

# Machine Learning

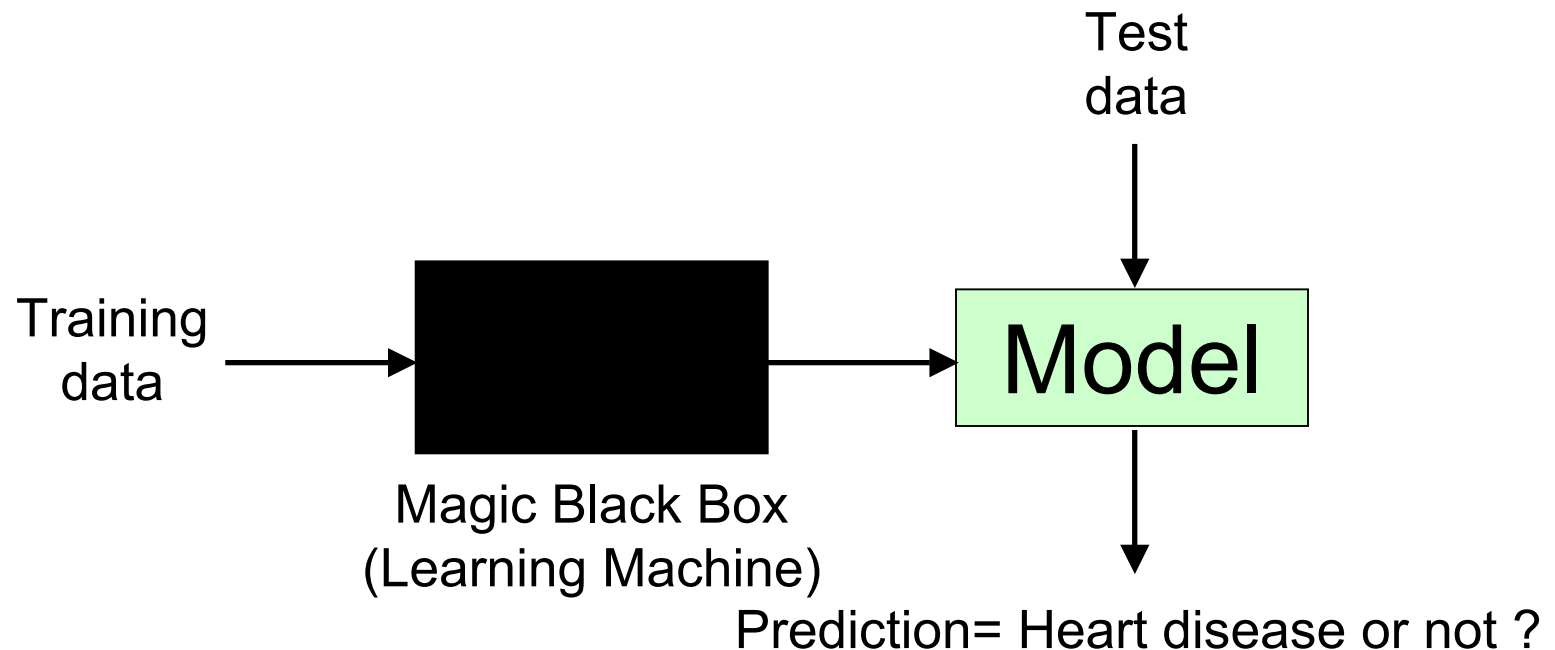
“... a computer program that can learn from **experience** with respect to some class of **tasks** and **performance** measure ...” (Mitchell, 1997)

“... the capacity of a computer to learn from **experience**, i.e. to **modify** its processing on the basis of newly acquired **information** ...” (Oxford Dictionary, 1989)

Examples are:

- Support Vector Machines
- Artificial Neural Networks
- Fuzzy Logic
- Evolutionary Computation

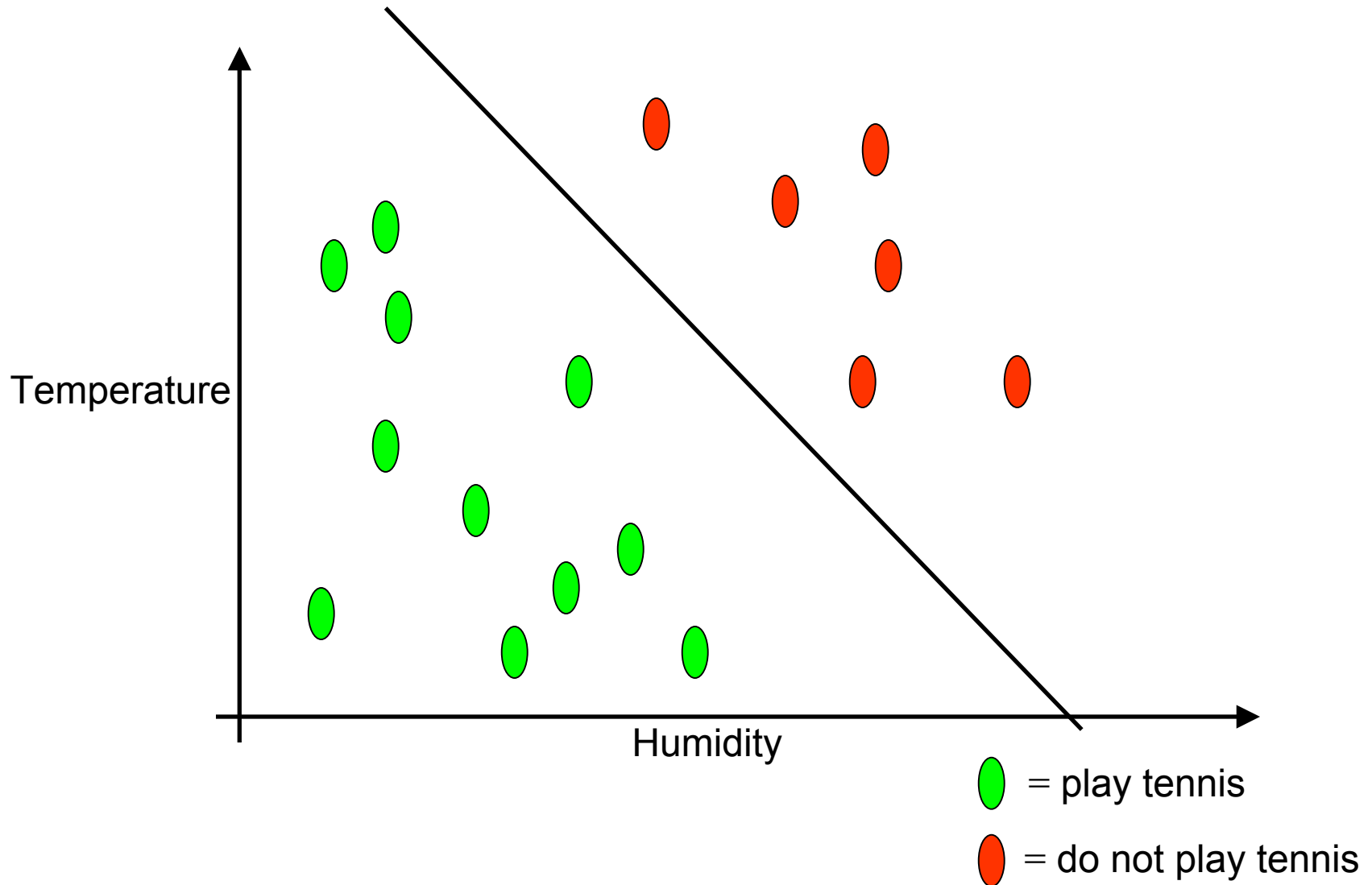
# Black Box View of Machine Learning



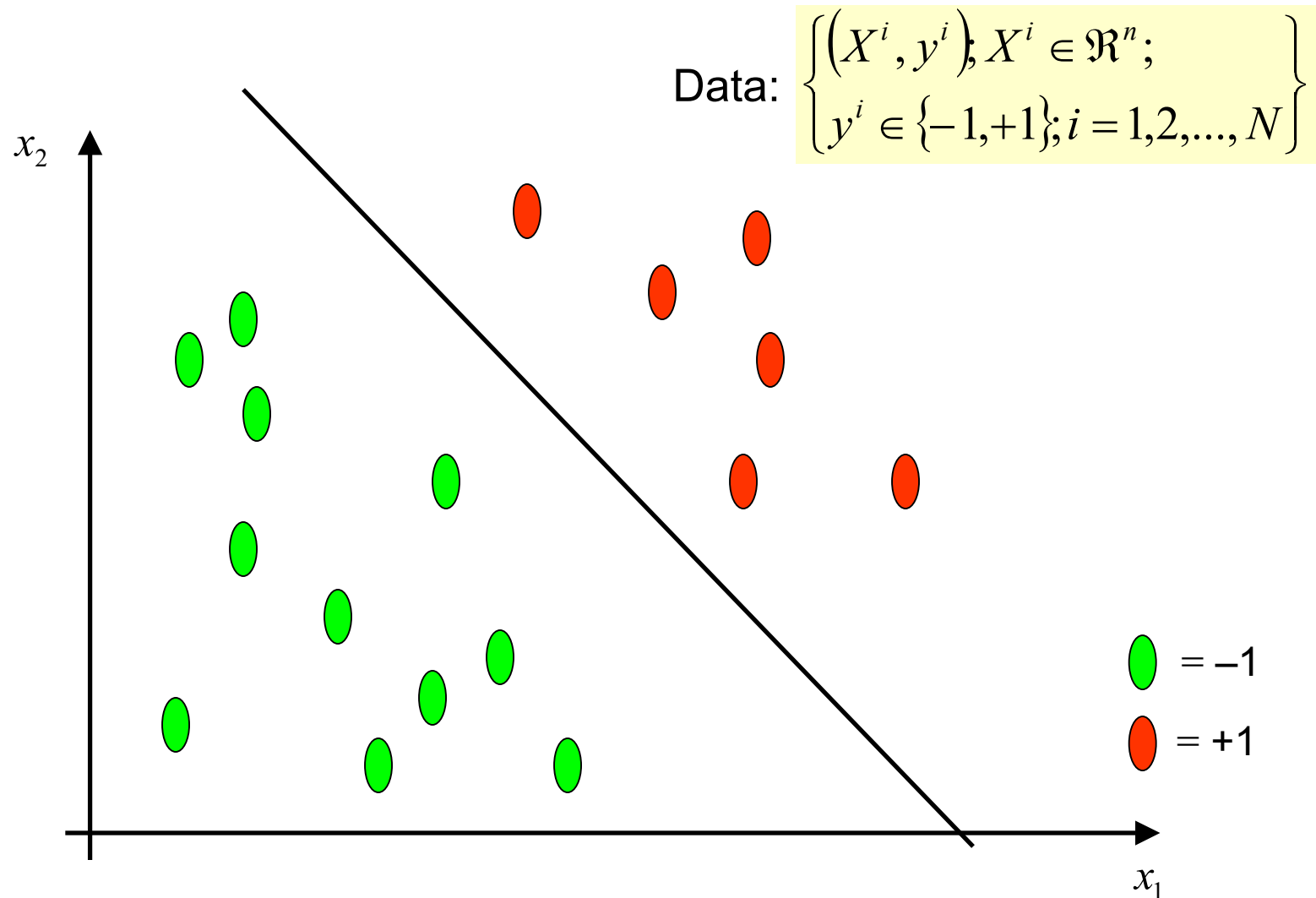
Training data: Expression pattern of some disease and healthy person

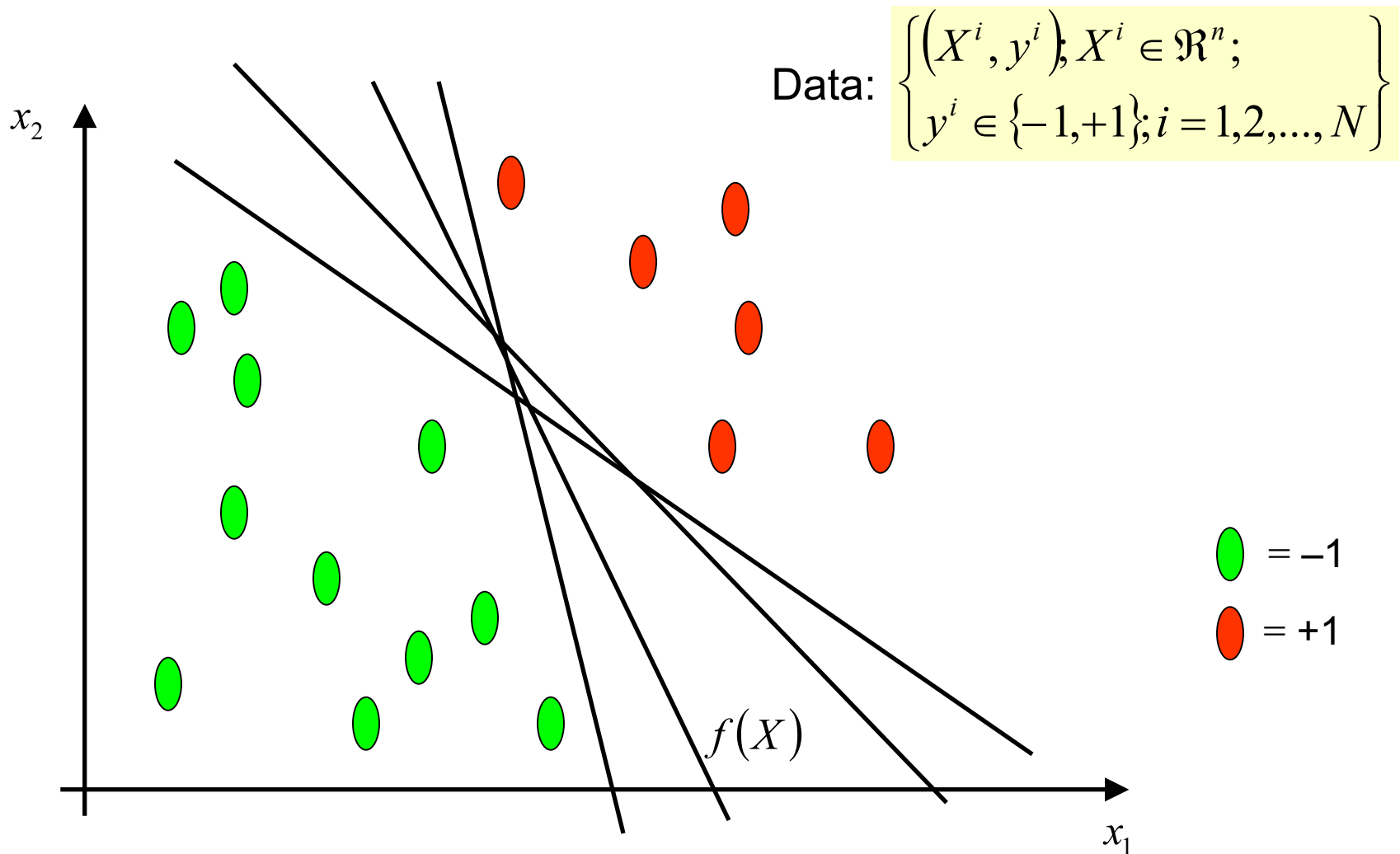
Model: Can distinguish between healthy and sick person. Can be used for prediction

## Tennis Example



# Linear Support Vector Machine





All hyperplanes in  $\mathbb{R}^n$  are parameterized by  $w$  and a constant  $b$

Objective: To find a hyperplane  $f(X) = \text{sign}(\langle w, X \rangle + b)$  that correctly classify the data

The optimal hyperplane  $H$  is such that

$$\langle w.X^i \rangle + b \geq +1 \quad \text{when } y^i = +1$$

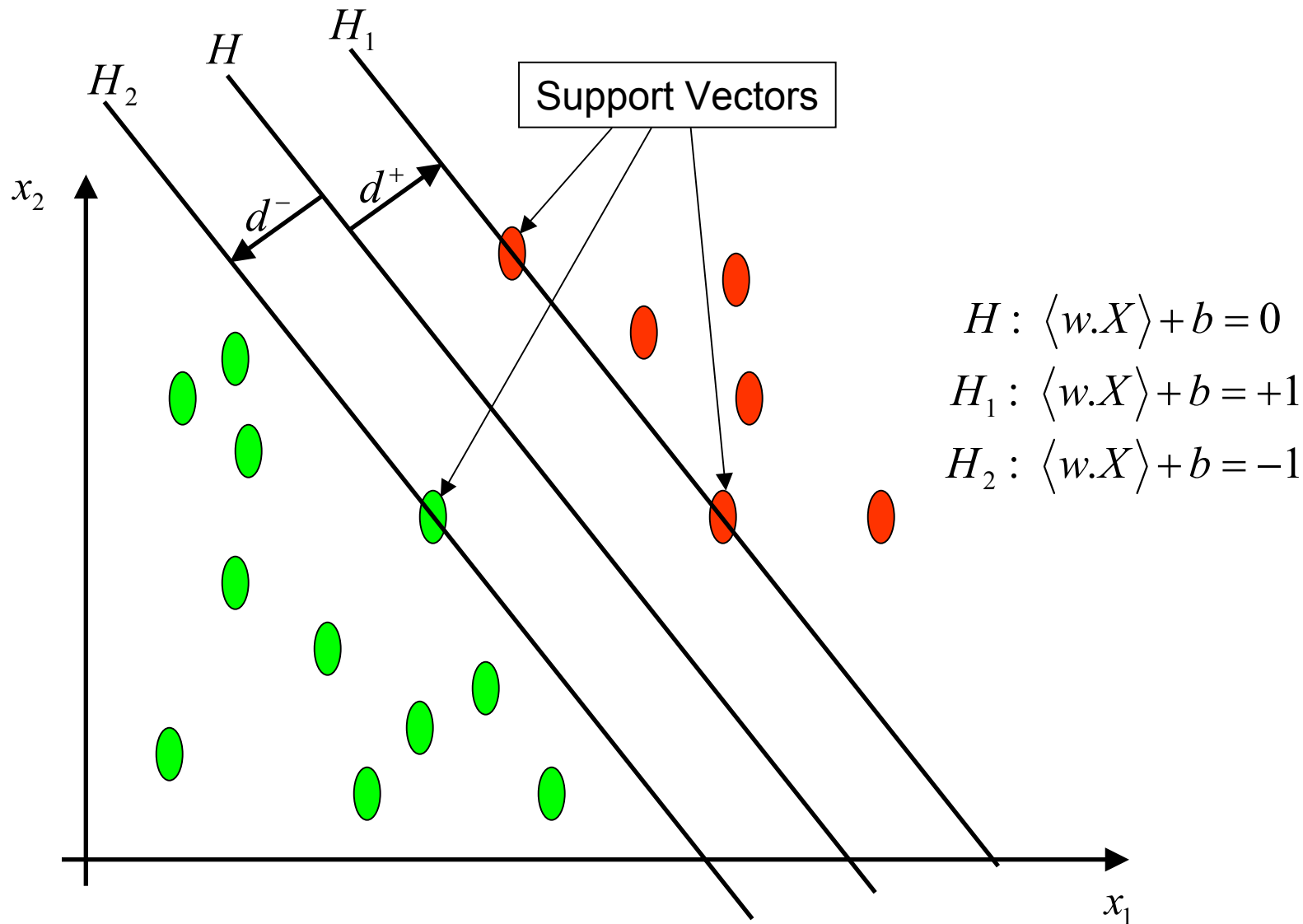
$$\langle w.X^i \rangle + b \leq -1 \quad \text{when } y^i = -1$$

$$H : \langle w.X \rangle + b = 0$$

Define

$$H_1 : \langle w.X \rangle + b = +1$$

$$H_2 : \langle w.X \rangle + b = -1$$



$d^+ / d^-$  = shortest distance to the closest positive / negative point

Margin of separation =  $d^+ + d^-$

Desired: A Classifier with as big margin as possible

Distance between  $H$  and  $H_1 = \frac{1}{\|w\|_2}$  (geometric margin)

Distance between  $H_2$  and  $H_1 = \frac{2}{\|w\|_2}$  (total margin)

In order to maximize the margin we need to minimize  $\|w\|$

Under this condition there are no data points between  $H_1$  and  $H_2$

$$\left. \begin{array}{l} \langle w.X^i \rangle + b \geq +1 \text{ when } y^i = +1 \\ \langle w.X^i \rangle + b \leq -1 \text{ when } y^i = -1 \end{array} \right\} \text{ Can be combined as } y^i (\langle w.X^i \rangle + b) \geq 1$$



The maximum Margin Classifier can be obtained by solving the following optimization problem:

$$\begin{aligned} & \underset{w,b}{\text{Minimize}} \quad \langle w.w \rangle \\ & y^i (\langle w.X^i \rangle + b) \geq 1, \quad i = 1, 2, \dots, N \end{aligned}$$

Formulating the Lagrangian, the primal Lagrangian is obtained as,

$$L(w, b, \alpha) = \frac{1}{2} \langle w.w \rangle - \sum_{i=1}^N \alpha_i [y^i (\langle w.X^i \rangle + b) - 1]$$

where  $\alpha_i \geq 0$  are Lagrangian multipliers

The corresponding dual is found by differentiating w.r.t  $w$  and  $b$  and imposing stationarity,

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^N y^i \alpha_i X^i = 0$$
$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N y^i \alpha_i = 0$$

Simplifying the relations,

$$w = \sum_{i=1}^N y^i \alpha_i X^i$$
$$\sum_{i=1}^N y^i \alpha_i = 0$$

Resubstituting back in the primal Lagrangian we get,

$$\begin{aligned}
 L(w, b, \alpha) &= \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^N \alpha_i \left[ y^i (\langle w, X^i \rangle + b) - 1 \right] \\
 &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y^i y^j \alpha_i \alpha_j \langle X^i, X^j \rangle - \sum_{i=1}^N \sum_{j=1}^N y^i y^j \alpha_i \alpha_j \langle X^i, X^j \rangle + \sum_{i=1}^N \alpha_i \\
 &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y^i y^j \alpha_i \alpha_j \langle X^i, X^j \rangle
 \end{aligned}$$

Hence the dual problem is

$$\begin{aligned}
 &\text{Max}_{\alpha} \quad \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y^i y^j \alpha_i \alpha_j \langle X^i, X^j \rangle \\
 &\text{subject to} \quad \sum_{i=1}^N y^i \alpha_i = 0, \\
 &\quad \quad \quad \alpha_i \geq 0, i = 1, 2, \dots, N
 \end{aligned}$$

The dual problem can be easily solved by readily available quadratic programming solvers to give,

$$\alpha^*$$

The weight vector that realises the maximal margin hyperplane is given by

$$w^* = \sum_{i=1}^N y^i \alpha_i^* X^i$$

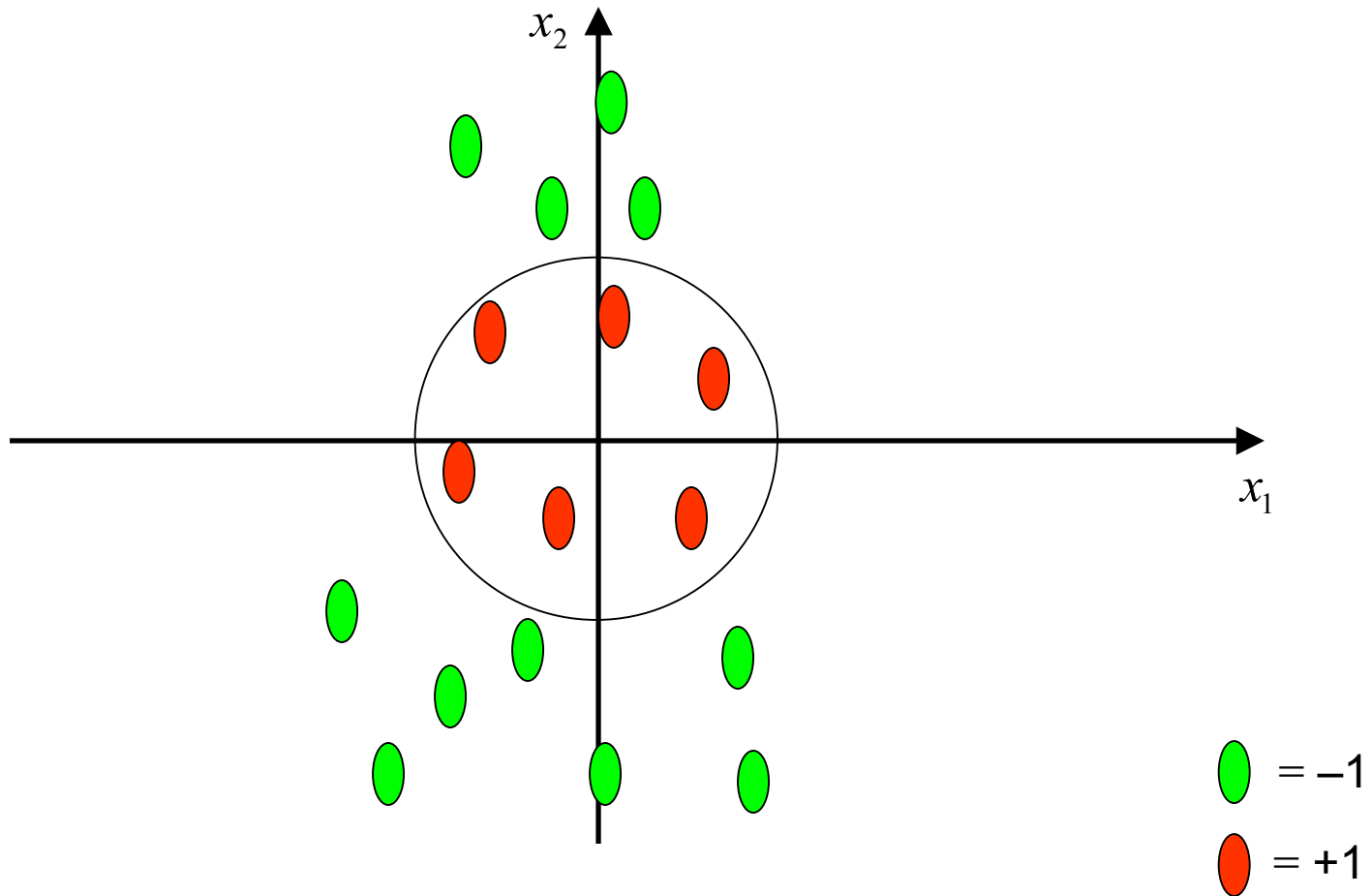
The value of  $b$  does not appear in the dual problem and so  $b^*$  is found from primal constraints

$$b^* = - \frac{\max_{y^i=-1} (\langle w^* . X^i \rangle) + \min_{y^i=1} (\langle w^* . X^i \rangle)}{2}$$

In the dual solution it turns out that most of the  $\alpha_i^*$  are zero. Non-zero  $\alpha_i^*$  occur for the points that are closest to the hyperplane. These are known as the **support vectors** ( $SV$ ). The optimal hyperplane is given by

$$\begin{aligned} f(X) &= \sum_{i=1}^N y^i \alpha_i^* \langle X^i . X \rangle + b^* \\ &= \sum_{i \in SV} y^i \alpha_i^* \langle X^i . X \rangle + b^* \end{aligned}$$

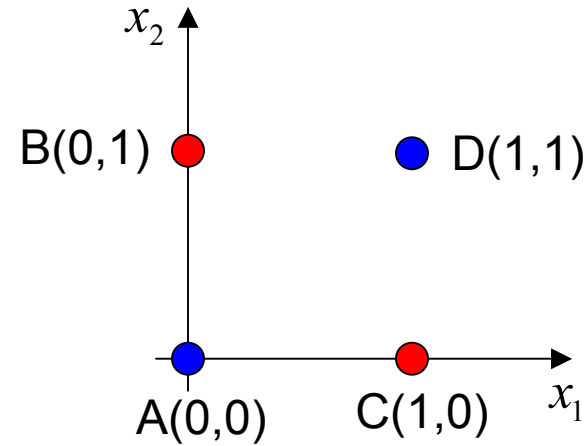
# Problems with Linear SVM



What if the decision function is not linear ?

# XOR Problem

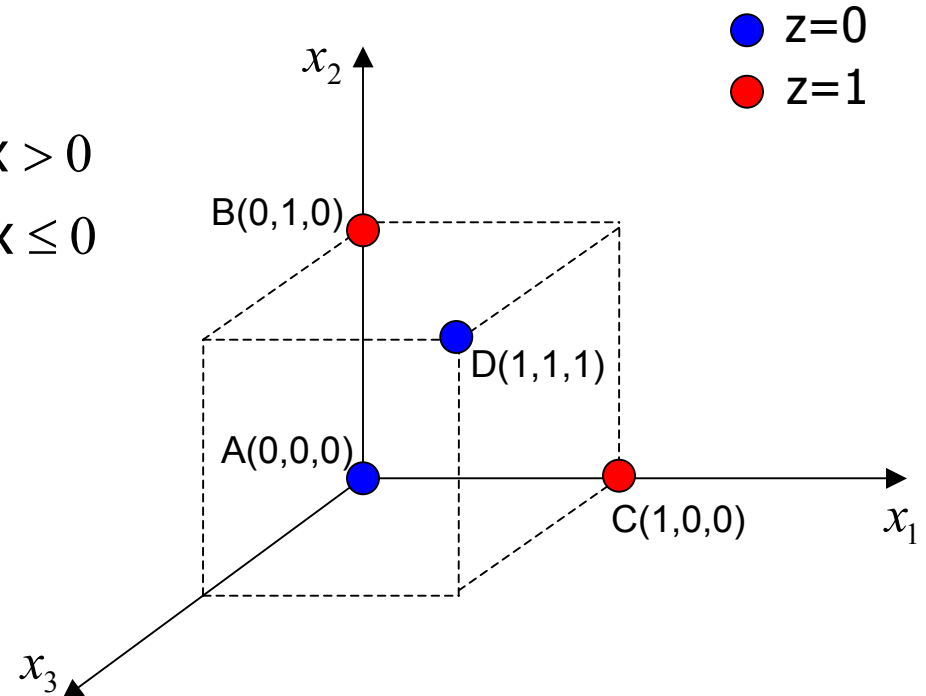
$x_1$	$x_2$	$z$
0	0	0
0	1	1
1	0	1
1	1	0

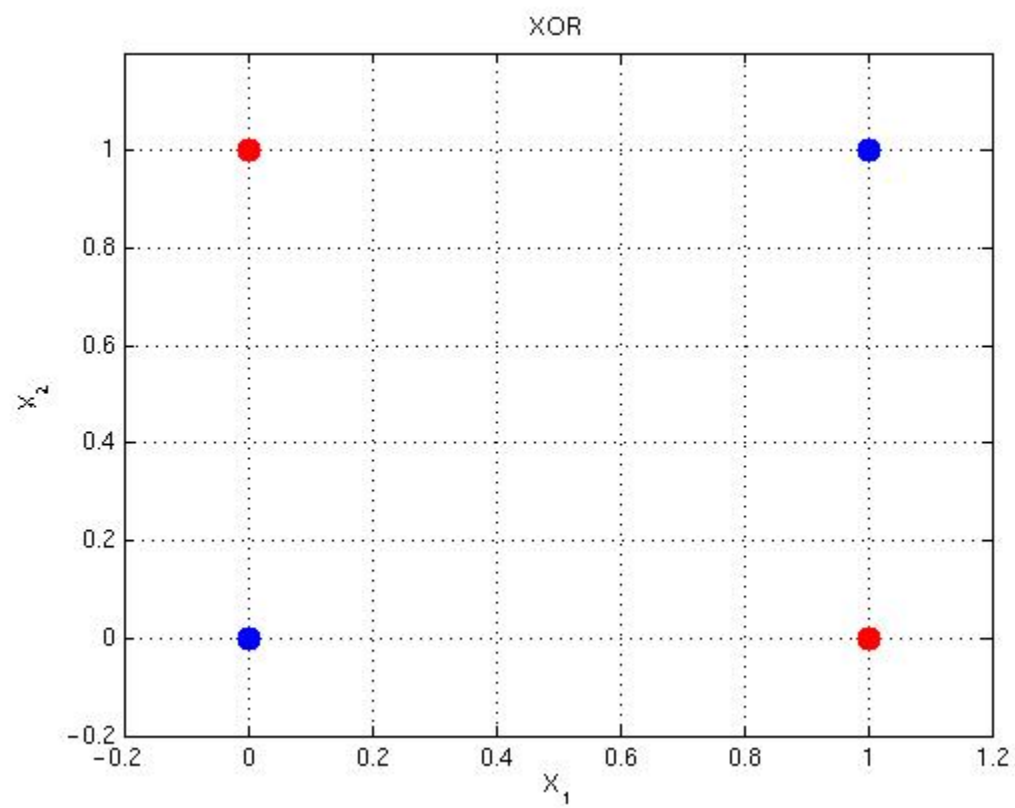


$$[x_1 \ x_2] \Rightarrow [x_1 \ x_2 \ x_3]$$

$$x_3 = T(x_1 x_2 - 0.5) \text{ where } T(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

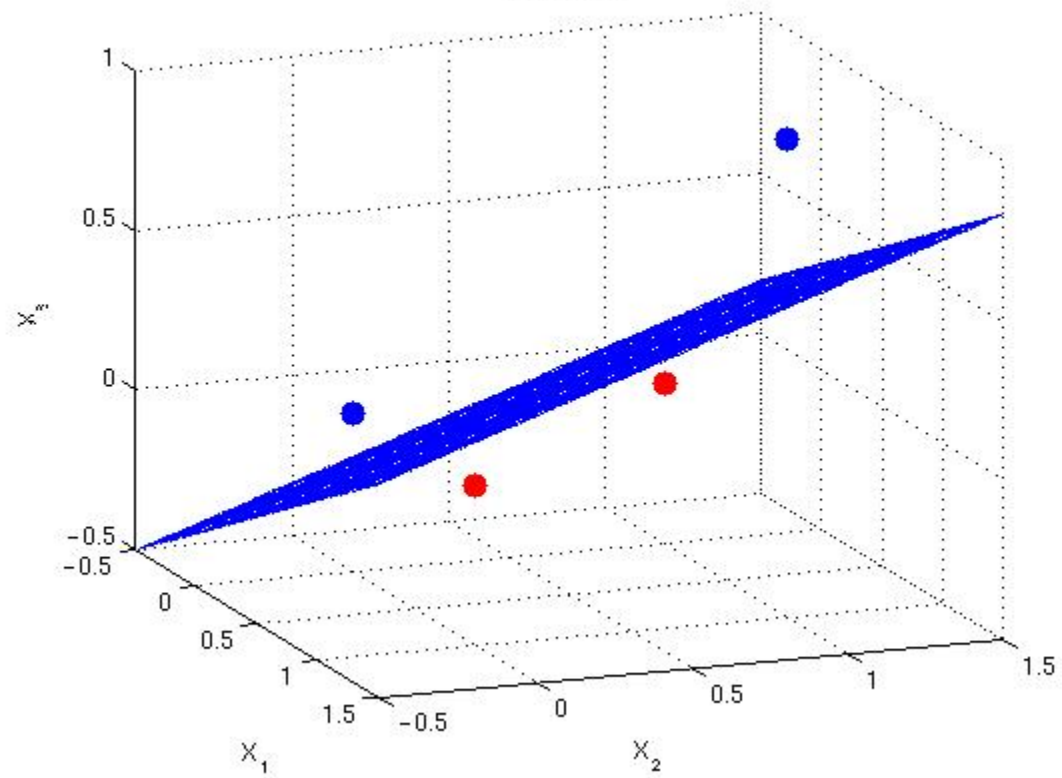
$x_1$	$x_2$	$x_3$	$z$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0







XOR in 3D



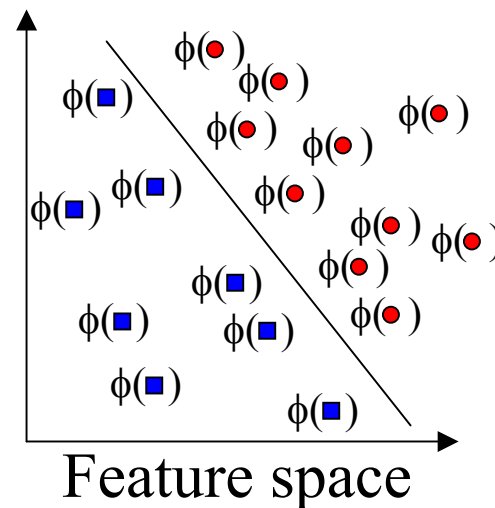
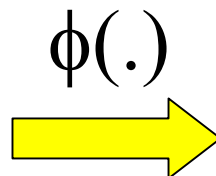
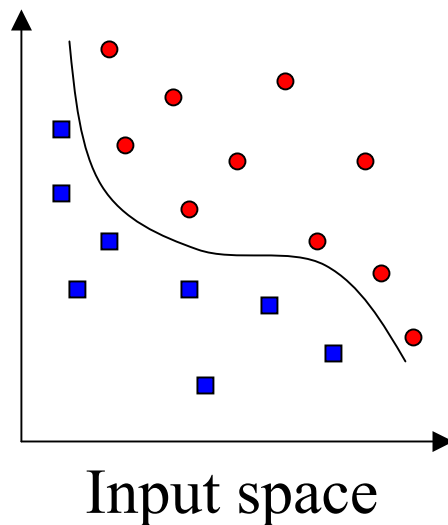


# Extension to Non-linear Decision Boundary

---

- So far, we only consider large-margin classifier with a linear decision boundary
- How to generalize it to become nonlinear?
- Key idea: transform  $\mathbf{x}_i$  to a higher dimensional space to “make life easier”
  - Input space: the space the point  $\mathbf{x}_i$  are located
  - Feature space: the space of  $\phi(\mathbf{x}_i)$  after transformation
- Why transform?
  - Linear operation in the feature space is equivalent to non-linear operation in input space
  - Classification can become easier with a proper transformation. In the XOR problem, for example, adding a new feature of  $x_1x_2$  make the problem linearly separable

# Transforming the Data



- Computation in the feature space can be costly because it is high dimensional
  - The feature space is typically infinite-dimensional!
- The kernel trick comes to rescue

# The Kernel Trick

- Recall the SVM optimization problem

$$\begin{aligned} \max. \quad W(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{subject to } C &\geq \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- The data points only appear as **inner product**
- As long as we can calculate the inner product in the feature space, we do not need the mapping explicitly
- Many common geometric operations (angles, distances) can be expressed by inner products
- Define the kernel function  $K$  by

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



## An Example for $\phi(\cdot)$ and $K(\cdot, \cdot)$

- Suppose  $\phi(\cdot)$  is given as follows

$$\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- An inner product in the feature space is

$$\langle \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right), \phi\left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) \rangle = (1 + x_1y_1 + x_2y_2)^2$$

- So, if we define the kernel function as follows, there is no need to carry out  $\phi(\cdot)$  explicitly

$$K(\mathbf{x}, \mathbf{y}) = (1 + x_1y_1 + x_2y_2)^2$$

- This use of kernel function to avoid carrying out  $\phi(\cdot)$  explicitly is known as the **kernel trick**



# Kernel Functions

---

- In practical use of SVM, only the kernel function (and not  $\phi(\cdot)$ ) is specified
- Kernel function can be thought of as a similarity measure between the input objects
- Not all similarity measure can be used as kernel function, however
  - The kernel function needs to satisfy the Mercer function, i.e., the function is “positive-definite”
  - This has the consequence that the kernel matrix, where the  $(i,j)$ -th entry is the  $K(\mathbf{x}_i, \mathbf{x}_j)$ , is always positive definite
- Note that  $\mathbf{x}_i$  needs not be vectorial for the kernel function to exist. This opens up enormous opportunities for classification with sequences, graphs, etc., by SVM



# Examples of Kernel Functions

---

- Polynomial kernel with degree  $d$

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

- Radial basis function kernel with width  $\sigma$

$$K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2))$$

- Closely related to radial basis function neural networks

- Sigmoid with parameter  $\kappa$  and  $\theta$

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x}^T \mathbf{y} + \theta)$$

- It does not satisfy the Mercer condition on all  $\kappa$  and  $\theta$



# Modification Due to Kernel Function

- Change all inner products to kernel functions
- For training,

Original

$$\max. W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{subject to } C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$$

With kernel  
function

$$\max. W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to } C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$$





# Modification Due to Kernel Function

- For testing, the new data  $\mathbf{z}$  is classified as class 1 if  $f \geq 0$ , and as class 2 if  $f < 0$

Original

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}$$
$$f = \mathbf{w}^T \mathbf{z} + b = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}^T \mathbf{z} + b$$

With kernel  
function

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \phi(\mathbf{x}_{t_j})$$
$$f = \langle \mathbf{w}, \phi(\mathbf{z}) \rangle + b = \sum_{j=1}^s \alpha_{t_j} y_{t_j} K(\mathbf{x}_{t_j}, \mathbf{z}) + b$$

# Example

- Suppose we have 5 1D data points
  - $x_1=1, x_2=2, x_3=4, x_4=5, x_5=6$ , with 1, 2, 6 as class 1 and 4, 5 as class 2  $\Rightarrow y_1=1, y_2=1, y_3=-1, y_4=-1, y_5=1$
- We use the polynomial kernel of degree 2
  - $K(x,y) = (xy+1)^2$
  - $C$  is set to 100
- We first find  $\alpha_i$  ( $i=1, \dots, 5$ ) by

$$\max. \quad \sum_{i=1}^5 \alpha_i - \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \alpha_i \alpha_j y_i y_j (x_i x_j + 1)^2$$

$$\text{subject to } 100 \geq \alpha_i \geq 0, \quad \sum_{i=1}^5 \alpha_i y_i = 0$$

# Example

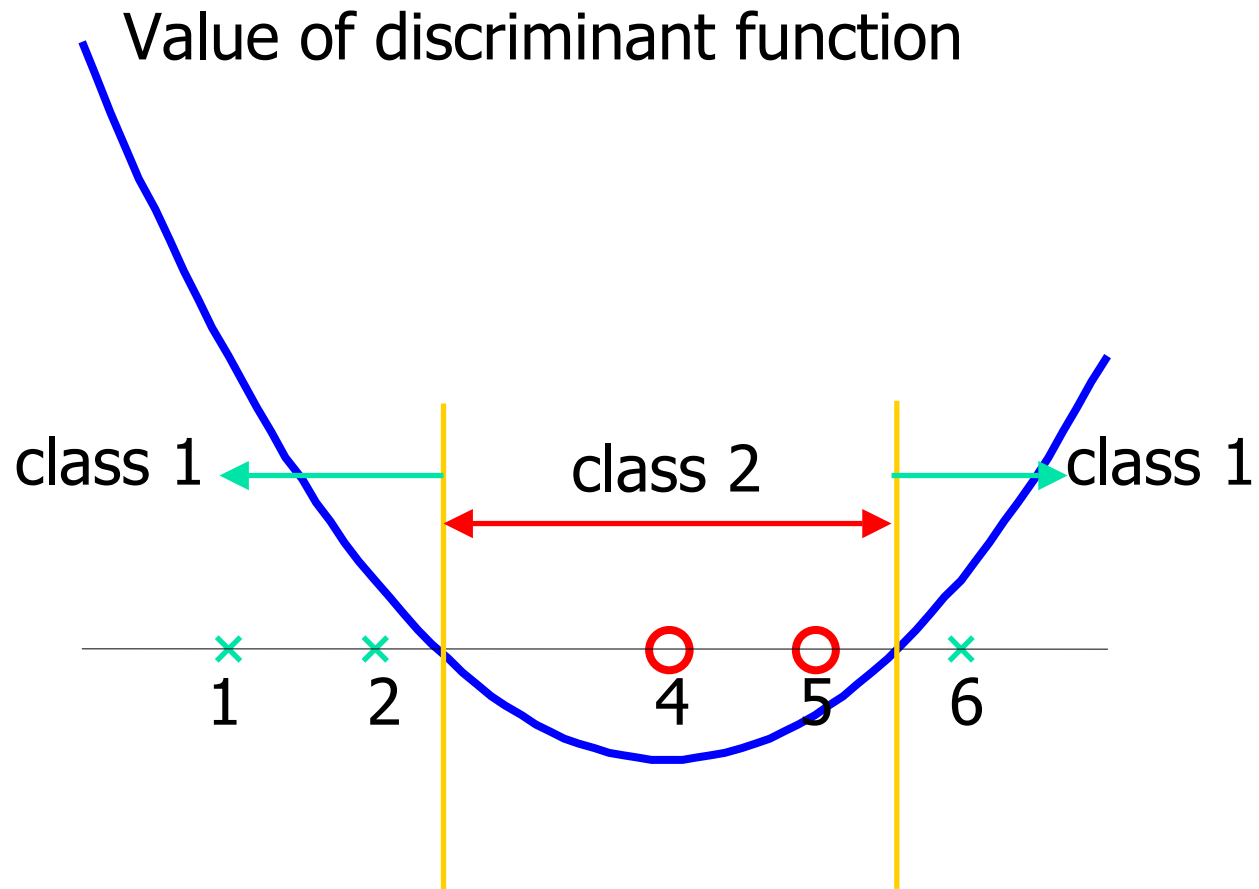
- By using a QP solver, we get
  - $\alpha_1=0, \alpha_2=2.5, \alpha_3=0, \alpha_4=7.333, \alpha_5=4.833$
  - Note that the constraints are indeed satisfied
  - The support vectors are  $\{x_2=2, x_4=5, x_5=6\}$
- The discriminant function is

$$\begin{aligned} f(y) \\ &= 2.5(1)(2y + 1)^2 + 7.333(-1)(5y + 1)^2 + 4.833(1)(6y + 1)^2 + b \\ &= 0.6667x^2 - 5.333x + b \end{aligned}$$

- $b$  is recovered by solving  $f(2)=1$  or by  $f(5)=-1$  or by  $f(6)=1$ , as  $x_2, x_4, x_5$  lie on  $y_i(w^T \phi(z) + b) = 1$  and all give  $b=9$

$$\longrightarrow f(y) = 0.6667x^2 - 5.333x + 9$$

# Example





# Summary: Steps for Classification

---

- Prepare the pattern matrix
- Select the kernel function to use
- Select the parameter of the kernel function and the value of  $C$ 
  - You can use the values suggested by the SVM software, or you can set apart a validation set to determine the values of the parameter
- Execute the training algorithm and obtain the  $\alpha_i$
- Unseen data can be classified using the  $\alpha_i$  and the support vectors



# Choosing the Kernel Function

---

- Probably the most tricky part of using SVM.
- The kernel function is important because it creates the kernel matrix, which summarize all the data
- Many principles have been proposed (diffusion kernel, Fisher kernel, string kernel, ...)
- There are even research to estimate the kernel matrix from available information
- In practice, a low degree polynomial kernel or RBF kernel with a reasonable width is a good initial try
- Note that SVM with RBF kernel is closely related to RBF neural networks



# Multi-class Classification

---

- SVM is basically a two-class classifier
- One can change the QP formulation to allow multi-class classification
- More commonly, the data set is divided into two parts “intelligently” in different ways and a separate SVM is trained for each way of division
- Multi-class classification is done by combining the output of all the SVM classifiers
  - Majority rule
  - Error correcting code
  - Directed acyclic graph



# Why SVM Work?

---

- The feature space is often very high dimensional. Why don't we have the curse of dimensionality?
- A classifier in a high-dimensional space has many parameters and is hard to estimate
- Vapnik argues that the fundamental problem is not the number of parameters to be estimated. Rather, the problem is about the flexibility of a classifier
- Typically, a classifier with many parameters is very flexible, but there are also exceptions
  - Let  $x_i = 10^i$  where  $i$  ranges from 1 to  $n$ . The classifier  $y = \text{sign}(\sin(\alpha x))$  can classify all  $x_i$  correctly for all possible combination of class labels on  $x_i$
  - This 1-parameter classifier is very flexible





# Why SVM works?

---

- Vapnik argues that the flexibility of a classifier should not be characterized by the number of parameters, but by the flexibility (capacity) of a classifier
  - This is formalized by the “VC-dimension” of a classifier
- The addition of  $\frac{1}{2}||w||^2$  has the effect of restricting the VC-dimension of the classifier in the feature space
- The SVM objective can also be justified by structural risk minimization: the empirical risk (training error), plus a term related to the generalization ability of the classifier, is minimized
- Another view: the SVM loss function is analogous to ridge regression. The term  $\frac{1}{2}||w||^2$  “shrinks” the parameters towards zero to avoid overfitting