



CS 4104

APPLIED MACHINE LEARNING

Dr. Hashim Yasin

**National University of Computer
and Emerging Sciences,
Faisalabad, Pakistan.**

SVM ... PRELIMINARIES



Linear Separable Data

3

- **Given:** m examples $(x_1, c_1), \dots, (x_m, c_m)$
- **Goal:** Learn classification!
- Most simple case: *binary classification*, where each example shows
 - ▣ n -dimensional input data vector $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,n})$ and
 - ▣ its binary classification $c_i \in \{+1, -1\}$
- e.g., classification of all web pages into “*related to computer science*” and “*not related to computer science*”:
 - ▣ **Given:** data vectors \mathbf{x}_i with binary elements x_{ij} for appearance or missing appearance of a relevant keyword.
 - ▣ **Goal:** classification of new web pages with small prediction error

Linear Separable Data

4

- Classification of n -dimensional input data

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

is possible using a linear separation function

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

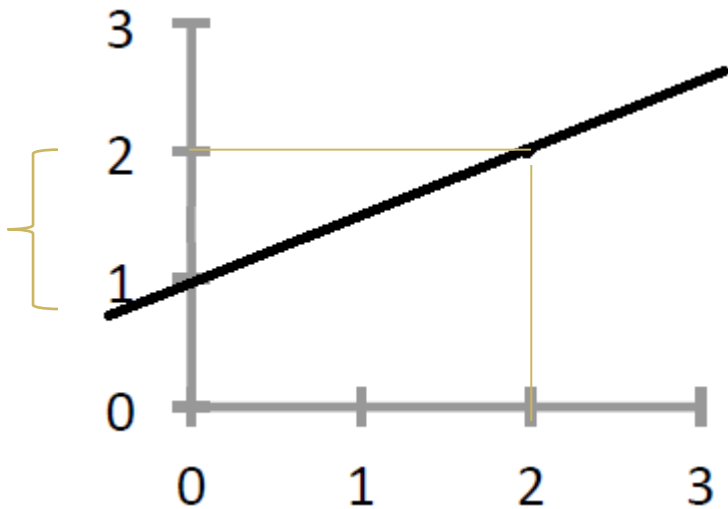
- \mathbf{x} is classified as positive ($c = +1$), if $f(\mathbf{x}) \geq 0$.
- \mathbf{x} is classified as negative ($c = -1$), if $f(\mathbf{x}) < 0$.

Linear Separable Data

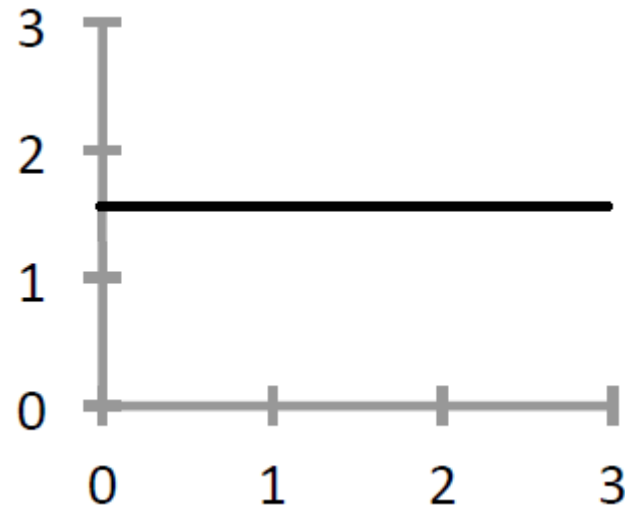
5

$$y = \theta_0 + \theta_1 x$$

$$\theta = \frac{\text{change in } Y}{\text{change in } X}$$



$$\theta_1 = 0.5, \theta_0 = 1$$

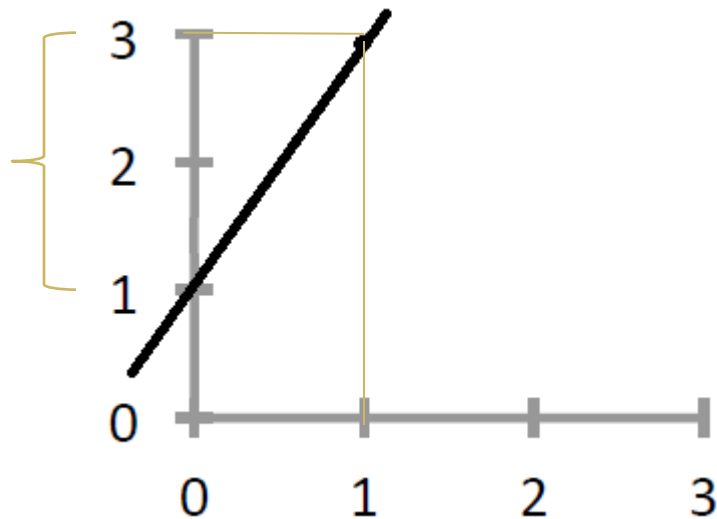


$$\theta_1 = 0, \theta_0 = 1.5$$

Linear Separable Data

6

$$y = \theta_1 x + \theta_0$$



$$\theta_1 = 2, \theta_0 = 1$$

$$\theta = \frac{\text{change in } Y}{\text{change in } X}$$



$$\theta_1 = 1, \theta_0 = 0$$

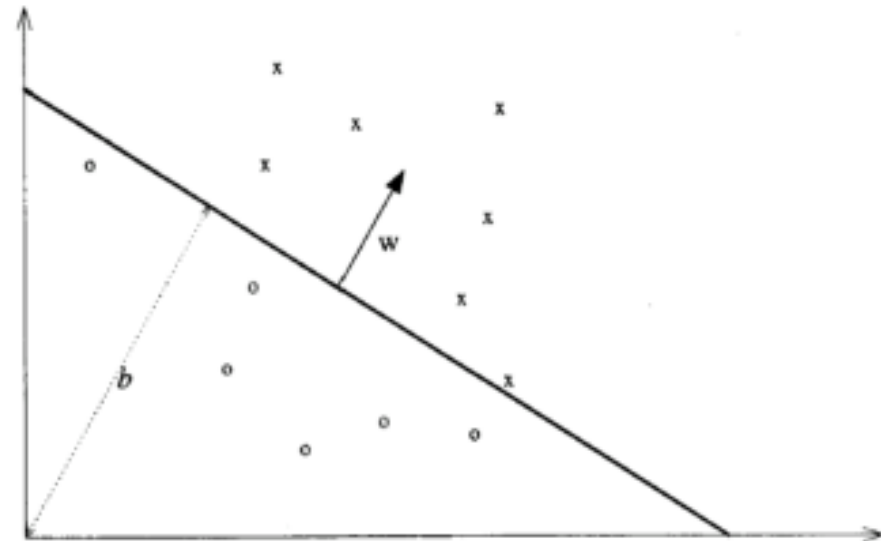
Linear Separable Data

7

The data can be separated in the n-dimensional data space by a planar hyper plane $\langle \mathbf{w}, \mathbf{x} \rangle - b = 0$.

Parameter \mathbf{w} defines the normal vector of the hyper plane, parameter b stands for the bias of the hyper plane.

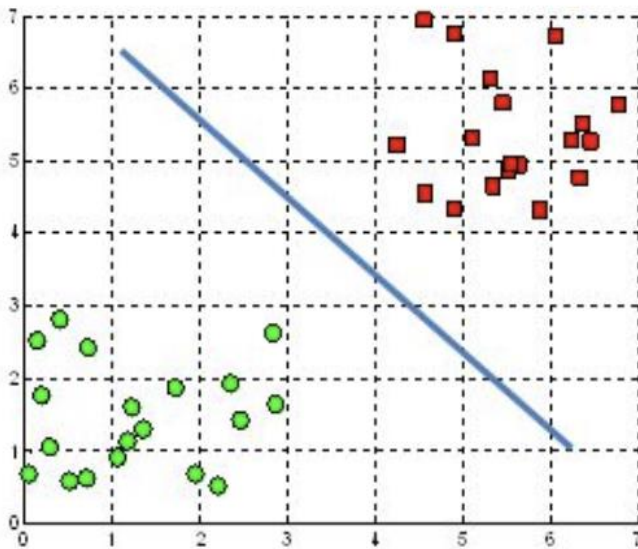
Planer hyper plane
 $\langle \mathbf{w}, \mathbf{x} \rangle - b = 0$



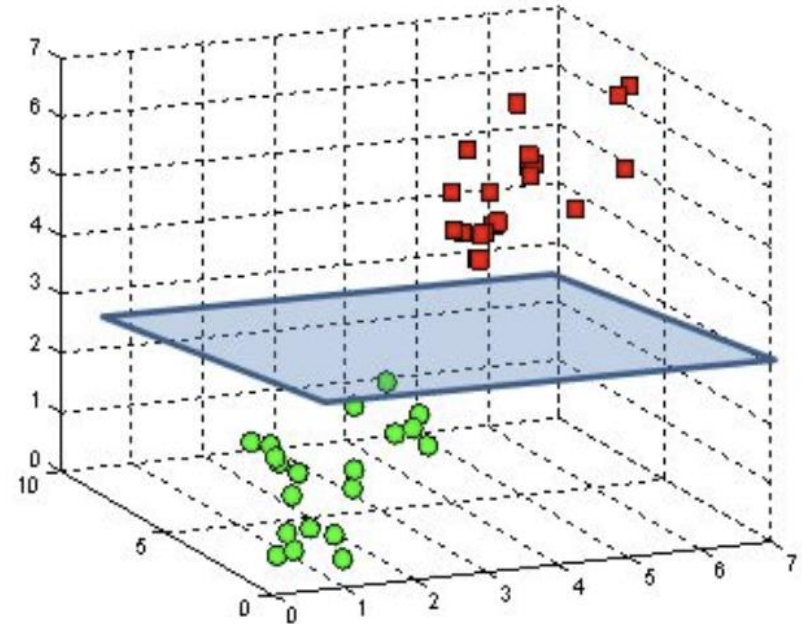
Hyperplane

8

A hyperplane in \mathbb{R}^2 is a line



A hyperplane in \mathbb{R}^3 is a plane

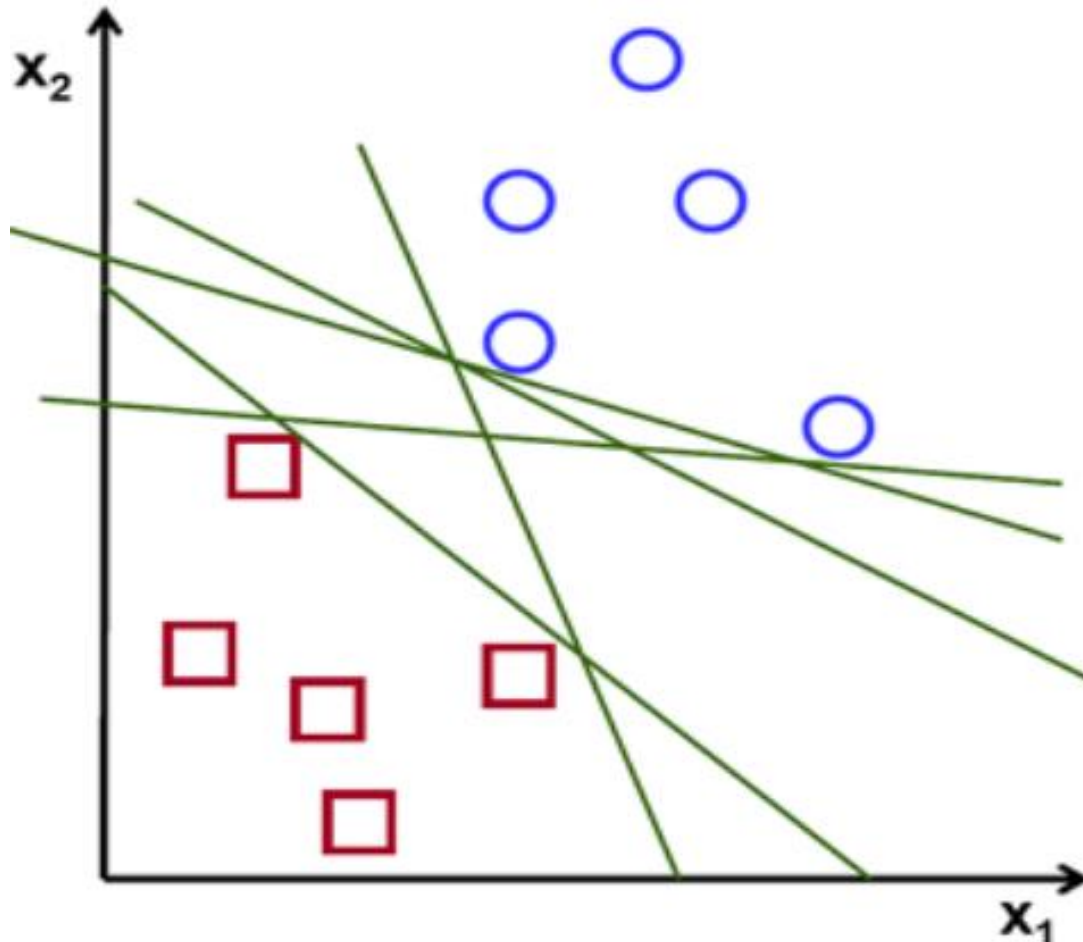


SUPPORT VECTOR MACHINE



Support Vector Machine (SVM)

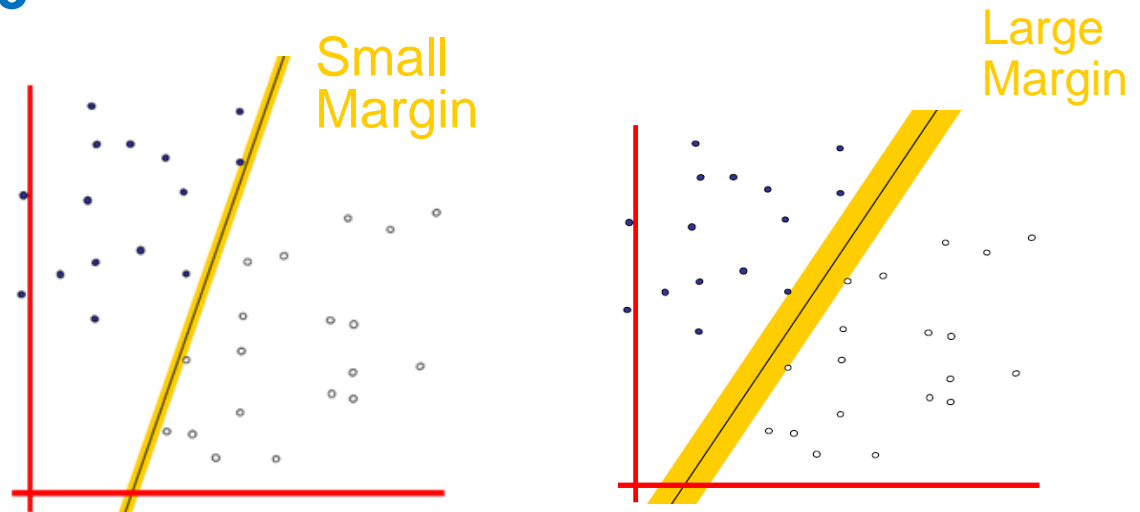
10



Support Vector Machine (SVM)

11

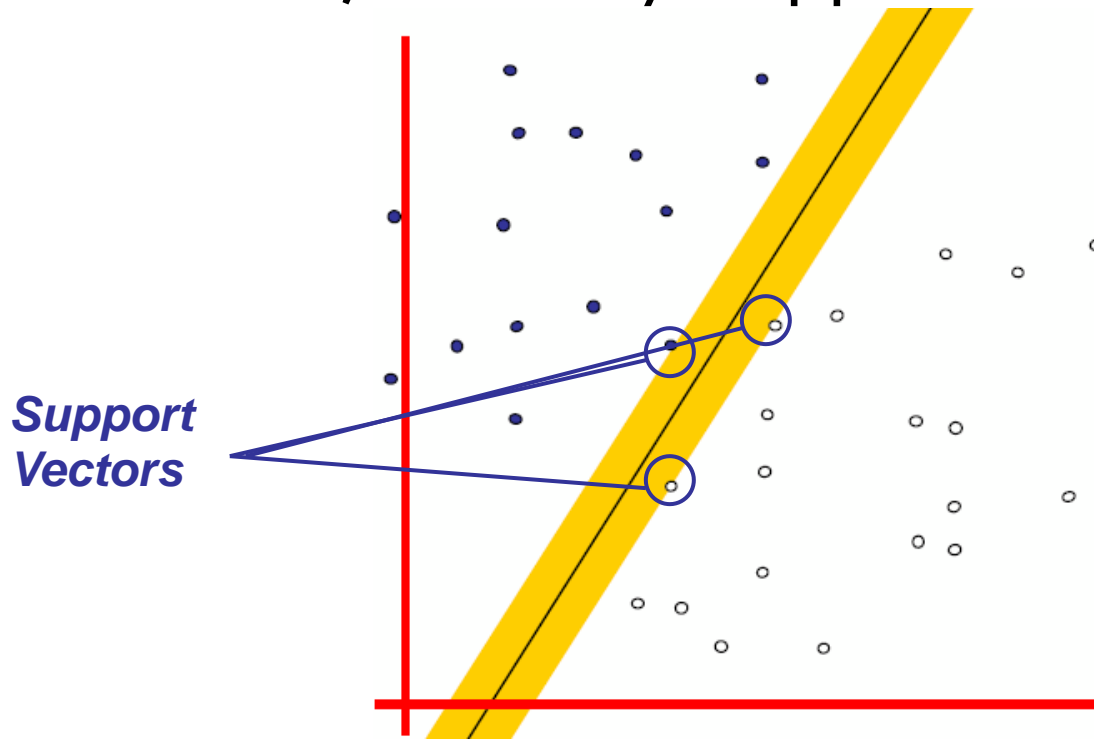
- Derives a linear separator
- informally chooses that separating hyper plane that maximizes the so-called *margin*, i.e., the space between positive and negative examples
- → **Maximum Margin Classifier**



Support Vector Machine (SVM)

12

- That points closest to the separator are called **Support Vectors**, since they “support” the separator.

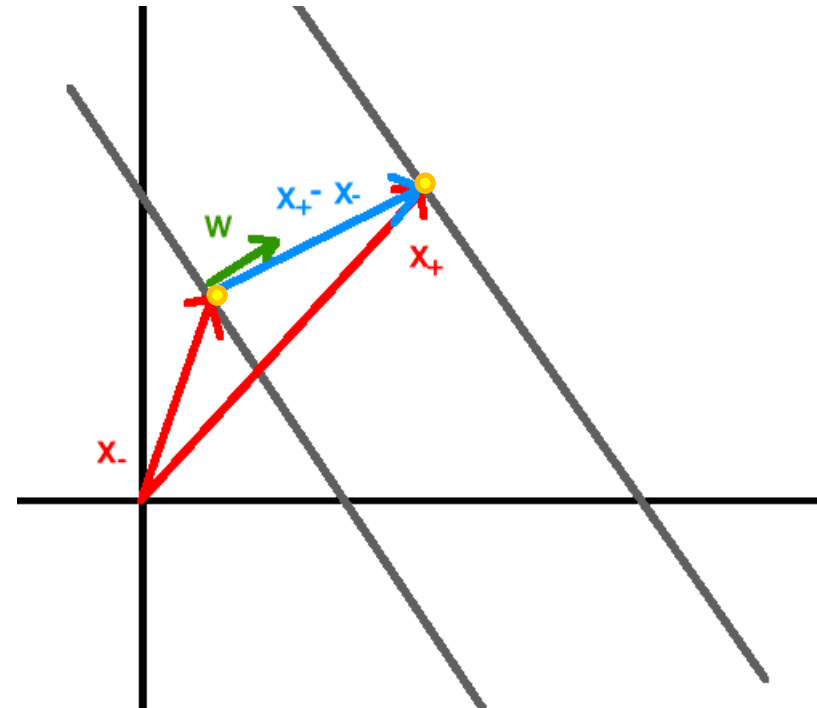


Support Vector Machine (SVM)

13

- To get an equation for the **width of the margin**,
 - ▣ **subtract** the first support vector from the second support vector and
 - ▣ **multiply** the result by the **unit vector** of $w = \frac{\vec{w}}{\|\vec{w}\|}$ which is always perpendicular to the decision boundary.

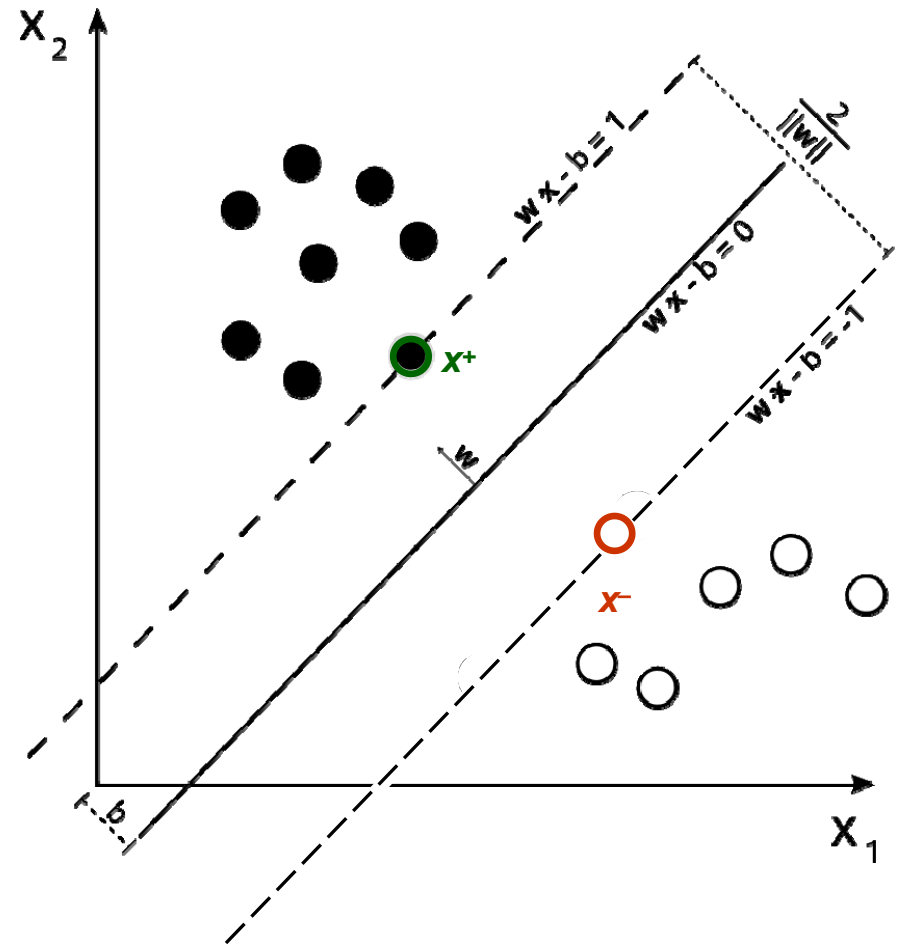
$$width = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$



Support Vector Machine (SVM)

14

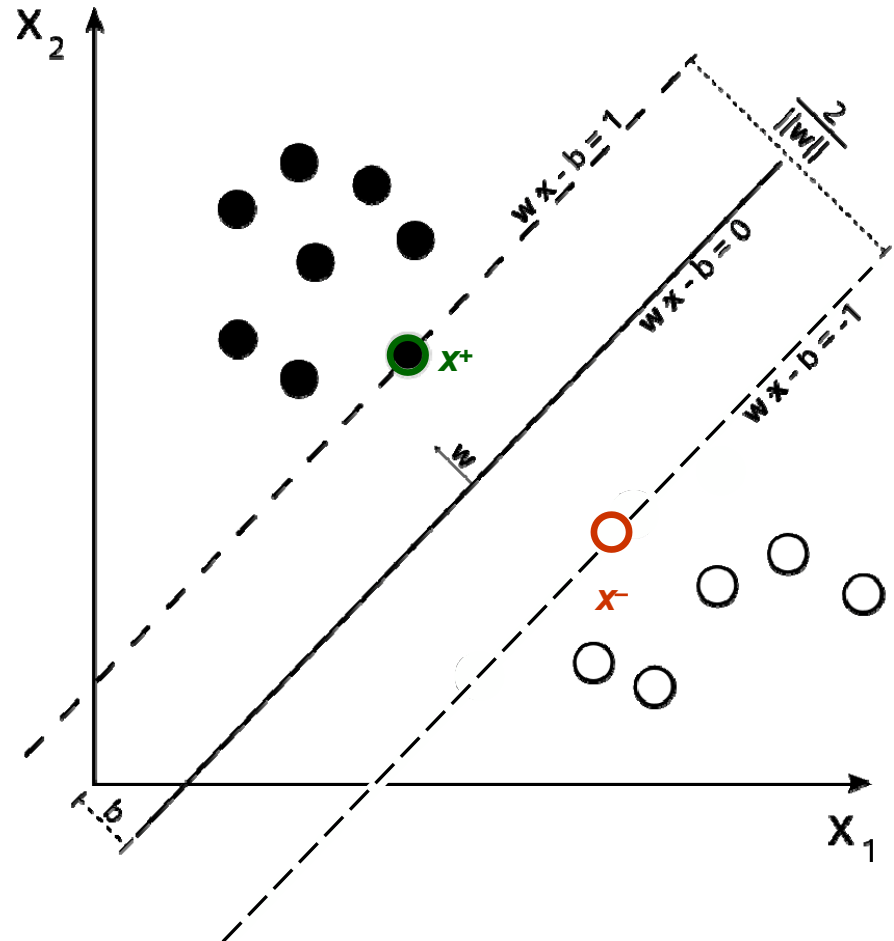
- First, we scale $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} - b$ that way, that the values of the support vectors are $+1$ and -1 , respectively.
- Thus, the width M of the margin can be expressed as a function of \mathbf{w} .
- $\langle \mathbf{w}, \mathbf{x}^+ \rangle - b = +1$ $\langle \mathbf{w}, \mathbf{x}^- \rangle - b = -1$
 $\sim \langle \mathbf{w}, (\mathbf{x}^+ - \mathbf{x}^-) \rangle = 2$
 $\sim M = \langle (\mathbf{w}/\|\mathbf{w}\|), (\mathbf{x}^+ - \mathbf{x}^-) \rangle$
 $= 2/\|\mathbf{w}\|$



Support Vector Machine (SVM)

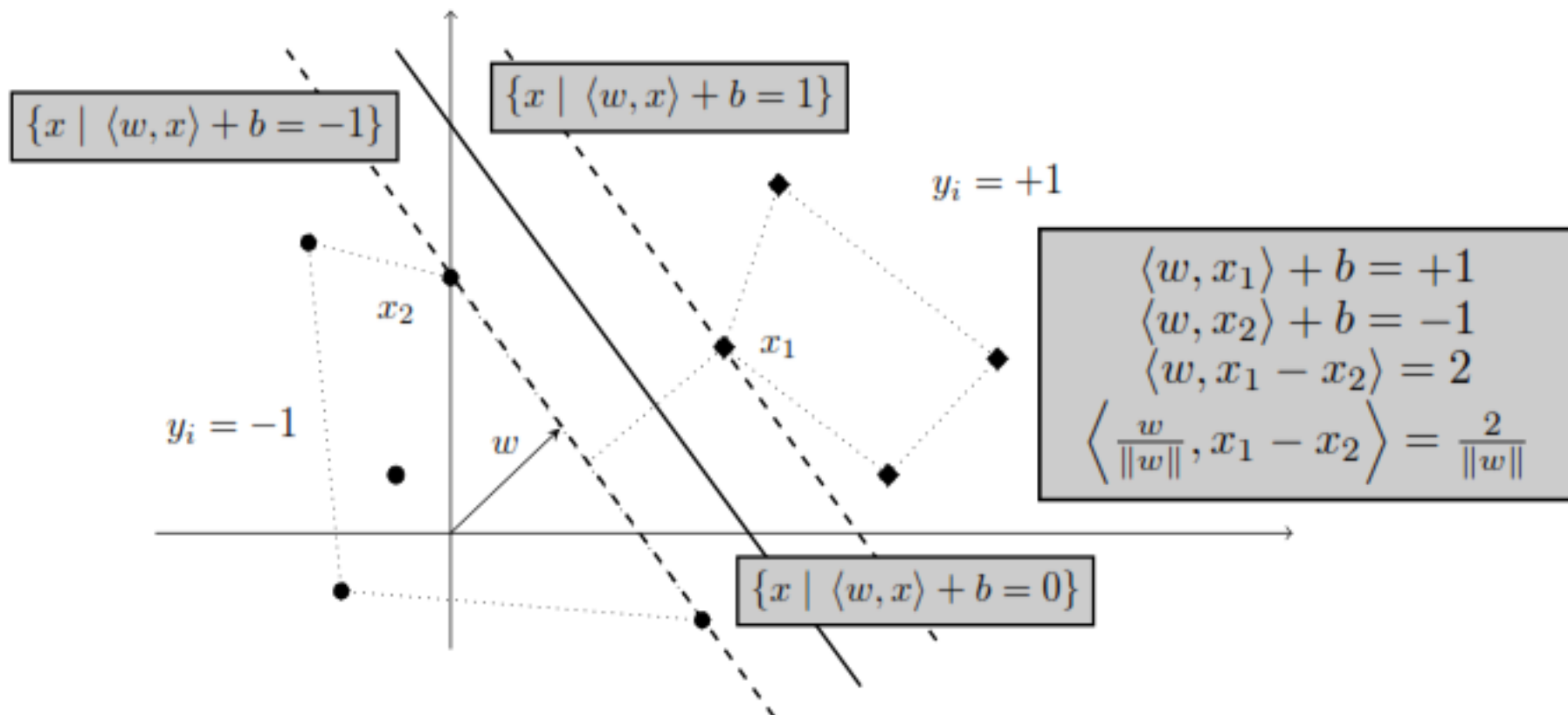
15

- $\langle \mathbf{w}, \mathbf{x}^+ \rangle - b = +1$ $\langle \mathbf{w}, \mathbf{x}^- \rangle - b = -1$
 $\sim \langle \mathbf{w}, (\mathbf{x}^+ - \mathbf{x}^-) \rangle = 2$
 $\sim M = \langle (\mathbf{w}/\|\mathbf{w}\|), (\mathbf{x}^+ - \mathbf{x}^-) \rangle$
 $= 2/\|\mathbf{w}\|$
- Maximizing $M=2/\|\mathbf{w}\|$ is equal to minimize $\|\mathbf{w}\|/2$.
- With subject to the constraint, that all examples of the training data are correctly classified.



Support Vector Machine (SVM)

16



Support Vector Machine (SVM)

17

- The problem of maximizing the margin therefore reduces to

$$\begin{aligned} \max_{w,b} \quad & \frac{1}{\|w\|} \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1 \text{ for all } i, \end{aligned}$$

Support Vector Machine (SVM)

18

- The problem of maximizing the margin therefore reduces to

$$\begin{aligned} \max_{w,b} \quad & \frac{1}{\|w\|} \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1 \text{ for all } i, \end{aligned}$$

- Which is equivalently to

$$\min \frac{1}{2} \|w\|^2 \quad \text{since} \quad \frac{d}{dx} \frac{1}{2} x^2 = x$$

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(\langle w, x_i \rangle + b) \geq 1 \text{ for all } i. \end{aligned}$$

Support Vector Machine (SVM)

19

The picture depicts

- the well classified points in black,

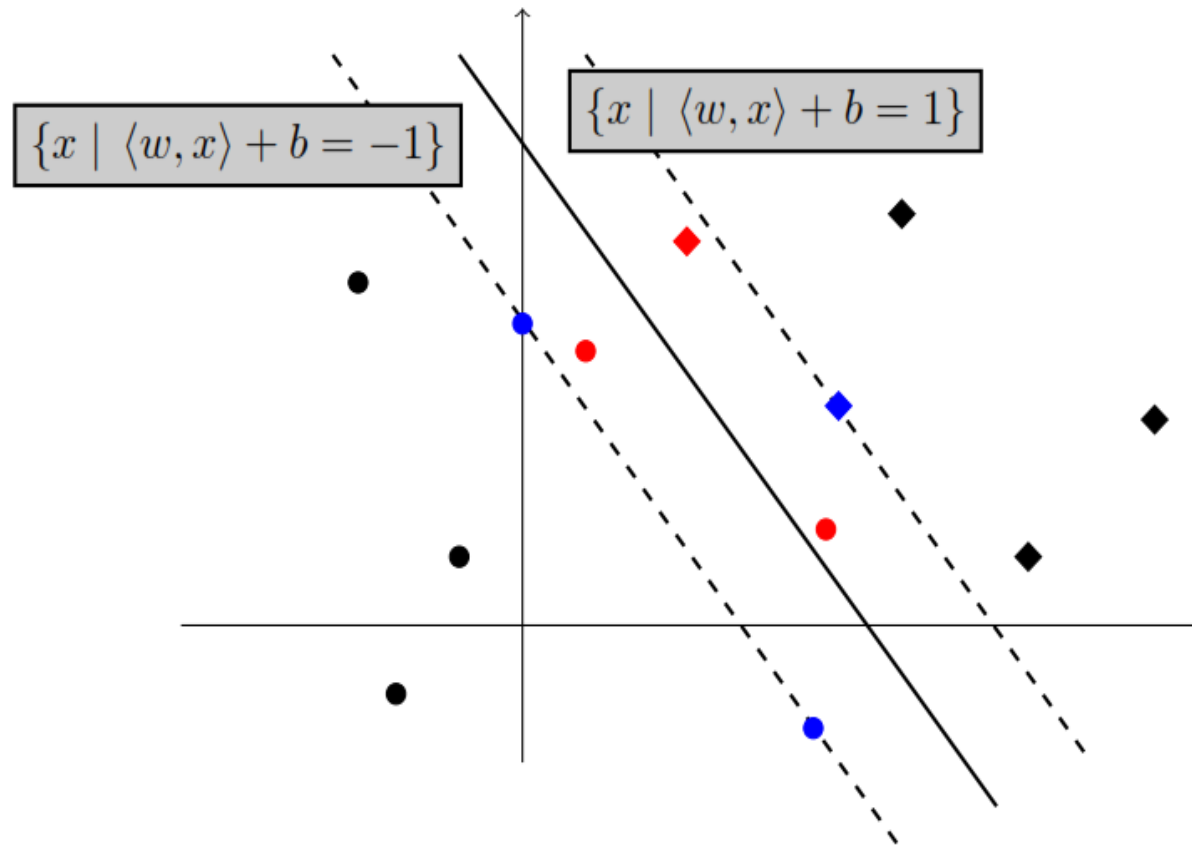
$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 1$$

- the support vectors in blue,

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1$$

- margin errors in red.

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) < 1$$



Support Vector Machine (SVM)

20

- This is a **constrained convex optimization problem** with a quadratic objective function and linear constraints.
- In deriving this equation, we implicitly assume that the **data is linearly separable**, that is, there is a hyperplane that correctly classifies the training data.
- Such a classifier is called a **hard margin classifier**.
- If the **data is not linearly separable**, then does not have a solution.

SVM ... KERNEL TRICK

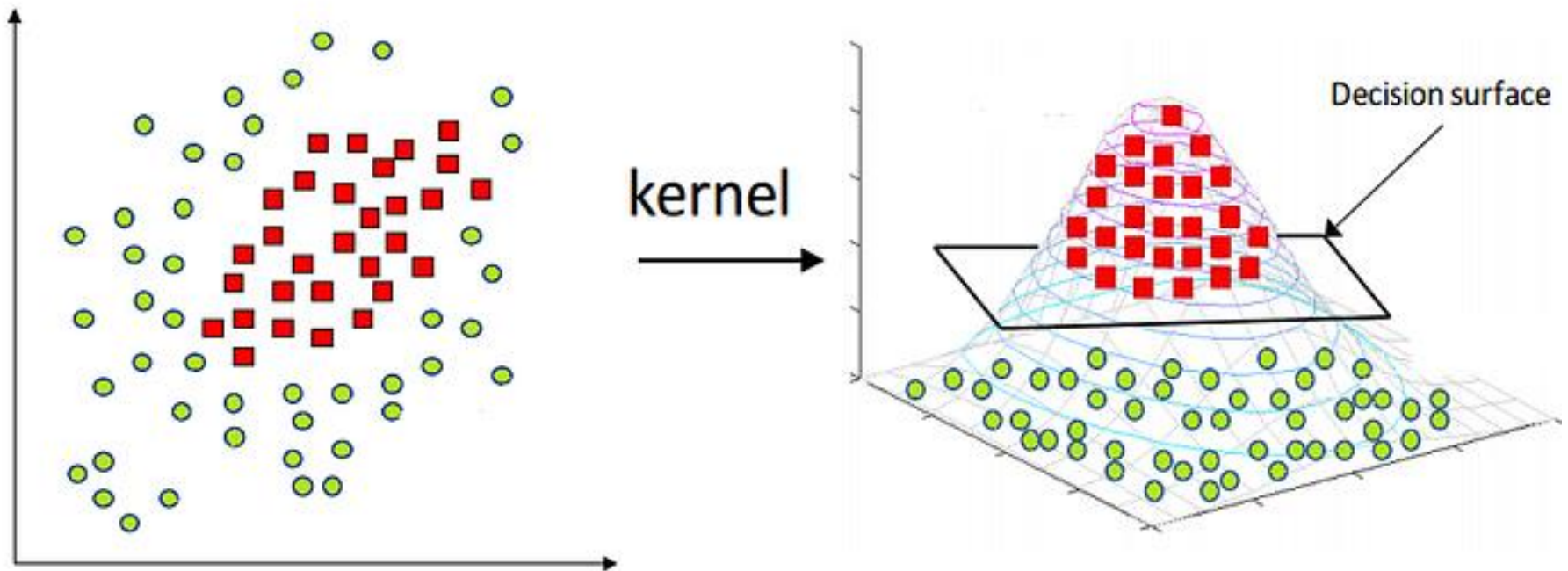
SVM ... Kernel Tricks

22

- Kernel Trick is widely used in the Support Vector Machines (SVM) model to bridge linearity and non-linearity.
- It converts non-linear lower-dimension space to a higher dimension space thereby we can get a linear classification.
- So, we are projecting the data with some extra features so that it can convert to a higher dimension space.

SVM ... Kernel Tricks

23



SVM ... Kernel Tricks

24

Kernel Function:

- A function that takes as its inputs vectors in the original space and returns the dot product of the vectors in the feature space is called a *kernel function*
- More formally, if we have data $\mathbf{x}, \mathbf{z} \in X$ and a map $\phi : X \rightarrow \Re^N$ then

$$k(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$$

is a kernel function

SVM ... Kernel Tricks

25

- Let us consider a simple kernel which is:

$$K(x, y) = \langle f(x), f(y) \rangle$$

where,

- ▣ K is the kernel function,
- ▣ X and Y are the dimensional inputs,
- ▣ f is the map from n-dimensional to m-dimensional space and,
- ▣ $\langle x, y \rangle$ is the dot product.

SVM ... Kernel Tricks

26

- Let us say that we have two points,
 - ▣ $x = (5, 6, 7)$ and $y = (8, 9, 10)$

- As we have seen, $K(x, y) = \langle f(x), f(y) \rangle$,
let us first calculate $\langle f(x), f(y) \rangle$

SVM ... Kernel Tricks

27

$$f(x) = (x_1x_1, x_1x_2, x_1x_3, x_2x_1, x_2x_2, x_2x_3, x_3x_1, x_3x_2, x_3x_3)$$

$$f(y) = (y_1y_1, y_1y_2, y_1y_3, y_2y_1, y_2y_2, y_2y_3, y_3y_1, y_3y_2, y_3y_3)$$

So,

$$f(5,6,7) = (25,30,35,30,36,42,35,42,49) \text{ and}$$

$$f(8,9,10) = (64,72,80,72,81,90,80,90,100)$$

So the dot product is,

$$f(x) \cdot f(y) = f(5,6,7) \cdot f(8,9,10) =$$

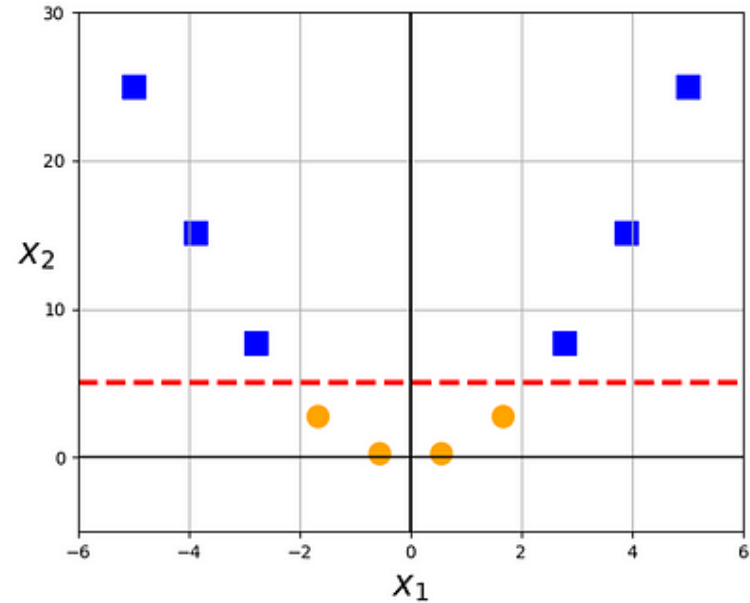
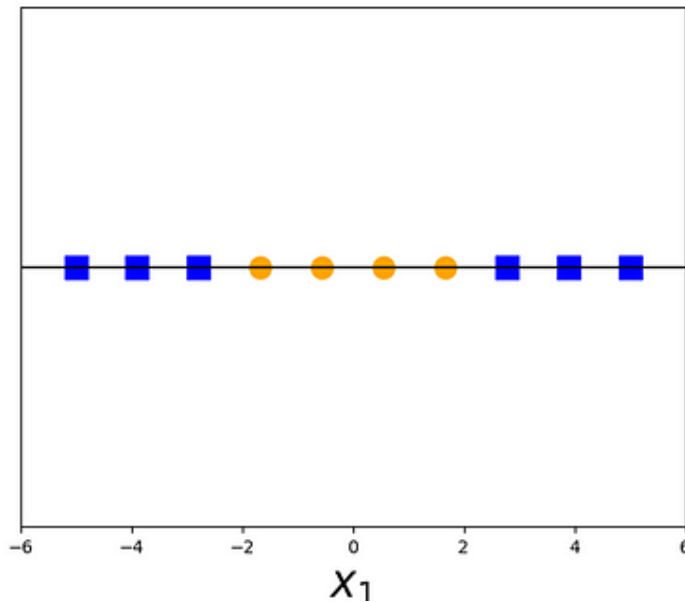
$$(1600 + 2160 + 2800 + 2160 + 2916 + 3780 + 2800 + 3780 + 4900) = 26,896$$

Using Kernel,

$$K(x, y) = (5*8 + 6*9 + 7*10) ^ 2 = (40 + 54 + 70) ^ 2 = 164*164 = 26,896$$

SVM ... Kernel Tricks

28

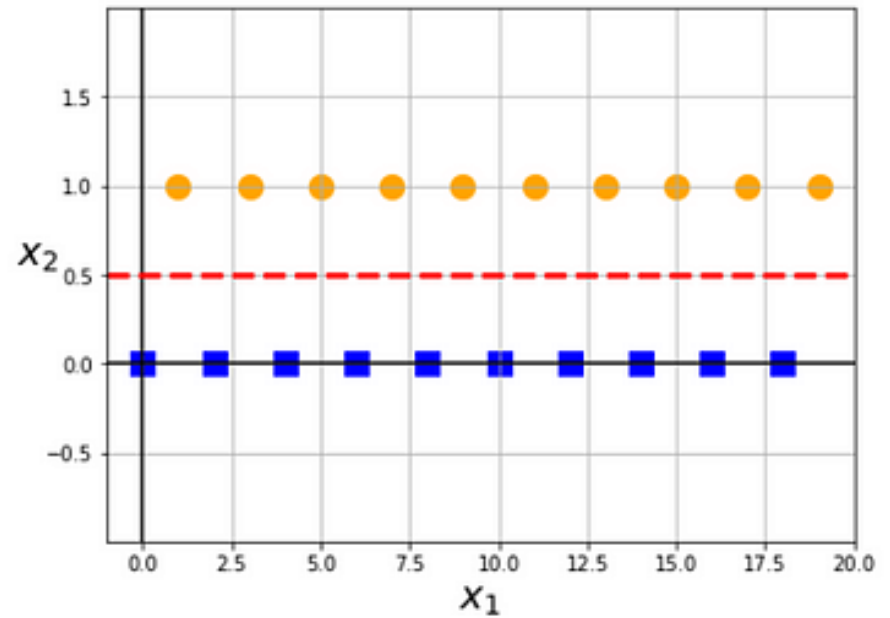
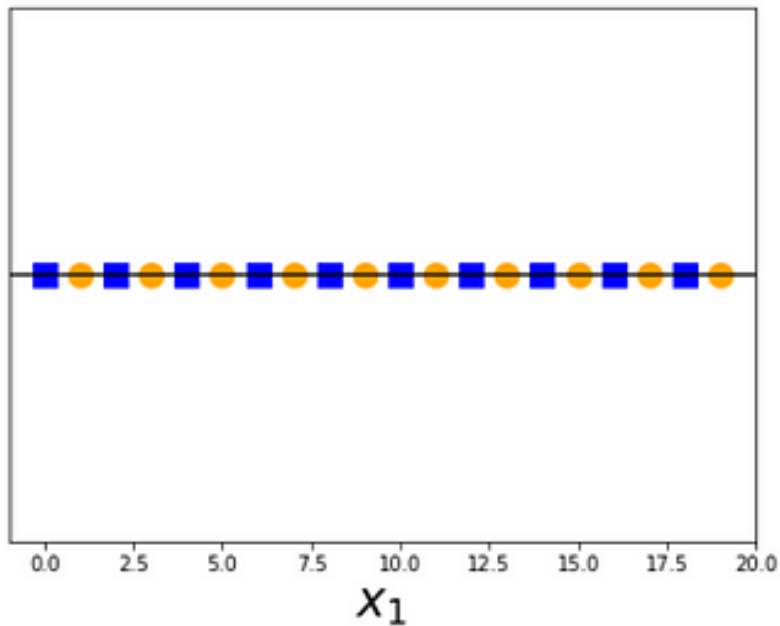


This data becomes linearly separable after a quadratic transformation to 2-dimensions.

In 1-dimension, this data is not linearly separable, but after applying the transformation $\phi(x) = x^2$ and adding this second dimension to our feature space, the classes become linearly separable.

SVM ... Kernel Tricks

29

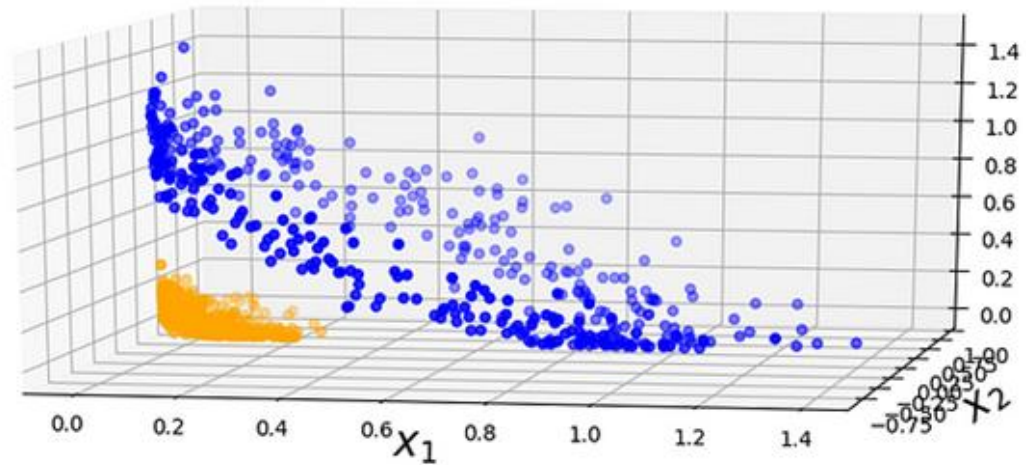
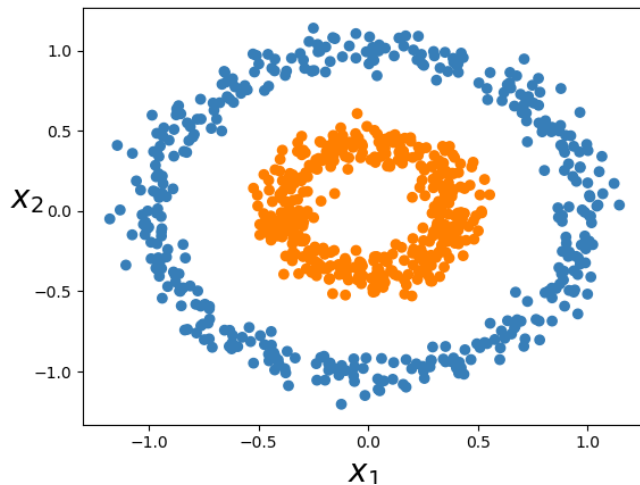


This transformation allows us to linearly separate the even and odd X_1 values in 2 dimensions.

In 1-dimension, this data is not linearly separable, but after applying the transformation $\phi(x) = x \bmod 2$ to our feature space, the classes become linearly separable.

SVM ... Kernel Tricks

30



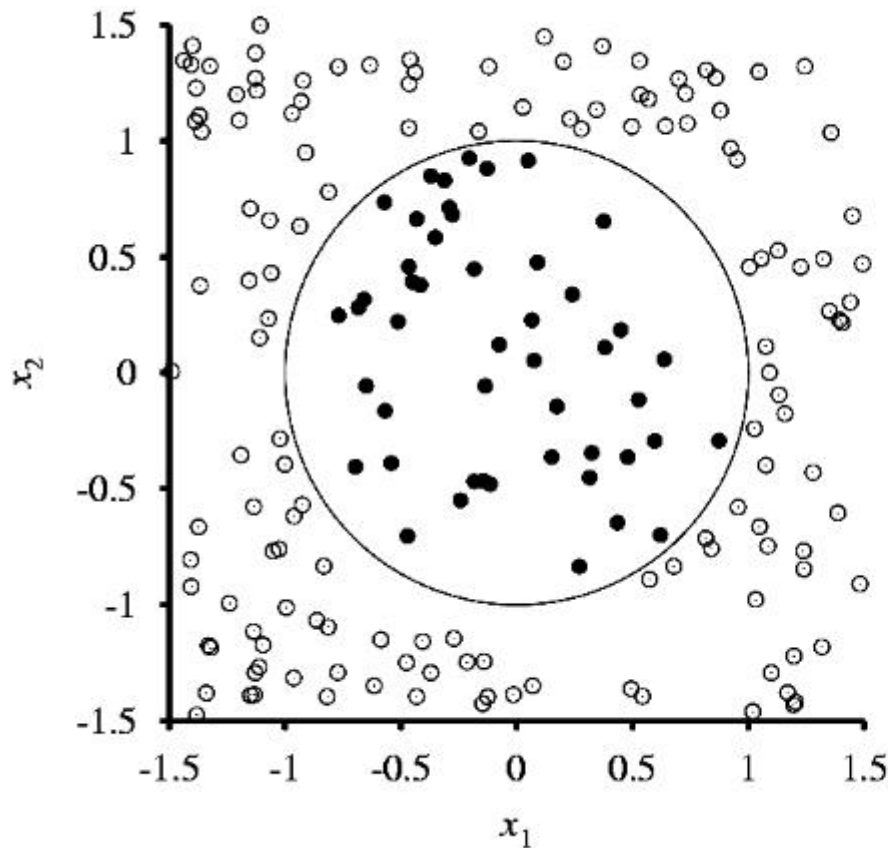
Linearly separable data in 3-d after applying the 2nd-degree polynomial transformation

In 2-dimension, this data is not linearly separable, but after applying the **second-degree polynomial transformation**, the classes become linearly separable.

$$\phi(\mathbf{x}) = \phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}$$

Support Vector Machine (SVM)

31



Example:

- Given: 2-dimensional **input space** defined by attributes $\mathbf{x} = (x_1, x_2)$.
- All positive examples ($y=+1$) are inside a circular region.
- All negative examples ($y=-1$) are outside that circular region.
- The separator is: $x_1^2 + x_2^2 \leq 1$.

⇒ **There is no linear separator!**

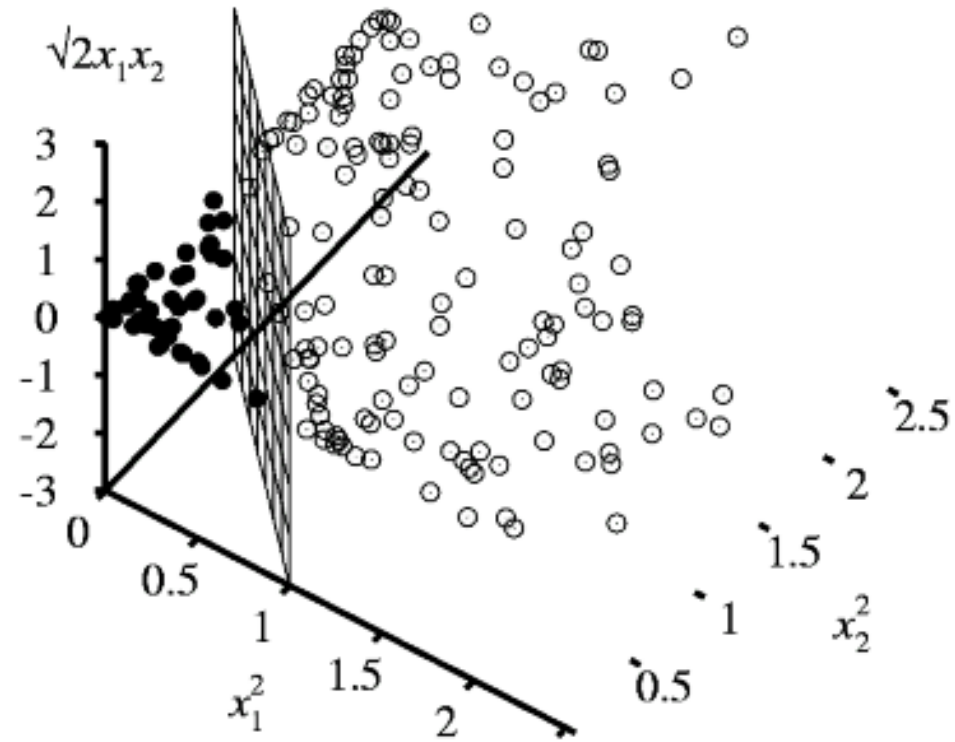
Support Vector Machine (SVM)

32

- Suppose we re-express the input data using some computed features – i.e. we map each input vector $\mathbf{x} = (x_1, x_2)$ to a new vector of feature values, $F(\mathbf{x})$.

- In particular, we use the three features $f_1 = x_1^2$, $f_2 = x_2^2$, $f_3 = \sqrt{2} x_1 x_2$.

⇒ the new vectors in the three-dimensional so-called feature space are linearly separable!



SVM ... Kernel Tricks

33

Steps involved in SVM:

- Collects the Data and plot it accordingly
- Apply the Kernel Trick
- Learns Linear Line that classifies the data
- Projects back the data

SVM ... Kernel Types

34

Linear Kernel:

- Let us say that we have two vectors, x and y , then the linear kernel is defined by the dot product of these two vectors:

$$K(x, y) = (x \cdot y)$$

SVM ... Kernel Types

35

Polynomial Kernel:

- Let us say that we have two vectors, x and y , then the polynomial kernel is defined by the dot product of these two vectors:

$$K(x, y) = (x \cdot y + 1)^d$$

where d is the degree of the polynomial

SVM ... Kernel Types

36

Gaussian RBF Kernel:

- Let us say that we have two vectors, x and y , then the Gaussian Radial Basis Function kernel is defined by the dot product of these two vectors:

$$K(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

The given sigma plays a very important role in the performance of the Gaussian kernel and should neither be overestimated nor underestimated, it should be carefully tuned according to the problem.

SVM ... Kernel Types

37

Laplacian Kernel:

- Let us say that we have two vectors with names x and y , then the Laplacian kernel is defined by the dot product of these two vectors:

$$K(x, y) = e^{-\frac{\|x-y\|}{\sigma}}$$

This type of kernel is less prone to changes and is totally equal exponential function kernel

SVM ... Kernel Types

38

Sigmoid Kernel:

- Let us say that we have two vectors, x and y , then the bipolar sigmoid function. The equation for the hyperbolic kernel function is:

$$K(x, y) = \tanh(ax^T y + c)$$

Support Vector Machine (SVM)

39

- Generally, we have the danger of overfitting!
- Therefore, we use the SVM as a maximum margin classifier.

