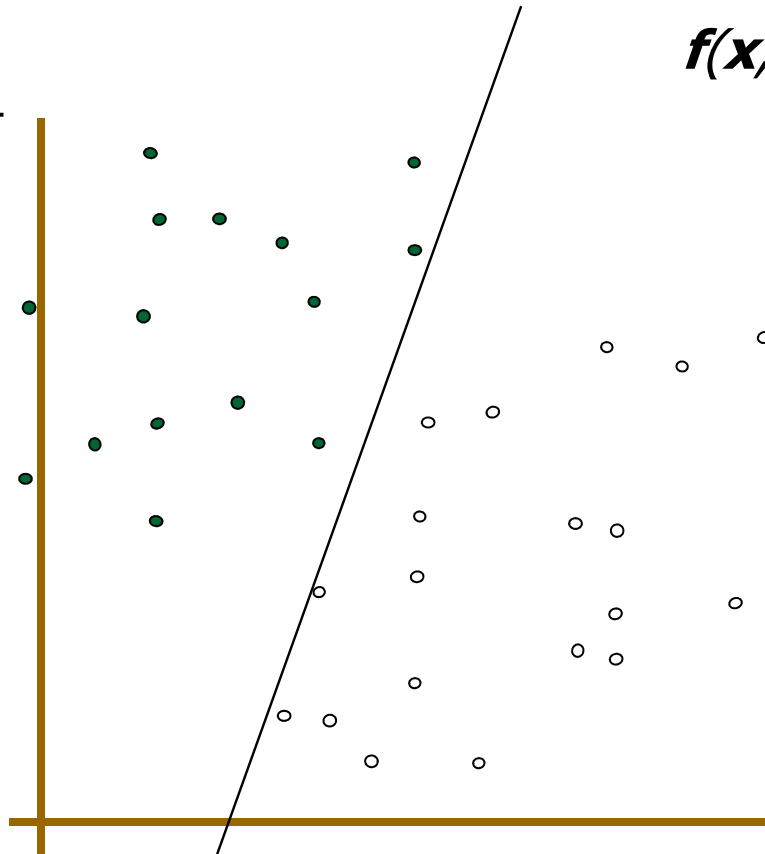


How would you classify this data?

# Linear Classifiers



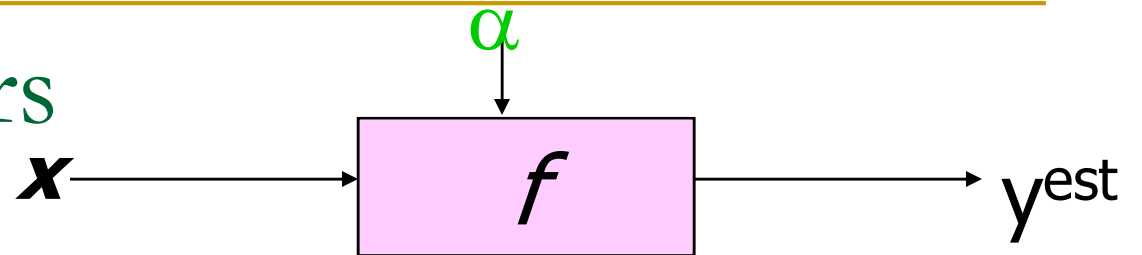
- denotes +1
- denotes -1



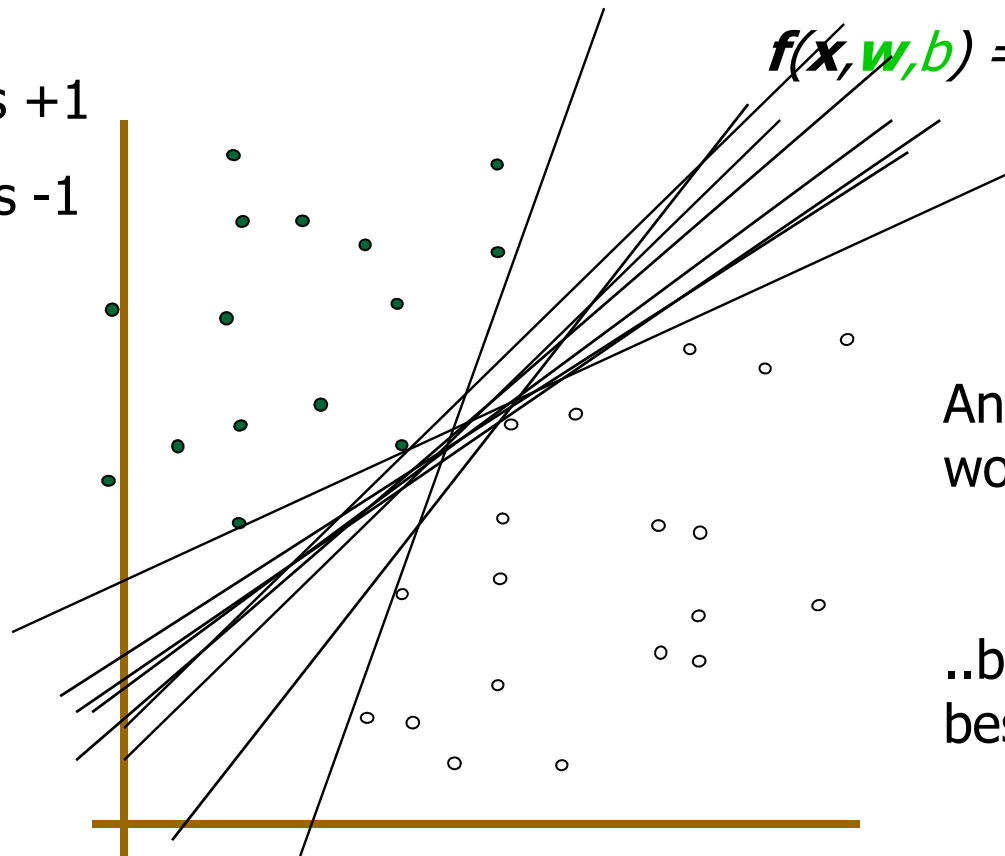
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

How would you classify this data?

# Linear Classifiers



- denotes +1
- denotes -1



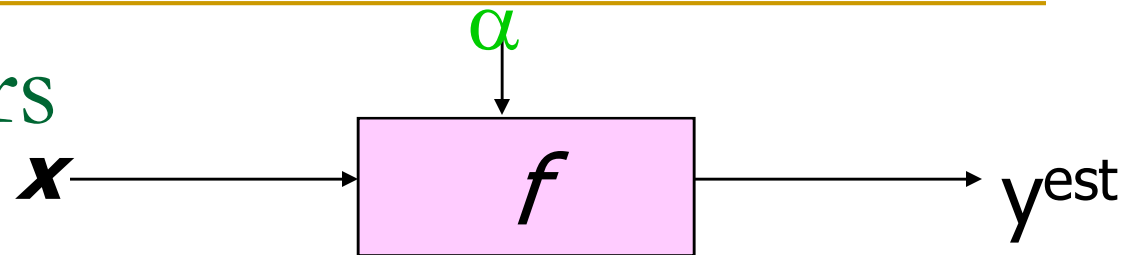
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

Any of these  
would be fine..

..but which is  
best?



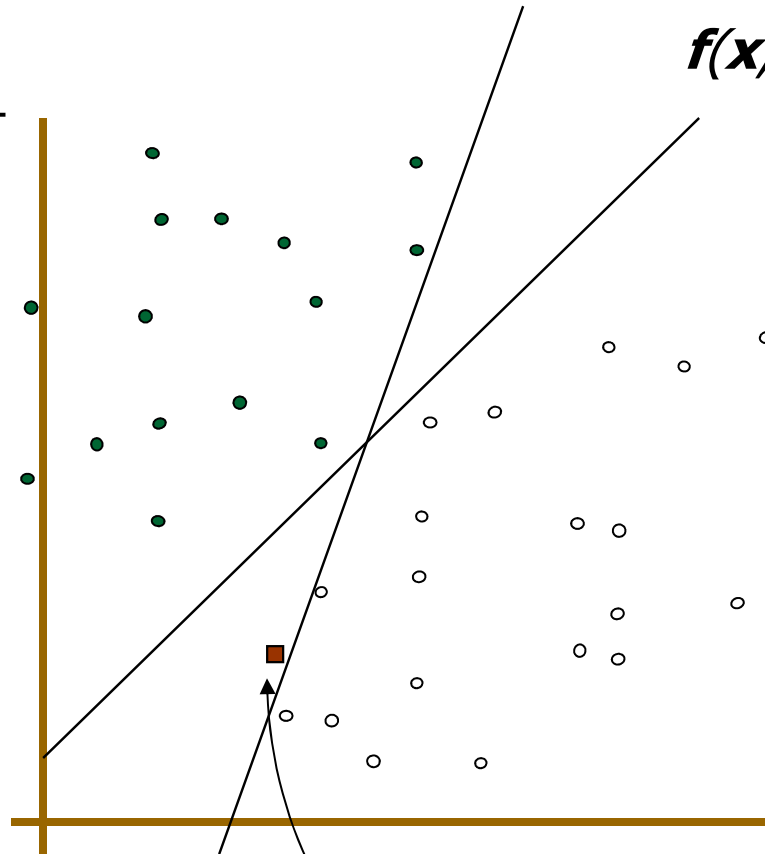
# Linear Classifiers



- denotes +1
- denotes -1

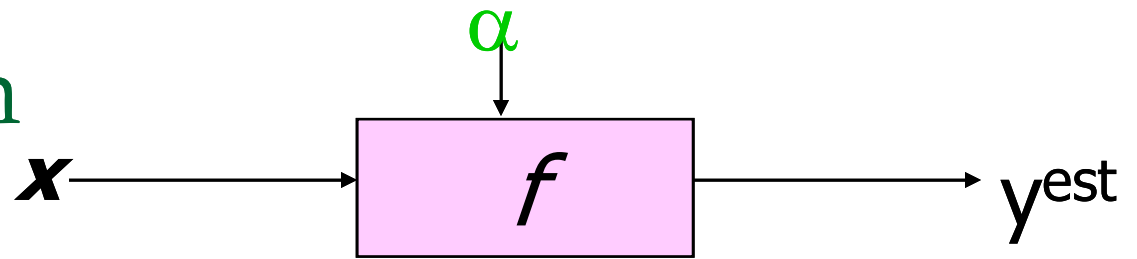
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

How would you classify this data?



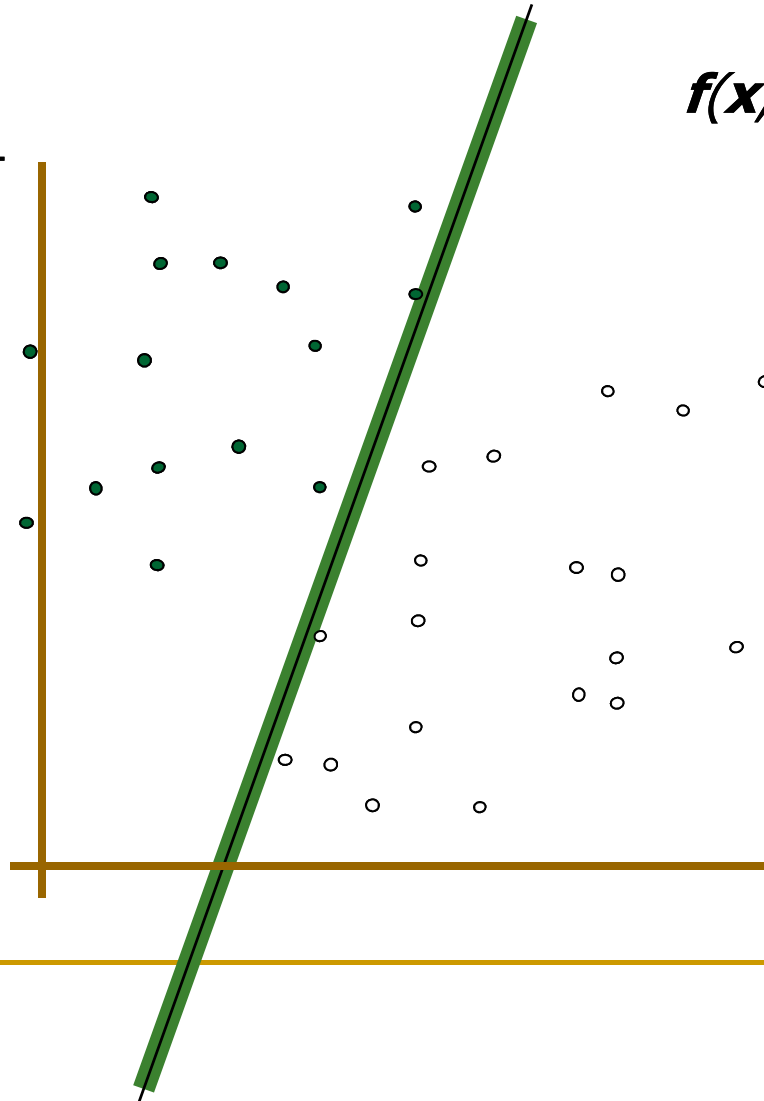
Misclassified  
to +1 class

# Classifier Margin



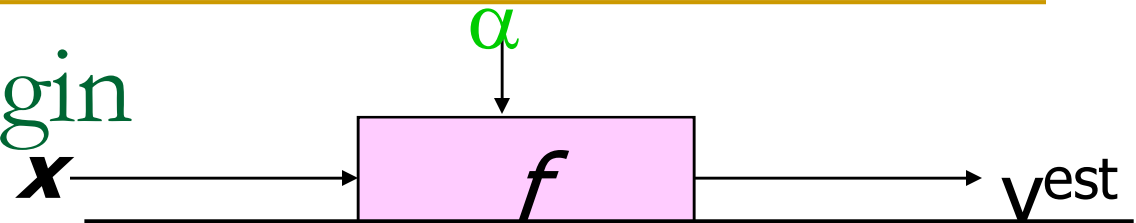
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

- denotes +1
- denotes -1



Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin



1. Maximizing the margin is good according to intuition and PAC\* theory
2. Implies that only support vectors are important; other training examples are ignorable.
3. Empirically it works very very well.

- denotes +1
- denotes -1

Support Vectors

are those datapoints that the margin pushes up against

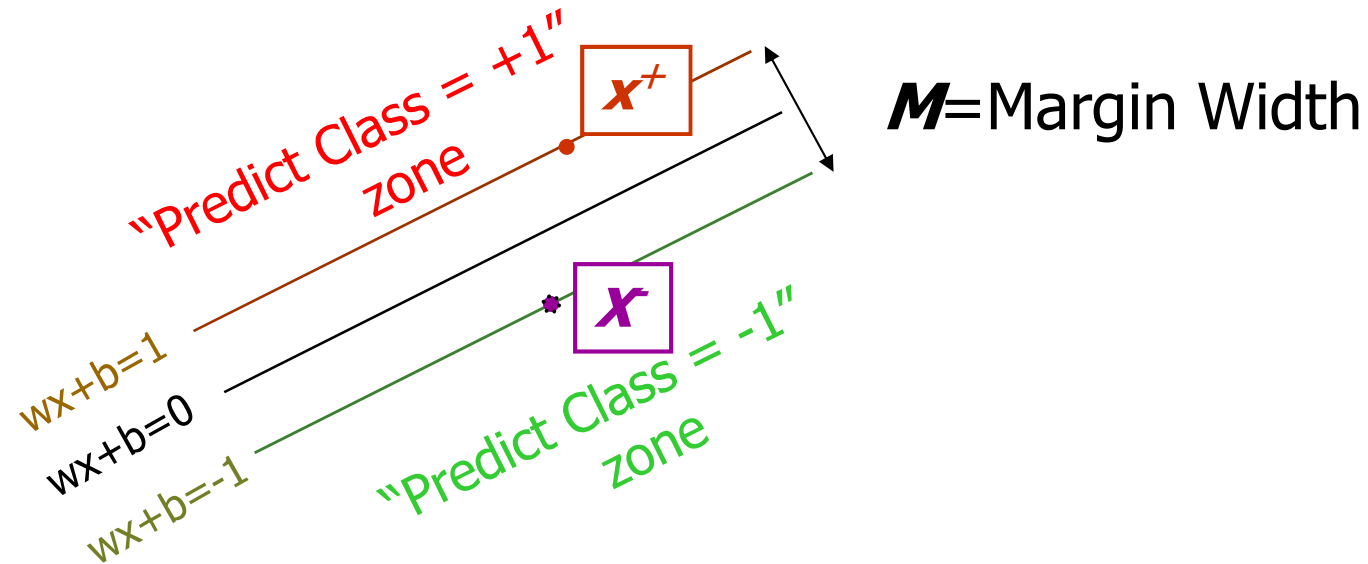
linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

\*Probability Approximately Correct

# Linear SVM Mathematically

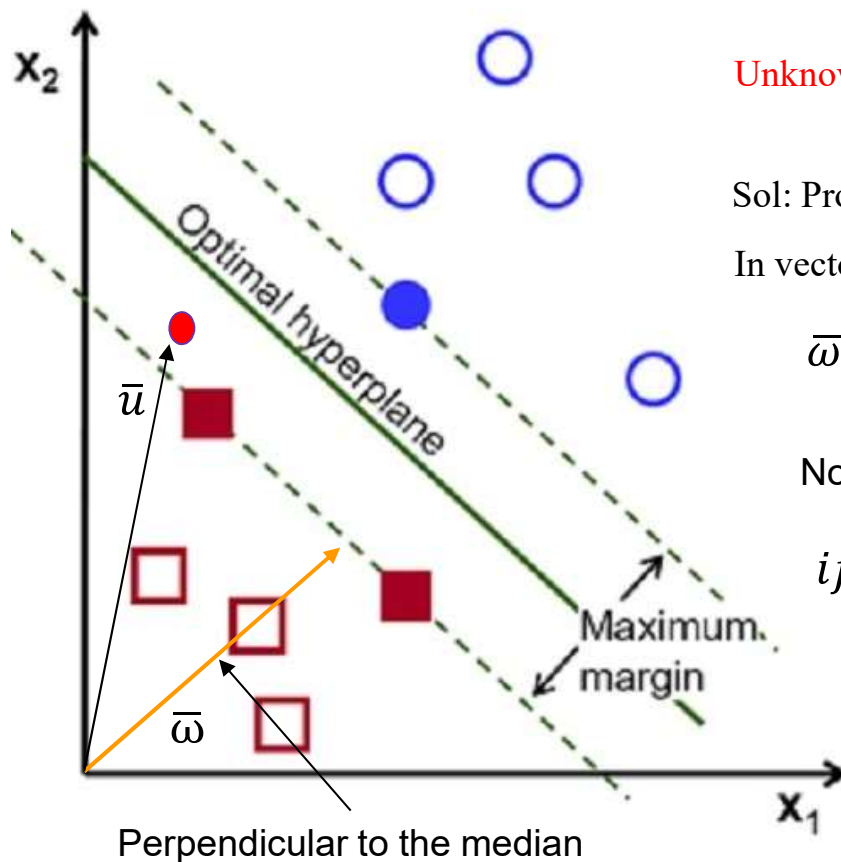


What we know:

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $w \cdot (x^+ - x^-) = 2$

$$M = \frac{(x^+ - x^-) \cdot w}{|w|} = \frac{2}{|w|}$$

# Linear SVM Mathematically



Unknown ( $\bar{u}$ ) is left side or right side of the DB?

Sol: Project the vector ( $\bar{u}$ ) on the perpendicular of the DB

In vector representation,

$$\bar{w} \cdot \bar{u} \geq C \quad \text{Here dot product is taking a projection of } \bar{u} \text{ on } \bar{w}$$

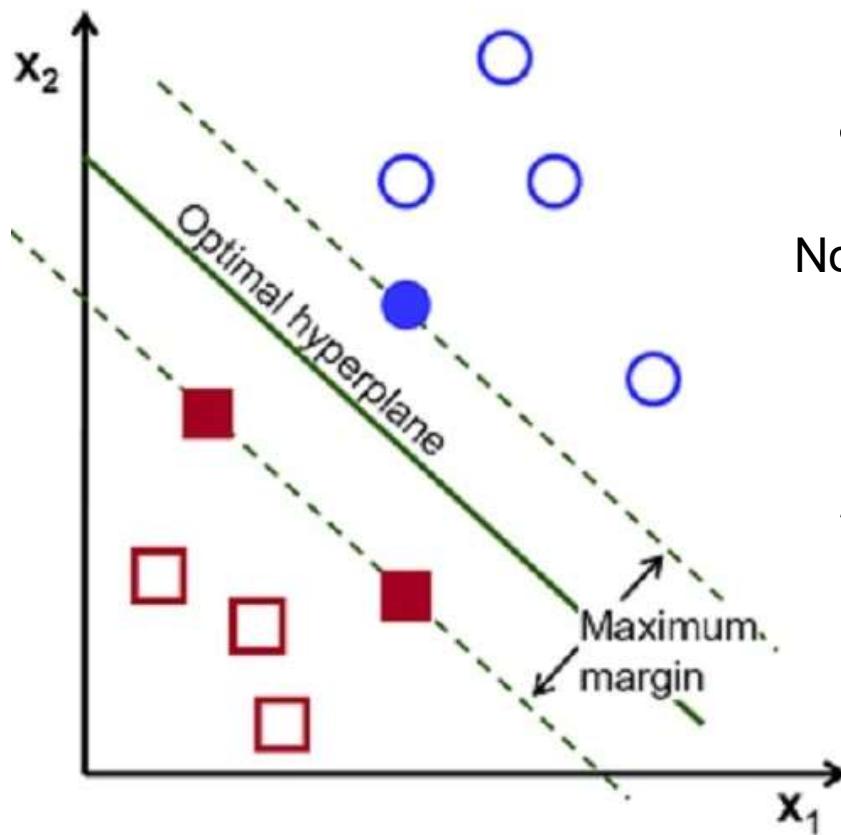
Now, without loss of generality, we can write

if  $\bar{w} \cdot \bar{u} + b \geq 0$  then + (○) : Decision rule

But the problem is, we don't know  $b$  and  $\bar{w}$  only we know that  $\bar{w}$  is perpendicular to DB, but there may be many such  $\bar{w}$  of different lengths.

$$C = -b$$

# Linear SVM Mathematically



$$\bar{\omega} \cdot \bar{x}_+ + b \geq 1$$

$$\bar{\omega} \cdot \bar{x}_- + b \geq -1$$

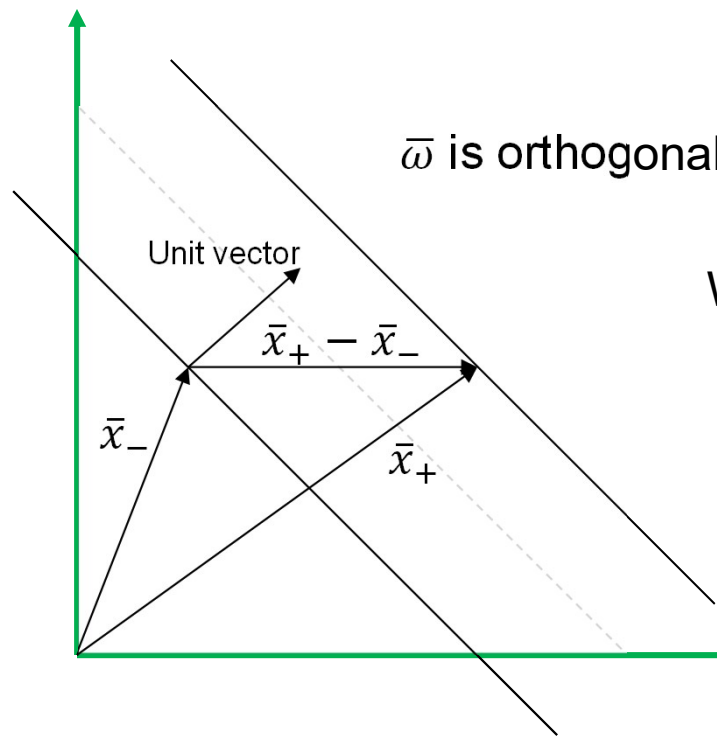
Now,  $y_i$  such that  $y_i = +1$  for + samples  
 $y_i = -1$  for - samples

$$y_i(\bar{\omega} \cdot \bar{x}_i + b) \geq 1$$

$$y_i(\bar{\omega} \cdot \bar{x}_i + b) - 1 \geq 0$$

$$\text{if } y_i(\bar{\omega} \cdot \bar{x}_i + b) - 1 = 0$$

Data  $(\bar{x}_i)$  is on the DB



$\bar{\omega}$  is orthogonal to the decision boundary

$$\begin{aligned} \text{Width} &= (\bar{x}_+ - \bar{x}_-) \cdot \frac{\bar{\omega}}{\|\omega\|} = \frac{\bar{\omega}x_+ - (\bar{\omega}x_-)}{\|\omega\|} \\ &= \frac{(1-b) - (-1-b)}{\|\omega\|} = \frac{1-b+1+b}{\|\omega\|} = \frac{2}{\|\omega\|} \end{aligned}$$

$$\text{Since } \bar{\omega} \cdot \bar{x}_+ + b \geq 1$$

$$\bar{\omega} \cdot \bar{x}_- + b \geq -1$$

# Linear SVM Mathematically

- Goal: 1) **Correctly classify all training data**

$$wx_i + b \geq 1 \quad \text{if } y_i = +1$$

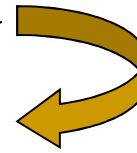
$$wx_i + b \leq -1 \quad \text{if } y_i = -1$$

$$y_i(wx_i + b) \geq 1 \quad \text{for all } i$$

- 2) **Maximize the Margin**

same as minimize

$$M = \frac{2}{|w|}$$
$$\frac{1}{2} w^t w$$



- We can formulate a Quadratic Optimization Problem and solve for  $w$  and  $b$

- Minimize  $\Phi(w) = \frac{1}{2} w^t w$
- subject to  $y_i(wx_i + b) \geq 1 \quad \forall i$



---

# Solving the Optimization Problem

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$  is minimized;

and for all  $\{(\mathbf{x}_i, y_i)\}: y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

- This is a optimization problem for *convex quadratic* objective subject to *linear constraints*.
  - Quadratic optimization problems are a well-known class of mathematical programming problems, and can be solved using Quadratic Programming (QP)
  - The solution involves constructing a *dual problem* where a *Lagrange multiplier*  $\alpha_i$  is associated with every constraint in the primary problem:
-

---

# Solving the Optimization Problem

- Lagrange duality to get the optimization problem's dual form
    - Allow us to use kernels to get optimal margin classifier to work efficiently in very high dimensional spaces
    - Allow us to derive a very efficient algorithm for solving the above optimization problems
-

---

# Lagrangian Duality

- The primal problem  $\min_w f(w)$   
 $\text{s. t. } g_i(w) \leq 0, i = 1, 2, \dots, k$   
 $h_i(w) = 0, i = 1, 2, \dots, l$

- The generalized Lagrangian

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

where the  $\alpha_i$ 's ( $\alpha_i \geq 0$ ) and  $\beta_i$ 's are called Lagrange multipliers

---

---

# Lagrangian Duality

- Lemma:

$$\begin{aligned} & \max_{\alpha, \beta, \alpha_{i \geq 0}} L(w, \alpha, \beta) \\ &= \begin{cases} f(w) & \text{if } w \text{ satisfies the primal constraints} \\ \infty & \text{otherwise} \end{cases} \end{aligned}$$

- So, the primal can be re-written

$$\min_w \max_{\alpha, \beta, \alpha_{i \geq 0}} L(w, \alpha, \beta)$$

---

---

# Lagrangian Duality: Dual Formation

- The primal problem:

$$p^* = \min_w \max_{\alpha, \beta, \alpha_{i \geq 0}} L(w, \alpha, \beta)$$

- The dual problem:

$$d^* = \max_{\alpha, \beta, \alpha_{i \geq 0}} \min_w L(w, \alpha, \beta)$$

- Theorem (weak duality):

$$\max_{\alpha, \beta, \alpha_{i \geq 0}} \min_w L(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta, \alpha_{i \geq 0}} L(w, \alpha, \beta)$$

i.e.  $d^* \leq p^*$

- Theorem (strong duality): iff there exists a saddle point of  $L(w, \alpha, \beta)$ , we have  $d^* = p^*$
-

---

## The KKT Conditions

- iff there exists some saddle point of  $L(w, \alpha, \beta)$ , then it satisfies “Karush-Kuhn-Tucker” (KKT) conditions:

$$\frac{\partial}{\partial w_i} L(w, \alpha, \beta) = 0, \quad i = 1 \dots k$$

$$\frac{\partial}{\partial \beta_i} L(w, \alpha, \beta) = 0, \quad i = 1 \dots l$$

$$\alpha_i g_i(w) = 0, i = 1, \dots, m$$

$$g_i(w) \leq 0, i = 1, \dots, m$$

$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

- **Theorem:** If  $w^*, \alpha^*$  and  $b^*$  satisfy the KKT conditions, then it is also a solution to the primal and the dual problem
-

---

# Solving the Optimization Problem

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$  is minimized;

and for all  $\{(\mathbf{x}_i, y_i)\}: y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

- This is a optimization problem for *convex quadratic* objective subject to *linear constraints*.
  - Quadratic optimization problems are a well-known class of mathematical programming problems, and can be solved using Quadratic Programming (QP)
  - The solution involves constructing a *dual problem* where a *Lagrange multiplier*  $\alpha_i$  is associated with every constraint in the primary problem:
-

---

# Lagrangian Duality for SVM

- The primal problem  $\min_w f(w)$   
s. t.  $g_i(w) \leq 0, i = 1, 2, \dots, k$   
 $h_i(w) = 0, i = 1, 2, \dots, l$

- The generalized Lagrangian

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

where the  $\alpha_i$ 's ( $\alpha_i \geq 0$ ) and  $\beta_i$ 's are called Lagrange multipliers

We don't have the equality constraints, so we only deal with  $\alpha_i$ s

---



---

## The KKT Conditions

- iff there exists some saddle point of  $L(w, \alpha, \beta)$ , then it satisfies “Karush-Kuhn-Tucker” (KKT) conditions:

$$\frac{\partial}{\partial w_i} L(w, \alpha, \beta) = 0, \quad i = 1 \dots k$$

$$\frac{\partial}{\partial \beta_i} L(w, \alpha, \beta) = 0, \quad i = 1 \dots l$$

$$\alpha_i g_i(w) = 0, i = 1, \dots, m$$

$$g_i(w) \leq 0, i = 1, \dots, m$$

$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

- **Theorem:** If  $w^*, \alpha^*$  and  $b^*$  satisfy the KKT conditions, then it is also a solution to the primal and the dual problem
-

---

# Support Vectors

$$\alpha_i g_i(w) = 0, i = 1, \dots, m$$

$$g_i(w) \leq 0, i = 1, \dots, m$$

- If  $\alpha_i > 0$ , then  $g_i(w) = 0$
  - Only a few  $\alpha_i$ 's are non-zero.
  - The training data points having non-zero  $\alpha_i$  are called support vectors
-

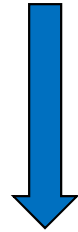
# Solving the SVM Optimization Problem

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$  is minimized;  
and for all  $\{(\mathbf{x}_i, y_i)\}: y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

*Quadratic programming  
with linear constraints.*

Lagrangian Function



$$\begin{aligned} \text{minimize } L_p(\mathbf{w}, b, \alpha_i) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \\ \text{s.t. } \alpha_i &\geq 0 \end{aligned}$$

Where  $p$  denotes *primal*

# Solving the SVM Optimization Problem

$$\begin{aligned} \text{minimize } L_p(\mathbf{w}, b, \alpha_i) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \\ \text{s.t. } \alpha_i &\geq 0 \end{aligned}$$

*Minimizing  $L_p$  w.r.t.  
 $\mathbf{w}$  and  $b$  with fixed  $\alpha$*

$$\begin{aligned} \frac{\partial L_p}{\partial \mathbf{w}} = 0 &\longrightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial L_p}{\partial b} = 0 &\longrightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

# Solving the SVM Optimization Problem

$$L_p(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) - b \sum_{i=1}^m \alpha_i y_i$$

$$L_p(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

$$\text{as } \sum_{i=1}^n \alpha_i y_i = 0$$

---

## The Dual Problem

- Now we have the following dual optimization problem

$$\begin{aligned} \max_{\alpha} J(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s. t. } \alpha_i &\geq 0, i = 1, \dots, m \\ \sum_{i=1}^m \alpha_i y_i &= 0. \end{aligned}$$

- This is a **quadratic programming** problem
    - A global maximum of  $\alpha_i$  can always be found
-

---

# Support Vector Machine

- Once we have Lagrange multipliers  $\{\alpha_j\}$ , we can reconstruct the parameter vector  $\mathbf{w}$  as a weighted combination of the training examples

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \qquad \mathbf{w} = \sum_{i \in SV} \alpha_i y_i \mathbf{x}_i$$

- For testing with a new data  $\mathbf{z}$

Compute  $\mathbf{w}^T \mathbf{z} + b = \sum_{i \in SV} \alpha_i y_i (\mathbf{x}_i^T \mathbf{z}) + b$

and classify  $\mathbf{z}$  as class 1 if the sum is **positive**, and as class 2 **otherwise** (Note:  $\mathbf{w}$  need not to be formed explicitly)

---

---

# Solving the Optimization Problem

- The discriminant function is:

$$g(x) = \mathbf{w}^T \mathbf{x} + b = \sum_{i \in SV} \alpha_i \mathbf{x}_i^T \mathbf{x} + b$$

where  $\mathbf{x}_i^T \mathbf{x}$  = dot product of  $\mathbf{x}_i^T$  and  $\mathbf{x}$

- It relies on a dot product between a test point and the support vectors  $\mathbf{x}_i$
  - Solving the optimization problem involved computing the dot product of  $\mathbf{x}_i^T \mathbf{x}$  between all pairs of training points. It is scalar value.
  - The optimal  $\mathbf{w}$  is a linear combination of a small number of data points
-



---

# The Dual Problem

- Now we have the following dual optimization problem

$$\begin{aligned} \max_{\alpha} J(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s. t. } \alpha_i &\geq 0, i = 1, \dots, m \\ \sum_{i=1}^m \alpha_i y_i &= 0. \end{aligned}$$

- This is a **quadratic programming** problem
    - A global maximum of  $\alpha_i$  can always be found
-

---

to continue...

---

---

# Linear SVM Formulation

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$  is minimized;

and for all  $\{(\mathbf{x}_i, y_i)\}$ :  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

(1)  $\sum \alpha_i y_i = 0$

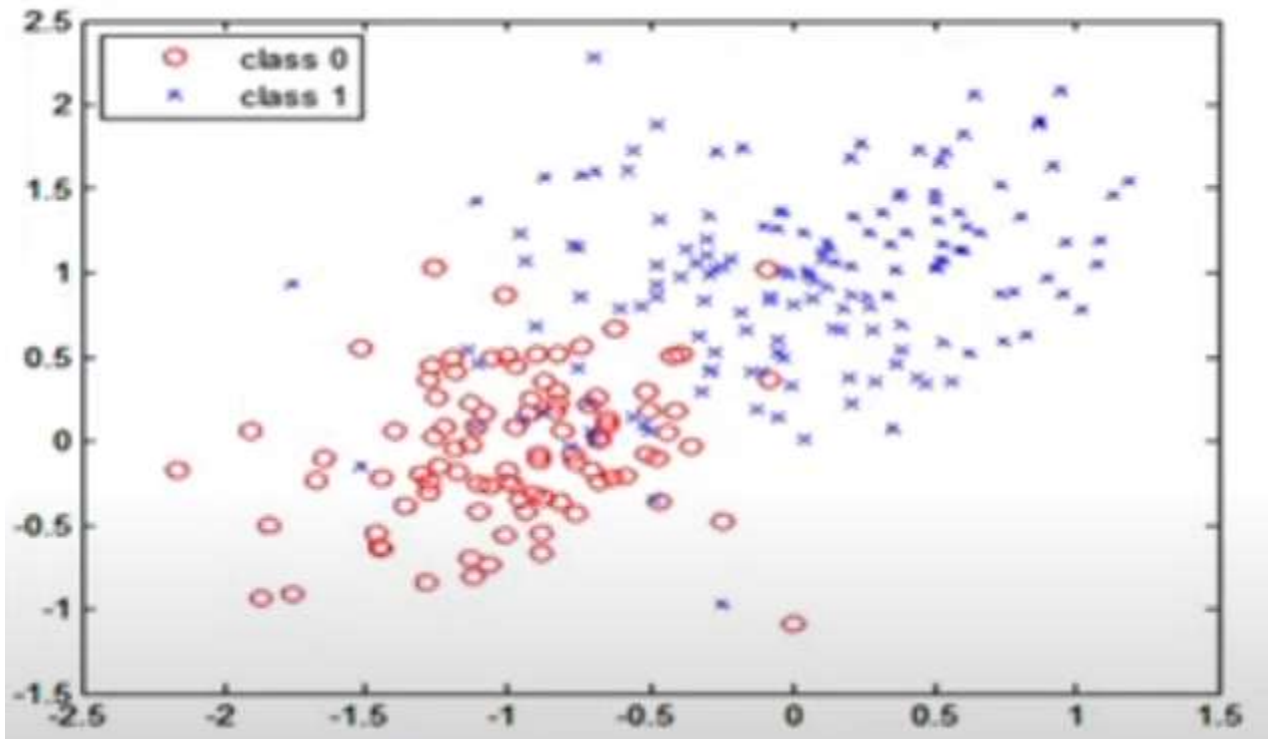
(2)  $\alpha_i \geq 0$  for all  $\alpha_i$

**Assumption:** Training examples are linearly separable

---

---

# Limitations of Linear SVM



- What if data is not linearly separable?

Or

- Noisy data points?

Extend the definition of maximum margin to allow the non-separating planes.

---

---

# The Optimization Problem Solution

- The solution has the form:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_k - \mathbf{w}^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ such that } \alpha_k \neq 0$$

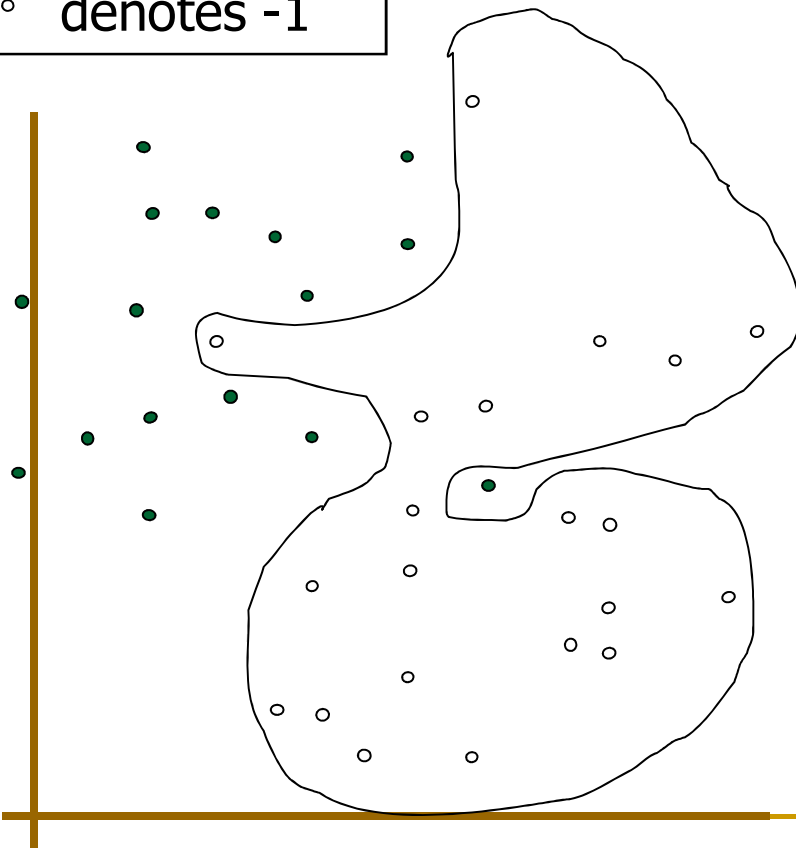
- Each non-zero  $\alpha_i$  indicates that corresponding  $\mathbf{x}_i$  is a support vector.
- Then the classifying function will have the form:

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

- Notice that it relies on an *inner product* between the test point  $\mathbf{x}$  and the support vectors  $\mathbf{x}_i$  – we will return to this later.
  - Also keep in mind that solving the optimization problem involved computing the inner products  $\mathbf{x}_i^T \mathbf{x}_j$  between all pairs of training points.
-

# Dataset with noise

- denotes +1
- denotes -1

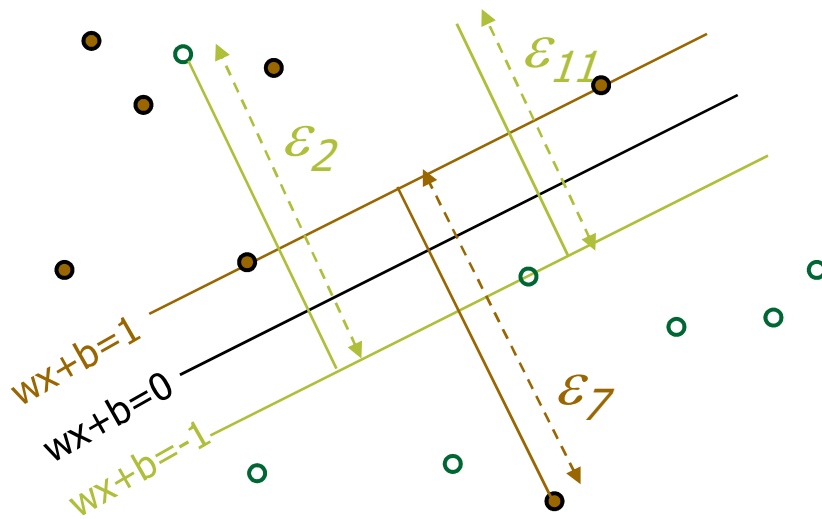


- **Hard Margin:** So far we require all data points be classified correctly
  - No training error
- **What if the training set is noisy?**
  - **Solution 1:** use very powerful kernels

**OVERFITTING!**

# Soft Margin Classification

**Slack variables  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.**



What should our quadratic optimization criterion be?

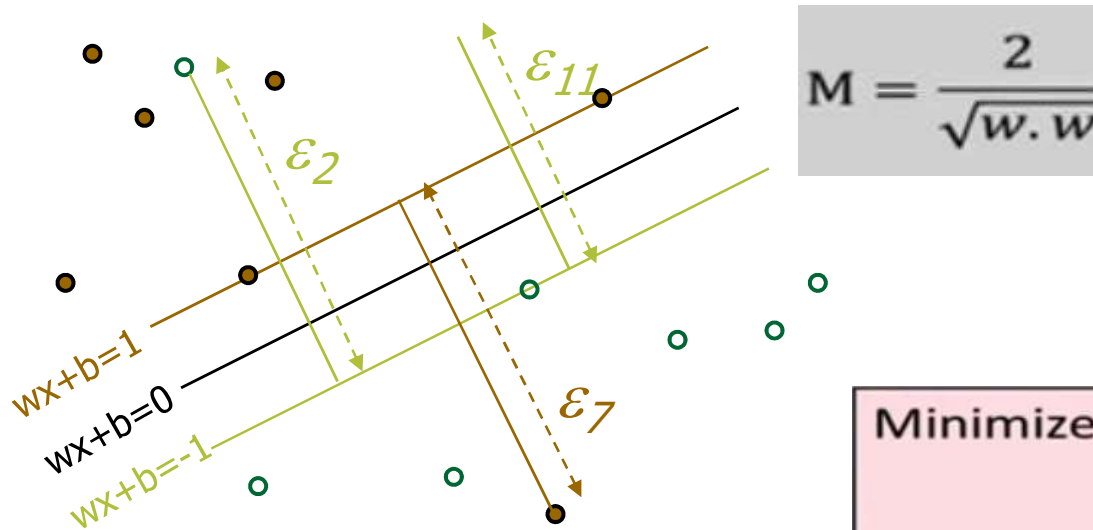
Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \cdot \# \text{ Training errors}$$

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \xi_k$$

The problem is that it is **no longer a quadratic optimization problem** and we **can't use QP solver** to solve it.

# Soft Margin Classification



- $C$  controls the relative importance of maximizing the margins and fitting the training data.
- Controls overfitting

Minimize

$$w \cdot w + C \sum_{k=1}^m \xi_k$$

$m$  constraints

$$\left[ \begin{array}{l} w \cdot x_k + b \geq 1 - \xi_k \text{ if } y_k = 1 \\ w \cdot x_k + b \leq -1 + \xi_k \text{ if } y_k = -1 \end{array} \right]$$

$\equiv$

$$y_k(w \cdot x_k + b) \geq 1 - \xi_k, \quad k=1, \dots, m$$

$$\xi_k \geq 0, \quad k=1, \dots, m$$

Activate Windows



# Lagrangian

$$\begin{aligned} L(w, b, \xi, \alpha, \beta) \\ = \frac{1}{2} w \cdot w + C \sum_{i=1}^m \xi_i \\ + \sum_{i=1}^m \alpha_i [y_i (x \cdot w + b) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i \end{aligned}$$

$\alpha_i$ 's and  $\beta_i$ 's are Lagrange multipliers ( $\geq 0$ ).

# Dual Formulation

Find  $\alpha_1, \alpha_2, \dots, \alpha_m$  s.t.

$$\max_{\alpha} J(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

**Linear SVM**

$$\begin{aligned} \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

**Noise Accounted**

$$\begin{aligned} \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

# Soft Margin Classifier

- $x_i$  with non-zero  $\alpha_i$  will be support vectors.
- Solution to the dual problem is:

$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$b = y_k(1 - \xi_k) - \sum_{i=1}^m \alpha_i y_i x_i x_k$$

for any  $k$  s.t.  $\alpha_k > 0$

For classification,

$$f(x) = \sum_{i=1}^m \alpha_i y_i x_i \cdot x + b$$

(no need to compute  $w$  explicitly)

---

# Hard Margin v.s. Soft Margin

- **The old formulation:**

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

- **The new formulation incorporating slack variables:**

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0 \text{ for all } i$$

- **Parameter  $C$  can be viewed as a way to control overfitting.**

---

# Linear SVMs: Overview

- The classifier is a *separating hyperplane*.
- Most “important” training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points  $\mathbf{x}_i$  are support vectors with non-zero Lagrangian multipliers  $\alpha_i$ .
- Both in the dual formulation of the problem and in the solution training points appear only inside dot products:

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

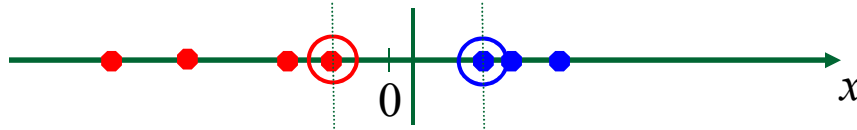
(1)  $\sum \alpha_i y_i = 0$

(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$

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

# Non-linear SVMs

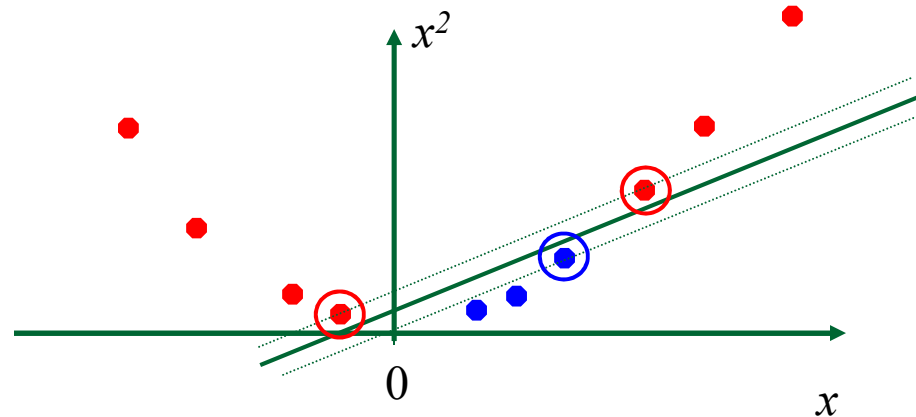
- Datasets that are linearly separable with some noise work out great:



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

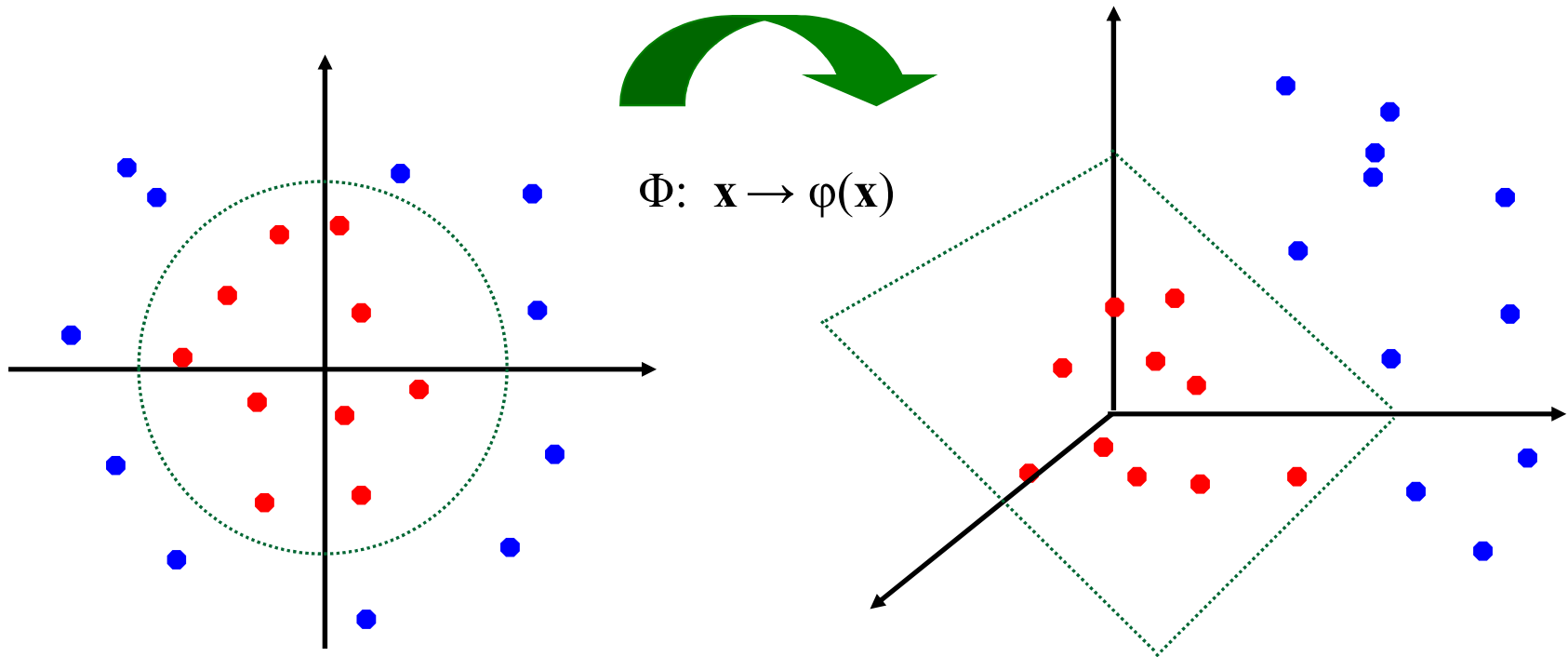


- How about... mapping data to a higher-dimensional space:



# Non-linear SVMs: Feature spaces

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



# Non- Linear SVM

## ■ Computational cost

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

(1)  $\sum \alpha_i y_i = 0$

(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$   $O(d^2)$

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \varphi(\mathbf{x}_i^T) \varphi(\mathbf{x}_j)$  is maximized and

(1)  $\sum \alpha_i y_i = 0$

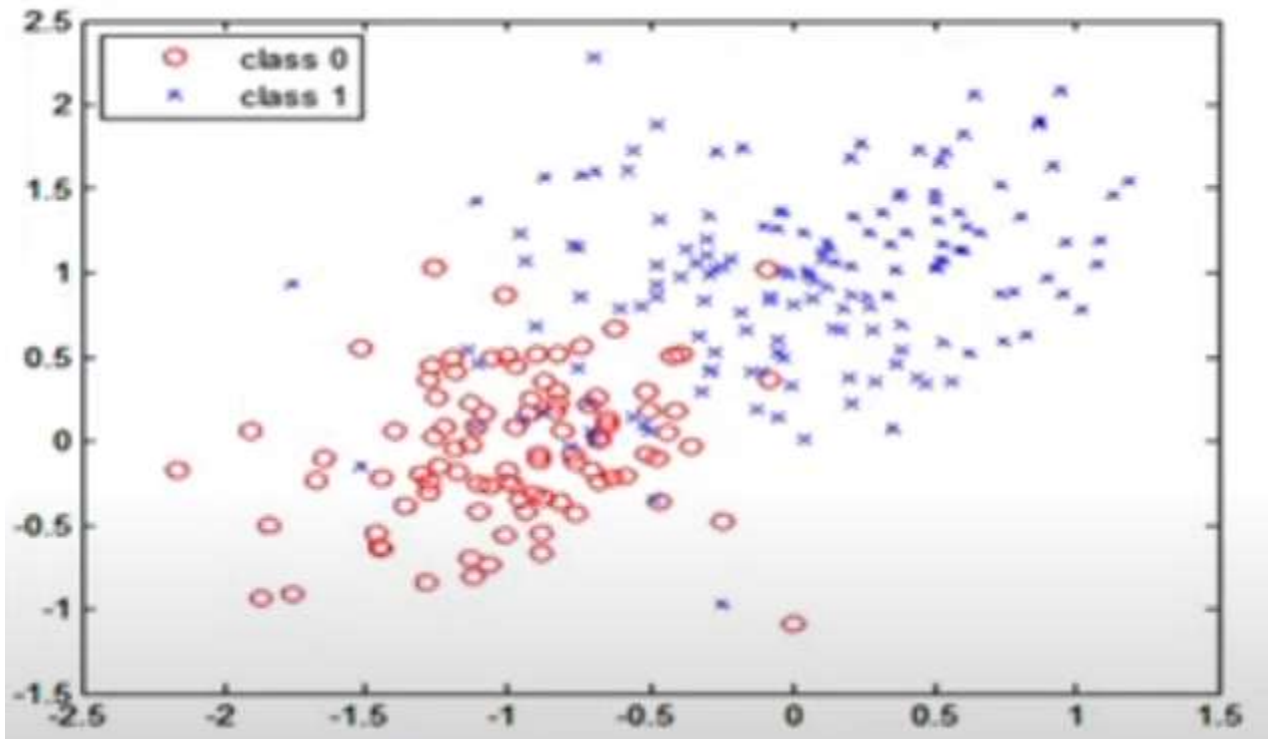
(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$   $O(D^2)$

If  $D \gg d$ , computational cost will increase substantially



---

# Limitations of Linear SVM



- What if data is not linearly separable?

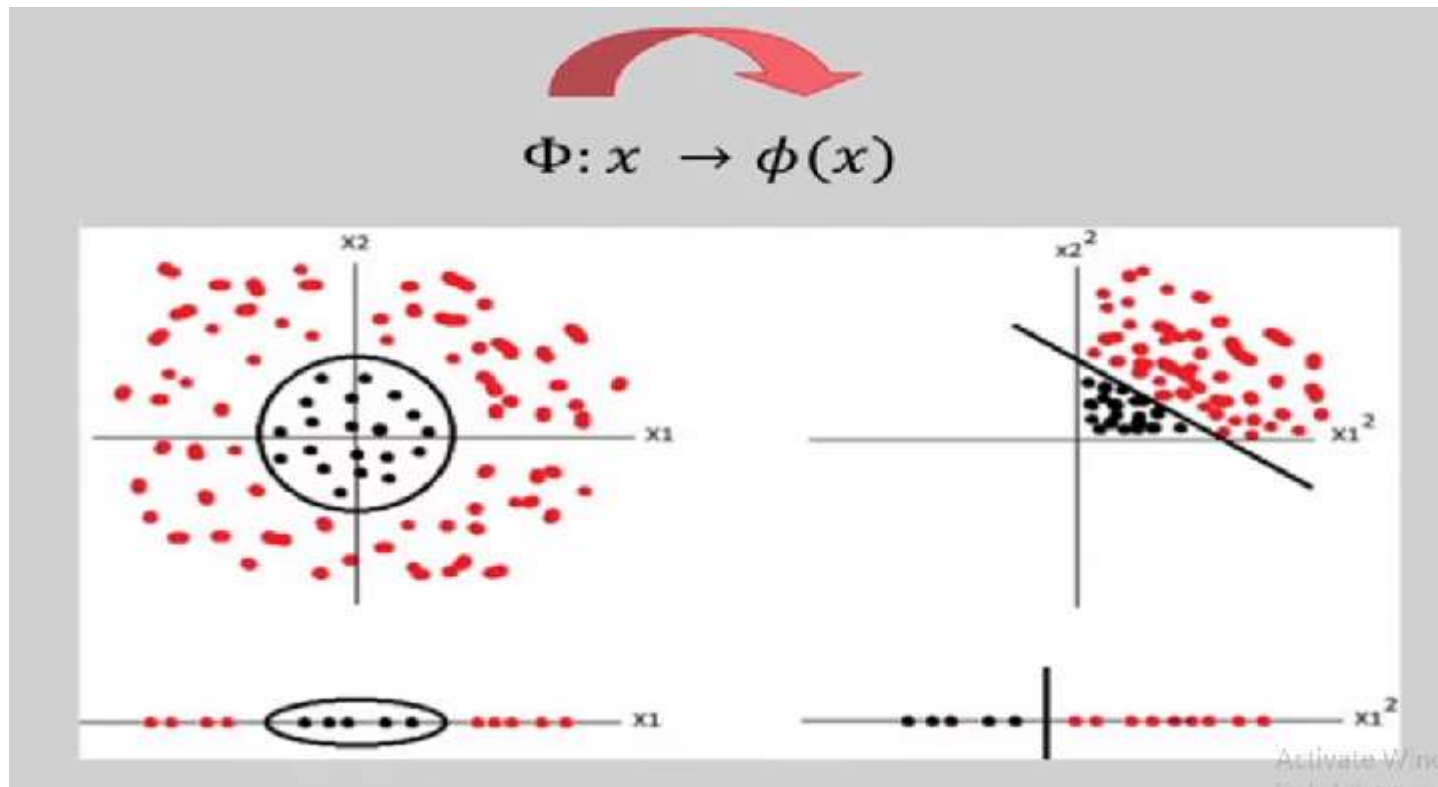
Or

- Noisy data points?

Extend the definition of maximum margin to allow the non-separating planes.

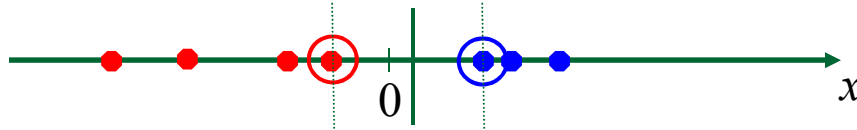
---

# Non Linear SVM: Feature Space

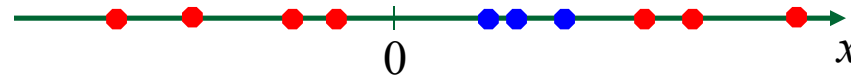


# Non-linear SVMs

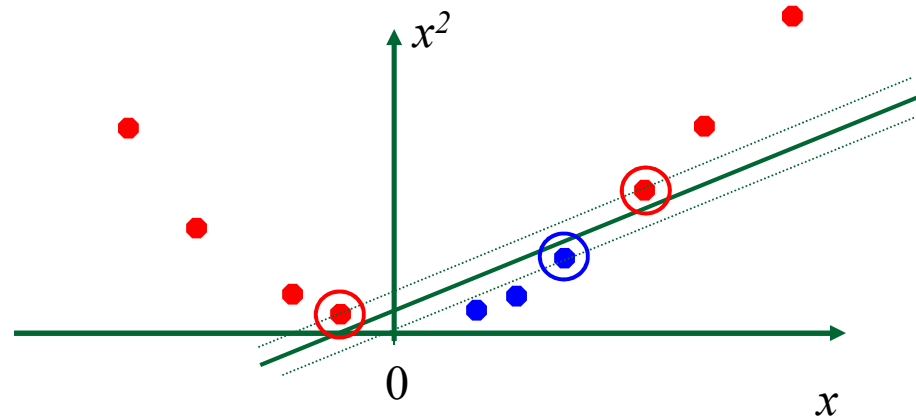
- Datasets that are linearly separable with some noise work out great:



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

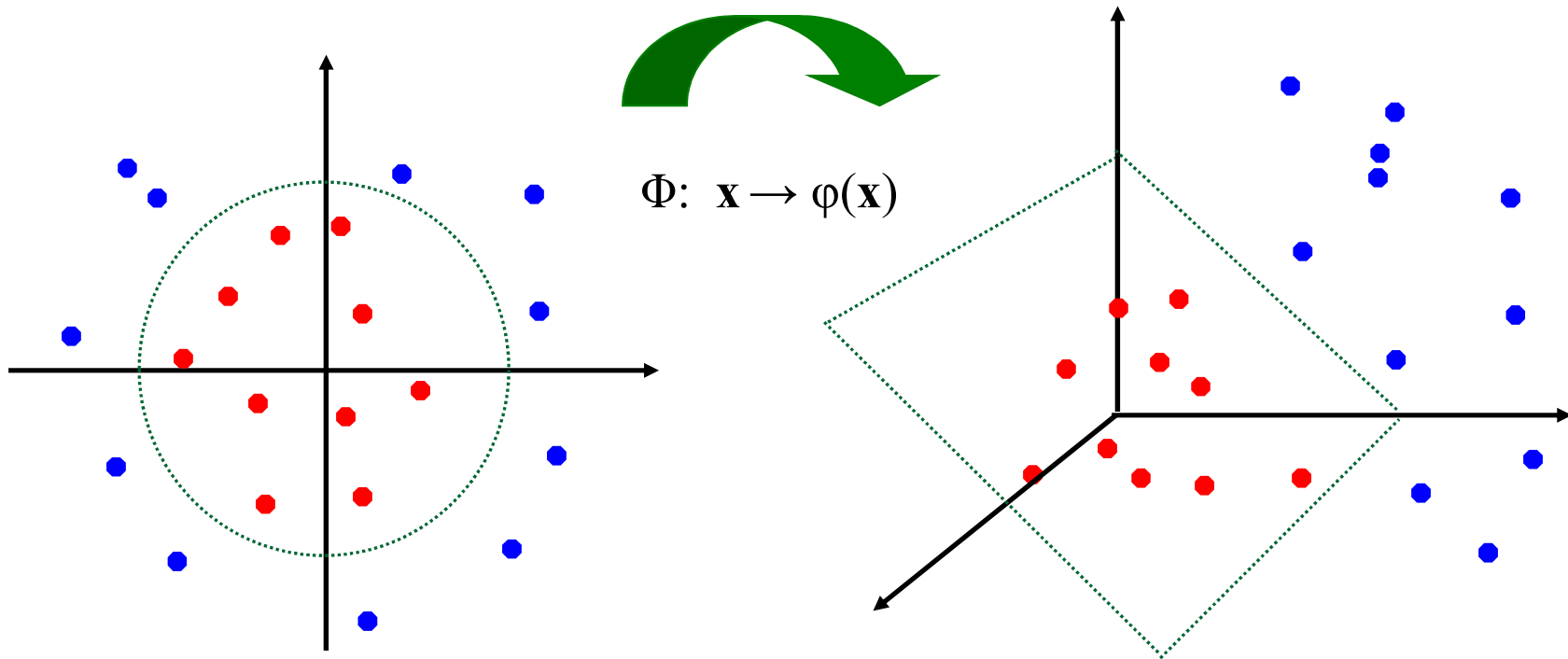


- How about... mapping data to a higher-dimensional space:



# Non-linear SVMs: Feature spaces

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



# The “Kernel Trick”

- The linear classifier relies on dot product between vectors  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- If every data point is mapped into high-dimensional space via some transformation  $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$ , the dot product becomes:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.
- Example:

2-dimensional vectors  $\mathbf{x} = [x_1 \ x_2]$ ; let  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ ,

Need to show that  $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ :

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2, \\ &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\ &= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \quad \text{where } \phi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2] \end{aligned}$$

# What Functions are Kernels?

- For some functions  $K(\mathbf{x}_i, \mathbf{x}_j)$  checking that

$K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^T \boldsymbol{\phi}(\mathbf{x}_j)$  can be cumbersome.

- Mercer's theorem:

*Every semi-positive definite symmetric function is a kernel*

- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

$K =$

$K(\mathbf{x}_1, \mathbf{x}_1)$	$K(\mathbf{x}_1, \mathbf{x}_2)$	$K(\mathbf{x}_1, \mathbf{x}_3)$	$\dots$	$K(\mathbf{x}_1, \mathbf{x}_N)$
$K(\mathbf{x}_2, \mathbf{x}_1)$	$K(\mathbf{x}_2, \mathbf{x}_2)$	$K(\mathbf{x}_2, \mathbf{x}_3)$		$K(\mathbf{x}_2, \mathbf{x}_N)$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$K(\mathbf{x}_N, \mathbf{x}_1)$	$K(\mathbf{x}_N, \mathbf{x}_2)$	$K(\mathbf{x}_N, \mathbf{x}_3)$	$\dots$	$K(\mathbf{x}_N, \mathbf{x}_N)$

---

# Examples of Kernel Functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial of power  $p$ :  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$

- Gaussian (radial-basis function network):

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

- Sigmoid:  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T \mathbf{x}_j + \beta_1)$
-

---

# Kernel Functions

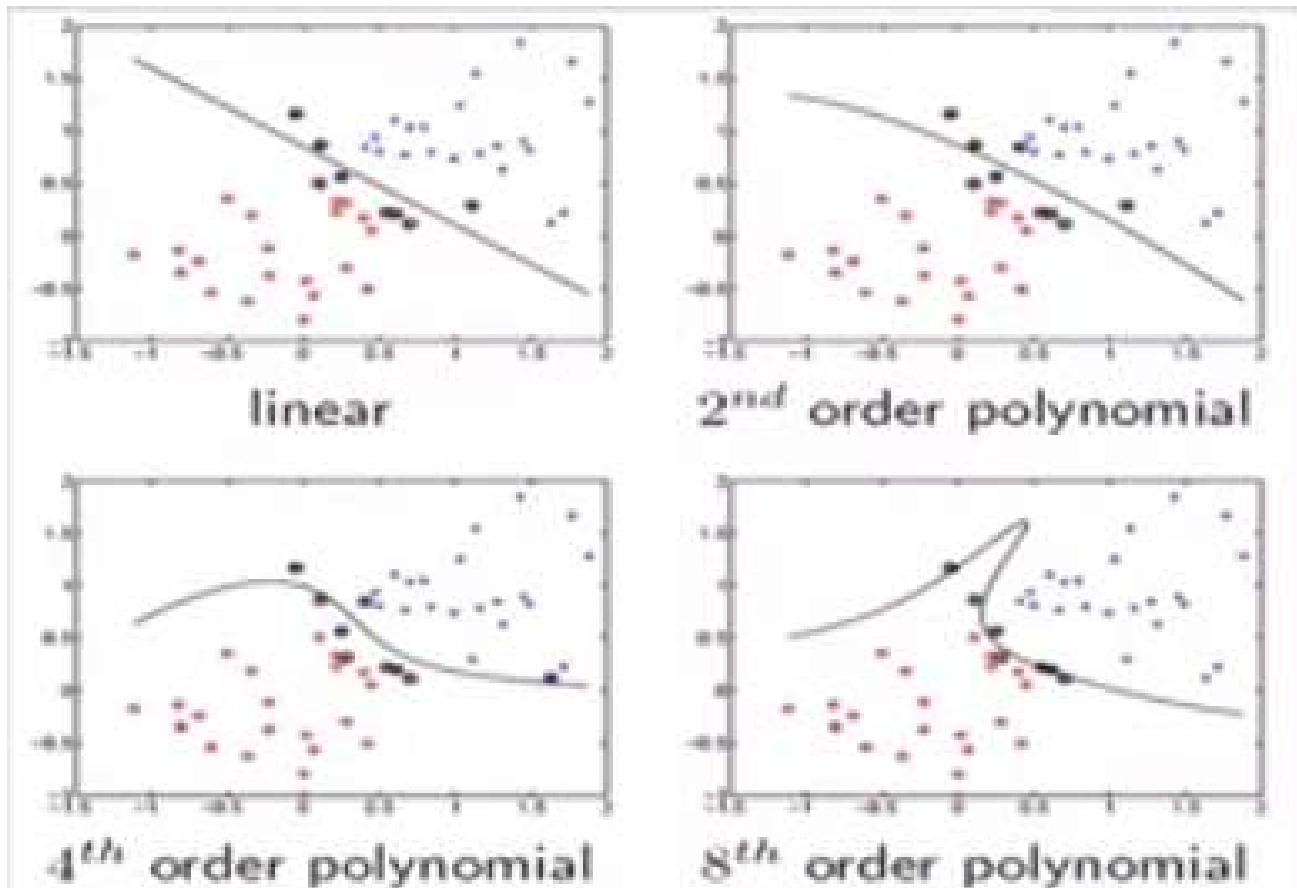
- Kernel function can be thought of a similarity measure between the input objects.
- Not all similarity measure can be used as Kernel function
- Mercer's condition state that any positive semi-definite kernel  $K(X,Y)$  i.e.

$$\sum_{i,j} K(x_i x_j) c_i c_j \geq 0$$

- Can be expressed as a dot product in a high dimensional space.
-

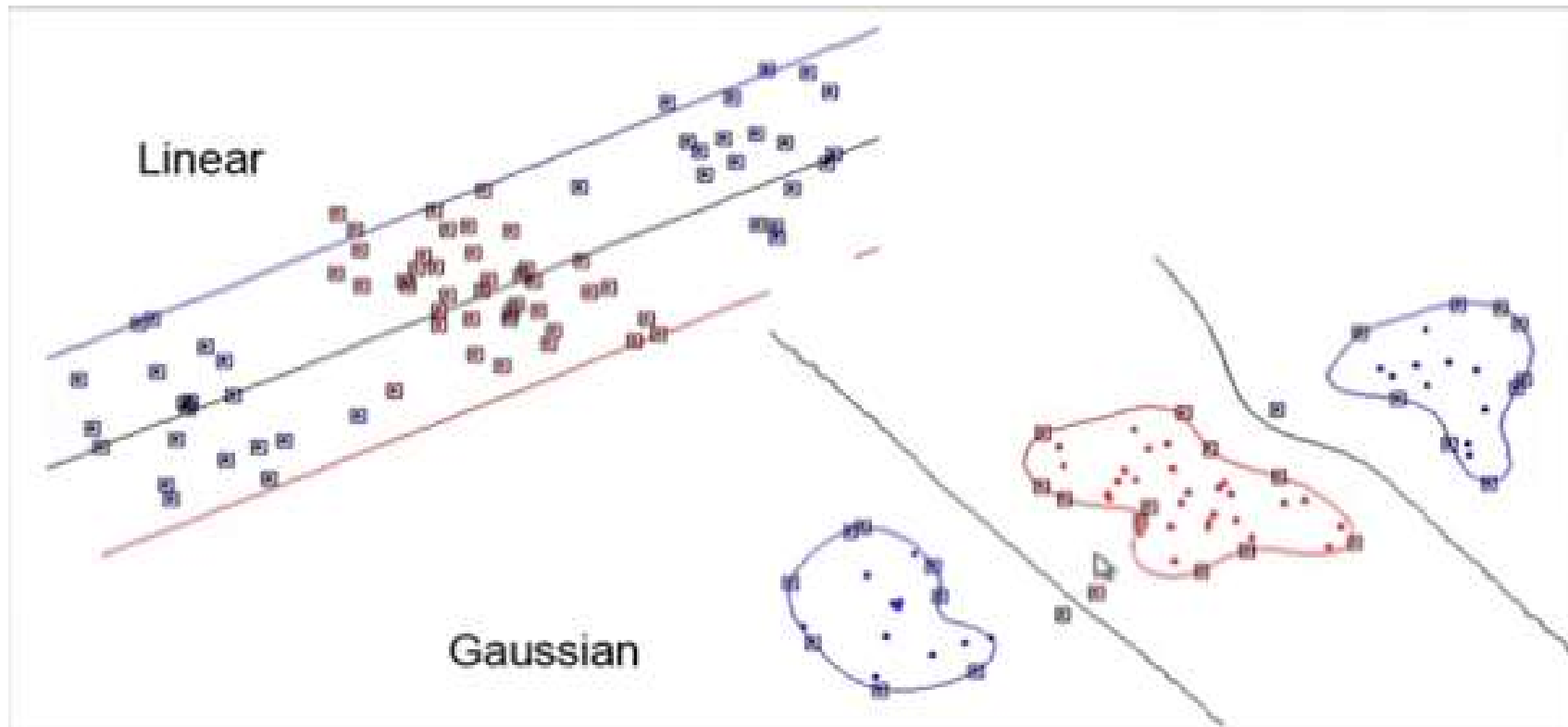


# SVM Examples



# Non Linear SVM

## Gaussian Kernels



---

# Non-linear SVMs Mathematically

- **Dual problem formulation:**

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$  is maximized and

(1)  $\sum \alpha_i y_i = 0$

(2)  $\alpha_i \geq 0$  for all  $\alpha_i$

- **The solution is:**

$$f(\mathbf{x}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$$

- **Optimization techniques for finding  $\alpha_i$ 's remain the same!**

---

---

# Nonlinear SVM - Overview

- SVM locates a separating hyperplane in the feature space and classify points in that space
  - It does not need to represent the space explicitly, simply by defining a kernel function
  - The kernel function plays the role of the dot product in the feature space.
-

---

# SVM: Summary

- Flexibility in choosing a similarity function
  - Sparseness of solution when dealing with large data sets
    - only support vectors are used to specify the separating hyperplane
  - Ability to handle large feature spaces
    - complexity does not depend on the dimensionality of the feature space
  - Overfitting can be controlled by soft margin approach
  - Nice math property: a simple convex optimization problem which is guaranteed to converge to a single global solution
  - Feature Selection
-

---

# SVM Applications

- SVM has been used successfully in many real-world problems
    - text (and hypertext) categorization
    - image classification
    - bioinformatics (Protein classification, Cancer classification)
    - hand-written character recognition
-

---

# Weakness of SVM

- It is sensitive to noise
    - A relatively small number of mislabeled examples can dramatically decrease the performance
  - It only considers two classes
    - how to do multi-class classification with SVM?
    - Answer:
      - 1) with output arity  $m$ , learn  $m$  SVM's
        - SVM 1 learns “Output==1” vs “Output != 1”
        - SVM 2 learns “Output==2” vs “Output != 2”
        - :
        - SVM  $m$  learns “Output== $m$ ” vs “Output !=  $m$ ”
      - 2) To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.
-

---

to continue...

---