



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Machine Learning Approaches - SVM

*B. Solaiman*  
*Basel.solaiman@imt-atlantique.fr*  
*Départ. Image & Traitement de l'Information*  
*Brest, France*

2

## *Introduction*



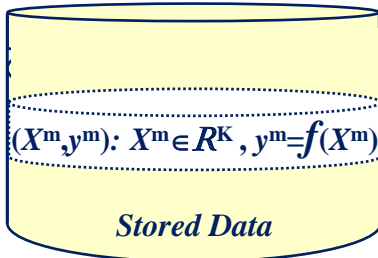
**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

*Machine Learning Approaches ----- B. Solaiman*



## Introduction

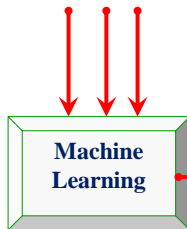
3



A *Pattern* (sample, instance, example, ...) is represented by a set of  $K$  **features**, or **attributes**, viewed as a  $K$ -dimensional feature vector  $X \in R^K$

The pattern “may” be associated with a target variable  $y = f(X)$  (in case of supervised learning)

$f$ : Unknown function associating  $X$  and  $y$



**Learning** : given the *training set*, estimate the prediction function  $\tilde{f}$  by minimizing “a” prediction error

$$\tilde{f} : X^m \in R^K, y^m = \tilde{f}(X^m)$$



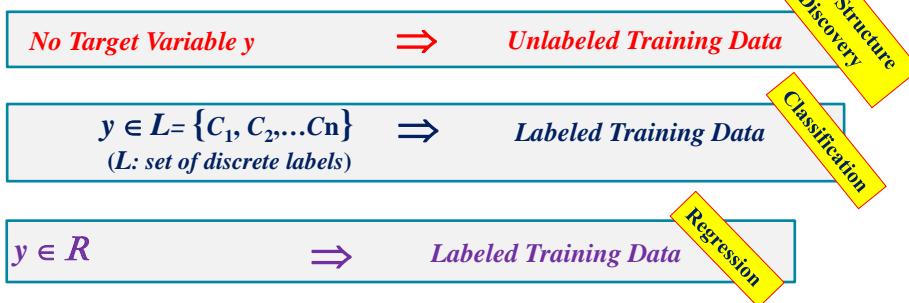
Machine Learning Approaches ----- B. Solaiman



## Introduction

4

**y**



### Regression problem

Instead of predicting the class of an input pattern, we want to predict continuous values



Machine Learning Approaches ----- B. Solaiman



## Introduction

5

### Natures of a feature



**Numerical, or quantitative, Feature**



**Symbolic, or qualitative, Feature**



**Nominal Feature:**

Feature values are symbols (= and  $\neq$  are the only operations we can conduct)



**Ordinal Feature:**

Feature values are “ordered” symbols (=,  $\neq$ ,  $<$ ,  $>$ ,  $\leq$  and  $\geq$  are the only operations we can conduct)



Machine Learning Approaches ----- B. Solaiman



## Introduction

6

### Distance between patterns



**Numerical Features**

$X = [x_1, x_2, \dots, x_n]$  &  $Y = [y_1, y_2, \dots, y_n]$   $x_i, y_i \in \mathbb{R}$

**→ Minkowski Distance**

$$d_p(X, Y) = \left\{ \sum_{i=1}^n |x_i - y_i|^p \right\}^{\frac{1}{p}}$$

**Manhattan distance :**  $d_1(X, Y) = \sum_{i=1}^n |x_i - y_i|$  ( $P = 1$ )

**Euclidian distance :**  $d_2(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$  ( $P = 2$ )

**Max distance :**  $d_\infty(X, Y) = \max_{i=1}^n |x_i - y_i|$  ( $P = \infty$ )



Machine Learning Approaches ----- B. Solaiman



## Introduction

7

### Distance between patterns



#### Nominal features

$$X = [x_1, x_2, \dots, x_n] \text{ \& } Y = [y_1, y_2, \dots, y_n]$$

$$x_i, y_i \in \Omega \subseteq \Omega_1 \times \Omega_2 \times \dots \times \Omega_n$$

$$\text{dist}(X, Y) = \frac{n - Q}{n}$$

**Q:** Number of features having the same modality in  $X$  and  $Y$



Machine Learning Approaches ----- B. Solaiman



## Introduction

8

### Distance between patterns



#### Nominal binary features

$$X = [x_1, x_2, \dots, x_n] \text{ \& } Y = [y_1, y_2, \dots, y_n] \quad x_i, y_i \in \Omega \subseteq \{0, 1\}$$

### Confusion Matrix

		Y	
		1	0
X	1	a	b
	0	c	d

$a$  : nb of feature assuming « 1 » in  $X$  and « 1 » in  $Y$

$b$  : nb of feature assuming « 1 » in  $X$  and « 0 » in  $Y$

$c$  : nb of feature assuming « 0 » in  $X$  and « 1 » in  $Y$

$d$  : nb of feature assuming « 0 » in  $X$  and « 0 » in  $Y$



Machine Learning Approaches ----- B. Solaiman



## Introduction

9

### Distance between patterns



### Nominal binary features

Case 1: Binary Symmetric Features (0 and 1 have the same importance)

#### Simple Matching Coefficient

$$\text{dist}(X, Y) = \frac{b + c}{a + b + c + d}$$

Example

	Y	
	1	0
X	1	a
	0	c
	b	d

X	1	1	1	0	1	0	0
Y	0	1	1	0	0	1	0

$$\text{dist}(X, Y) = \frac{2 + 1}{2 + 2 + 1 + 2} = 3 / 7 = 0.429$$



IMT Atlantique  
Stratégie d'Appui de la Loire  
École Mines-Télécom

Machine Learning Approaches ----- B. Solaiman



## Introduction

10

### Distance between patterns



### Nominal binary features

Case 2: Binary Non Symmetric Features  
(1 is more important than 0)

#### Jaccard' Coefficient

$$\text{dist}(X, Y) = \frac{b + c}{a + b + c}$$

Exemple

	Y	
	1	0
X	1	a
	0	c
	b	d

X	1	1	1	0	1	0	0
Y	0	1	1	0	0	1	0

$$\text{dist}(X, Y) = \frac{2 + 1}{2 + 2 + 1} = 3 / 5$$



IMT Atlantique  
Stratégie d'Appui de la Loire  
École Mines-Télécom

Machine Learning Approaches ----- B. Solaiman

## ***Support Vector Machines (SVM)***



*Machine Learning Approaches ----- B. Solaiman*

## ***Support Vector Machines (SVM)***

- ➔ *Linear Models for Classification*
- ➔ *Perceptron*
- ➔ *A Bit of Geometry*
- ➔ *Hard (Linear) Support Vector Machines (H-SVM)*
- ➔ *Soft Margin Support Vector Machines (SM-SVM)*
- ➔ *Kernel-based Support Vector Machines*



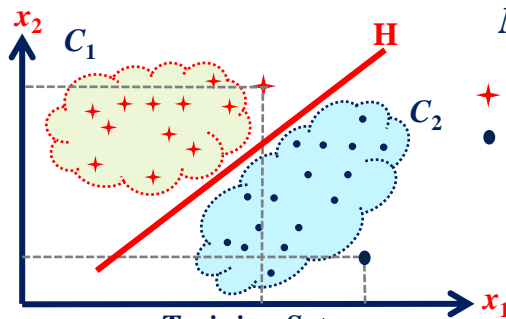
*Machine Learning Approaches ----- B. Solaiman*



## Support Vector Machines (SVM)

13

### Linear Models for Classification



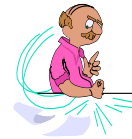
Training Set  
(Binary classification)

=  
Two classes)



Machine Learning Approaches ----- B. Solaiman

Model for Classification ?

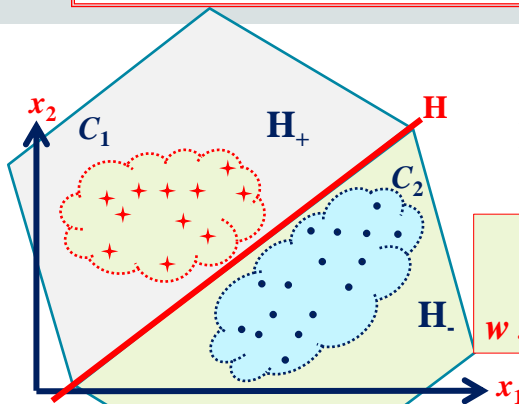


Linear Model :  
Hyperplane  
Separation Surface **H**



## Support Vector Machines (SVM)

14



IF :  $X$  "on"  $H$ :

IF :  $X \in H_+$ :

IF :  $X \in H_-$ :

How does it work ?

Hyperplane equation :

$$w \cdot X + b = 0$$

$$X = (x_1, x_2, \dots, x_n)$$

$$w = (\omega_1, \omega_2, \dots, \omega_n)$$

$$w \cdot X = x_1 \cdot \omega_1 + x_2 \cdot \omega_2 + \dots + x_n \cdot \omega_n$$

$$w \cdot X + b = 0$$

$$w \cdot X + b > 0$$

$$w \cdot X + b < 0$$

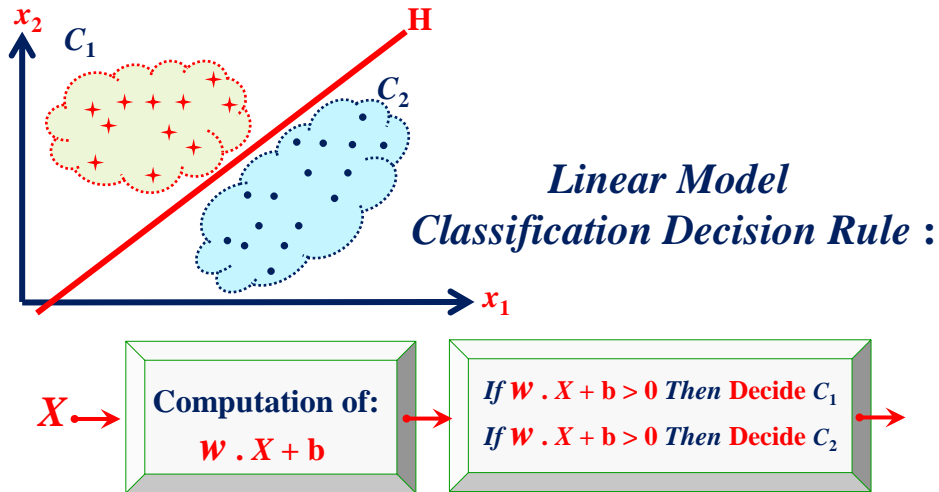


Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

15



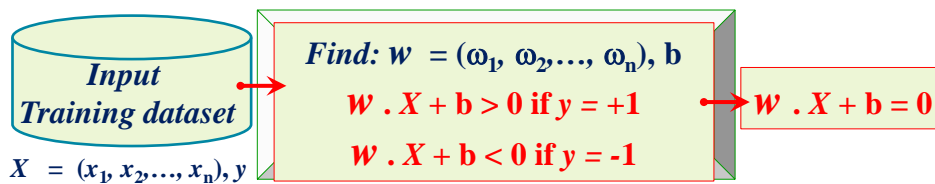
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

16

### Linear Model for Binary Classification

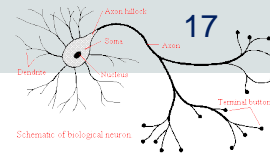


Machine Learning Approaches ----- B. Solaiman

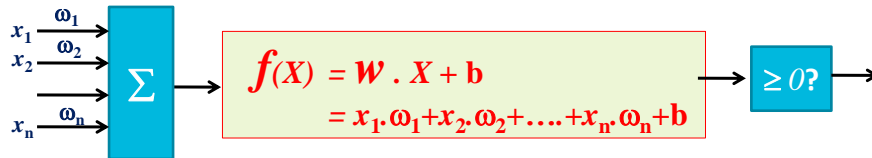




## Support Vector Machines (SVM)



### Perceptron, Rosenblatt - 1958



$$b = w_0, x_0 = 1 \rightarrow f(X) = w \cdot X \text{ (dot product)}$$

**Decision Rule:** If  $f(X) = w \cdot X > 0$  Then +1  
 If  $f(X) = w \cdot X < 0$  Then - 1



Machine Learning Approaches ----- B. Solaiman

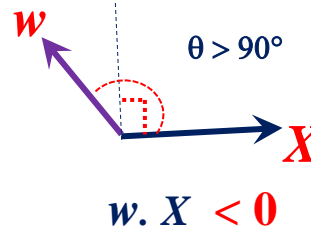
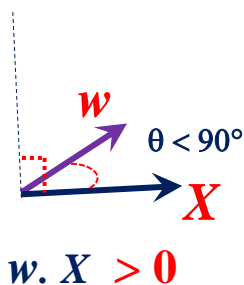


## Support Vector Machines (SVM)

18

### A Bit of Geometry

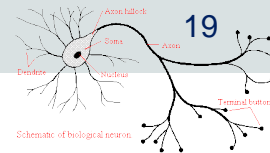
$$\begin{aligned} \text{Dot Product } \langle w, X \rangle &= w \cdot X \\ &= \|w\| \|X\| \cos \theta \\ &= x_1 \cdot w_1 + x_2 \cdot w_2 + \dots + x_n \cdot w_n \end{aligned}$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

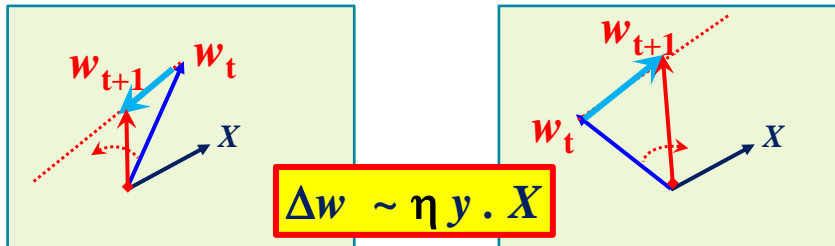


### Perceptron, Rosenblatt - 1958

**Decision Rule:** *If*  $f(X) = \mathbf{W} \cdot \mathbf{X} > 0$  *Then* +1

**Types of error:** *If*  $f(X) = \mathbf{W} \cdot \mathbf{X} < 0$  *Then* -1

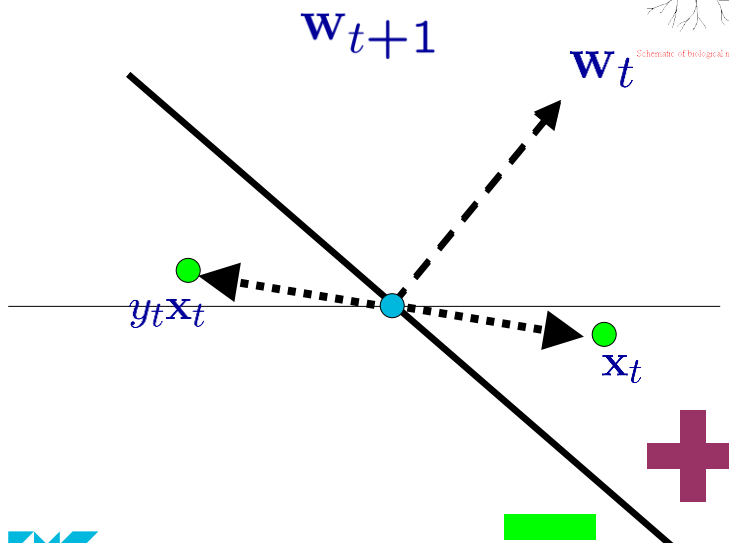
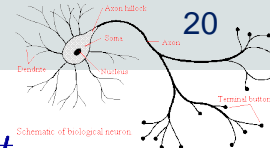
$f(X) = \mathbf{W} \cdot \mathbf{X} > 0$  *whereas*  $y = -1$      $f(X) = \mathbf{W} \cdot \mathbf{X} < 0$  *whereas*  $y = +1$



Machine Learning Approaches ----- B. Solaiman



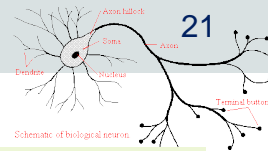
## Support Vector Machines (SVM)



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)



### Perceptron, Rosenblatt - 1958

#### Perceptron Algorithm :

Random initialization of  $(\omega_1, \omega_2, \dots, \omega_n), b$

Pick training samples  $X$  one by one

Predict class of  $X$  using current  $(\omega_1, \omega_2, \dots, \omega_n), b$

$$y' = \text{Sign}(W \cdot X)$$

If  $y'$  is correct (i.e.,  $y = y'$ ) : **No change:**  $w_{t+1} = w_t$

If  $y'$  is wrong: **Adjust**  $w_t$  :  $w_{t+1} = w_t + \eta \cdot y_t \cdot X_t$

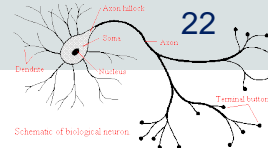
- $\eta$  is the learning rate parameter
- $X_t$  is the t-th training example
- $y_t$  is true t<sup>th</sup> class label ( $\{+1, -1\}$ )



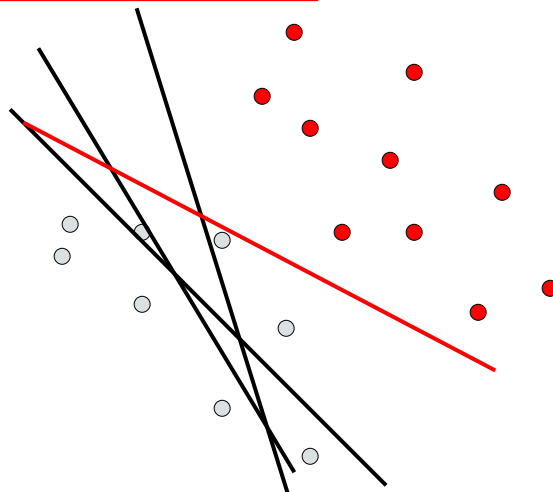
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)



### Perceptron, Rosenblatt - 1958



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

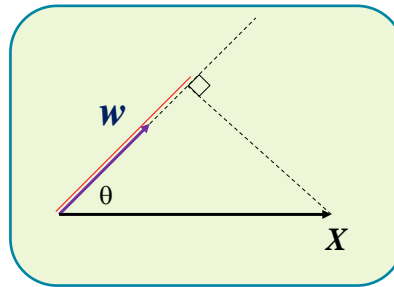
23

### A Bit of Geometry

**Projection of  $X$  onto  $w$**

$$= \|X\| \cos\theta$$

$$= w \cdot X / \|w\|$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

24

### A Bit of Geometry

### Parallel Hyperplanes

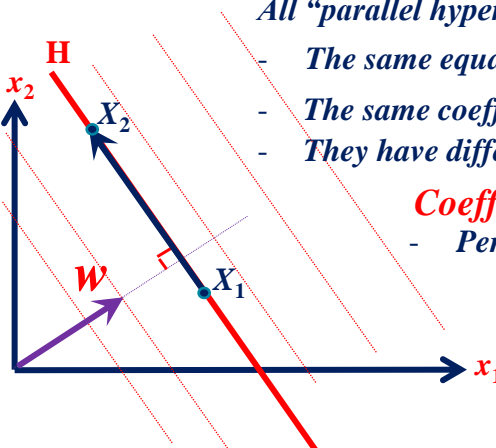
All “parallel hyperplanes” have:

- The same equation  $w \cdot X + b = 0$ ;
- The same coefficients vector  $w$ ;
- They have different origin translation values  $b$

### Coefficients vector : $w$

- Perpendicular to  $H$  ?

$$\begin{aligned} w \cdot (X_2 - X_1) &= w \cdot (X_2 - X_1) \\ &= w \cdot X_2 - w \cdot X_1 \\ &= -b - (-b) = 0 \end{aligned}$$



Machine Learning Approaches ----- B. Solaiman

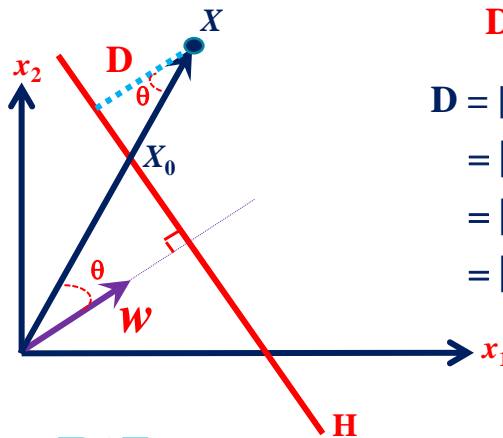


## Support Vector Machines (SVM)

25

### A Bit of Geometry

*Distance between an instance  $X$  and  $H$*



$$D = \text{dist}(X, H) = ??$$

$$\begin{aligned} D &= \|X_0X\| \cos\theta \\ &= \|w \cdot X_0X\| / \|w\| \\ &= \|w \cdot X - w \cdot X_0\| / \|w\| \\ &= \|w \cdot X + b\| / \|w\| \end{aligned}$$



Machine Learning Approaches ----- B. Solaiman

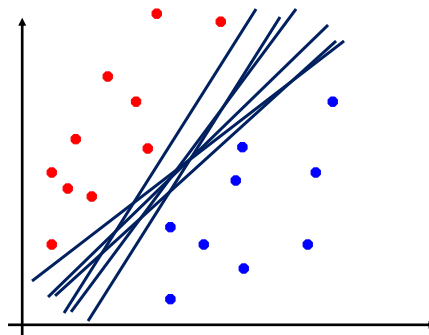


## Support Vector Machines (SVM)

26

### Hard (Linear) Support Vector Machines (H-SVM)

*Which of the linear separators is optimal?*



Machine Learning Approaches ----- B. Solaiman

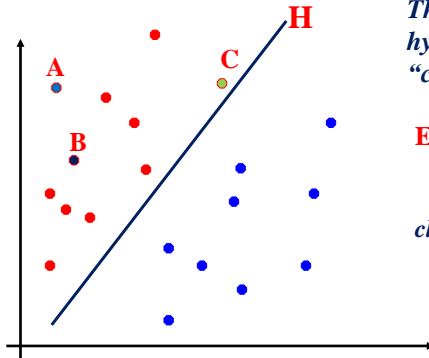


## Support Vector Machines (SVM)

27

### Hard (Linear) Support Vector Machines (H-SVM)

#### Largest Margin Concept:



The distance from the separating hyperplane,  $H$ , corresponds to the “confidence” of decision

#### Example:

We are more sure and confident about the class of **A** and **B** than of **C**



Machine Learning Approaches ----- B. Solaiman

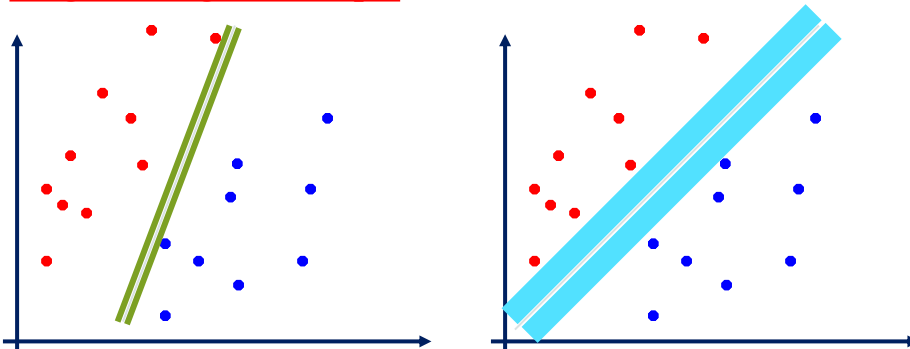


## Support Vector Machines (SVM)

28

### Hard (Linear) Support Vector Machines (H-SVM)

#### Largest Margin Concept:



The **margin** of a linear classifier is the width that the boundary could be increased by before hitting a datapoint instance



Machine Learning Approaches ----- B. Solaiman

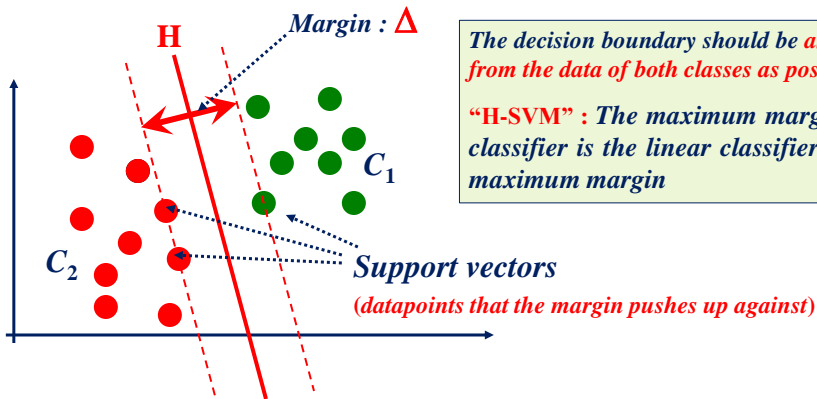


## Support Vector Machines (SVM)

29

### Hard (Linear) Support Vector Machines (H-SVM)

#### Largest Margin Concept :



Separating hyperplane :  $w \cdot X + b = 0$   
Machine Learning Approaches ----- B. Solaiman

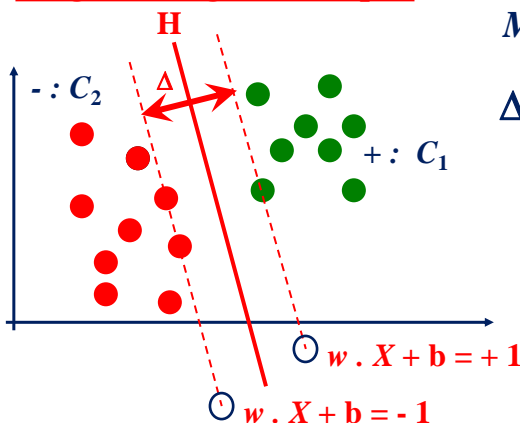


## Support Vector Machines (SVM)

30

### Hard (Linear) Support Vector Machines (H-SVM)

#### Largest Margin Concept :



Margin width :  $\Delta = ?$

$$\Delta = 2 \cdot |w \cdot X + b| / \|w\|$$



$$\Delta = 2 / \|w\|$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

31

### Hard (Linear) Support Vector Machines (H-SVM)

Vapnik, 1965; Vapnik, 1995 :

$$\{(X^m, y^m), m = 1, \dots, M, X^m \in \mathbb{R}^n, y^m \in \{+1, -1\}\}$$

**H-SVM formulation : Find  $w$  &  $b$  such that**

**Maximize :**

$$\Delta = 2 / \|w\|$$

**Subject to Support vectors constraints :**

$$y^m \cdot (w \cdot X^m + b) \geq +1 \quad \text{for all } m$$

**Minimize:**

$$\|w\|^2 / 2$$

**Quadratic Optimization problem, subject to linear constraints**



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

32

### 1D Example

Dataset :  $\{(-3, -1), (-1, -1), (+2, +1)\}$

**H :** All lines with a cross point between -1 and +1

**Margin constraints:**  $y^m \cdot (w \cdot X^m + b) \geq +1, m = 1, \dots, M :$

$$a \cdot (-3) + b < -1 \rightarrow b < 3a - 1$$

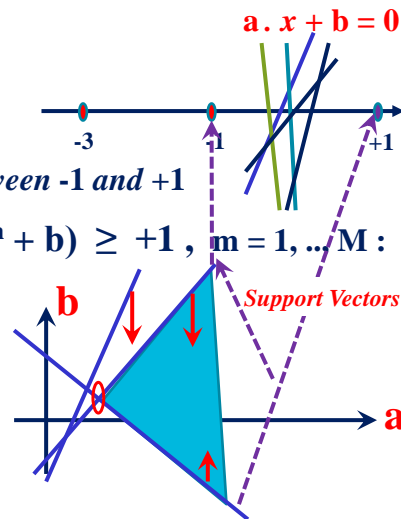
$$a \cdot (-1) + b < -1 \rightarrow b < a - 1$$

$$a \cdot (+2) + b > +1 \rightarrow b > -2a + 1$$

**Minimize**

$$\|w\|^2 = \|a\|^2 \rightarrow a = 0.66$$

$$\rightarrow b = 0.33$$



Machine Learning Approaches ----- B. Solaiman





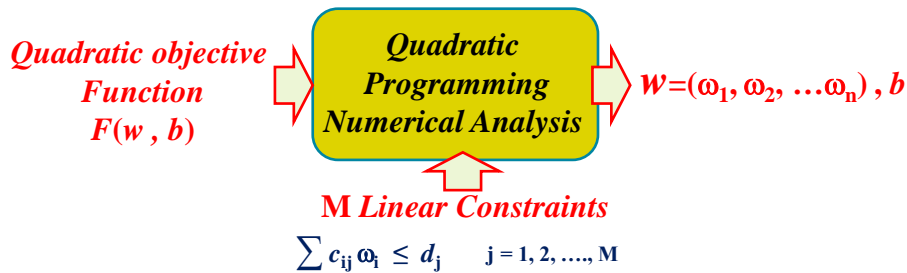
## Support Vector Machines (SVM)

33

### Hard (Linear) Support Vector Machines (H-SVM)

#### RAPPEL: QUADRATIC PROGRAMMING PROBLEM

$$F(w, b) = b + \sum a_i \omega_i + \sum \sum q_{ij} \omega_i \omega_j \quad w = (\omega_1, \omega_2, \dots, \omega_n), b$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

34

### Hard (Linear) Support Vector Machines (H-SVM)

#### Optimization using Lagrange Multipliers

$$\text{Lagrangian : } L = \|w\|^2 / 2 + \sum_{m=1, \dots, M} \lambda_m \{ 1 - [y^m \cdot (w \cdot X^m + b)] \}$$

*Lagrange Multipliers*

The constraints are added into the optimization by adding extra variables (called Lagrange multipliers)

$$\|w\|^2 = (\omega_1)^2 + (\omega_2)^2 + \dots + (\omega_n)^2$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

35

### Hard (Linear) Support Vector Machines (H-SVM)

#### Optimization using Lagrange Multipliers

**Objective Function:** 
$$L = \|w\|^2 / 2 + \sum_{m=1, \dots, M} \lambda_m \{ 1 - [y^m \cdot (w \cdot X^m + b)] \}$$

Setting the **gradient of  $L$  w.r.t.  $w$  and  $b$  to zero**, we have

$$\frac{\partial L}{\partial w} = w - \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m = 0$$

Linear sum  
of samples

$$w^* = \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m$$

$$\frac{\partial L}{\partial b} = - \sum_{m=1, \dots, M} \lambda_m y^m = 0$$

$$\sum_{m=1, \dots, M} \lambda_m y^m = 0$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

36

$$L = \|w\|^2 / 2 + \sum_{m=1, \dots, M} \lambda_m \{ 1 - [y^m \cdot (w \cdot X^m + b)] \} \quad w^* = \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m$$



$$L = (1/2) \cdot \left[ \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m \right] \cdot \left[ \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m \right] - \left[ \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m \right] \cdot \left[ \sum_{k=1, \dots, M} \lambda_k y^k \cdot X^k \right] - \sum_{m=1, \dots, M} \lambda_m y^m \cdot b + \sum_{m=1, \dots, M} \lambda_m$$



**The optimization problem becomes : (called : Dual Problem)**

Find  $\lambda_m$  such that:

$$L = - (1/2) \cdot \sum_{m=1, \dots, M} \sum_{k=1, \dots, M} \lambda_m \lambda_k y^m y^k \cdot [X^m \cdot X^k] + \sum_{m=1, \dots, M} \lambda_m$$

is  
maximized

where

$$\sum_{m=1, \dots, M} \lambda_m y^m = 0$$

$$\lambda_m \geq 0, m = 1, \dots, M$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

37

$$L = - (1/2) \cdot \sum_{m=1, \dots, M} \sum_{k=1, \dots, M} \lambda_m \lambda_k y^m y^k \cdot [X^m, X^k] + \sum_{m=1, \dots, M} \lambda_m$$

*Quadratic optimization problems are a well-known mathematical programming problems for which several algorithms exist;*

*Once resolved, most of  $\lambda_m$  are 0, only a small number have  $\lambda_m > 0$ ; the corresponding  $X^m$ s are the support vectors*



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

38

*Given a solution  $\lambda_1 \dots \lambda_M$  to the dual problem, solution to the primal (i.e., coefficients vector  $w^*$  and origin shift factor  $b^*$  solution) is given by :*

$$w^* = \sum_{m \in \text{Support Vectors}} \lambda_m y^m \cdot X^m$$

Learned weight
Support Vectors

$$b^*_{(k)} = y^k - \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m X^k$$

For any  $k$  such that  $\lambda_k > 0$



$$b^* = (1/K) \sum_k b^*_{(k)} \quad (K \text{ denotes the nb of SVs})$$



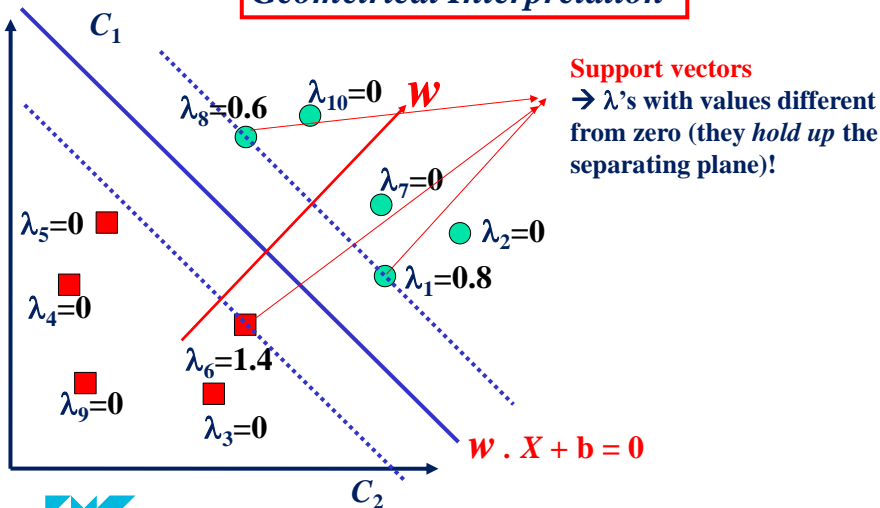
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

39

### Geometrical Interpretation



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

40



**Decision rule : Given a new data instance  $X$**

Then, the classifying function is  
 (note that we don't need  $w^*$  explicitly):

$$f(X) = \sum_{m \in \text{Support Vectors}} \lambda_m y^m \cdot X^m \cdot X + b^*$$



**If  $f(X) > 0$  Then decide  $C_1$ ;  
 Else decide  $C_2$**



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

41



Notice that  $f(X)$  relies on the dot product between the unknown instance  $X$  and the support vectors  $X^m$



Also keep in mind that solving the optimization problem involved computing the dot products  $X^m X^k$  between all training data instances



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

42

### Soft Margin Support Vector Machines (SM-SVM)

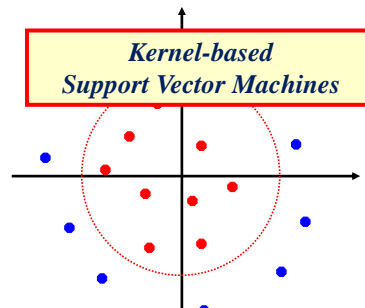
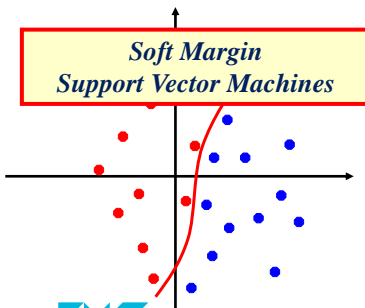
#### Non Linear Separability

Slightly non separable

Highly non separable

Soft Margin  
Support Vector Machines

Kernel-based  
Support Vector Machines



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

43

### Soft Margin Support Vector Machines (SM-SVM)

To allow errors in data, we relax the margin constraints by introducing slack variables,  $\xi_m (\geq 0)$  :

$$w \cdot X^m + b \geq +1 - \xi_m, \quad \text{if } y^m = +1$$

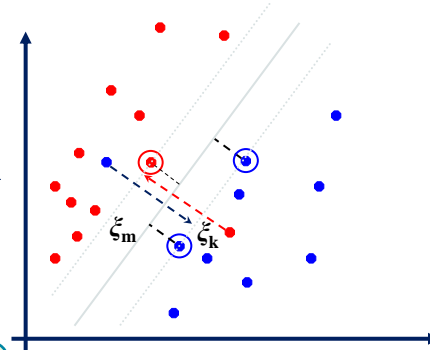
$$w \cdot X^m + b \leq -1 + \xi_m, \quad \text{if } y^m = -1$$



**The new constraints:**

$$y^m \cdot (w \cdot X^m + b) \geq 1 - \xi_m, \quad m=1, \dots, M$$

$$\xi_m \geq 0, \quad m=1, \dots, M$$



$\xi_m$ : Measure of margin violation by an erroneous instance  $X^m$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

44



Assigning an “extra cost” for errors in the objective function that we are going to optimize:

**Minimize:** 
$$L = \|w\|^2 / 2 + C \sum_{m=1, \dots, M} \xi_m$$

**Subject to:** 
$$y^m \cdot (w \cdot X^m + b) \geq 1 - \xi_m, \quad m=1, \dots, M$$
  

$$\xi_m \geq 0, \quad m=1, \dots, M$$

$$\sum_{m=1, \dots, M} \xi_m$$
 : Total margin violation

**C** : Large C means a higher penalty to errors



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

45

### Optimization using Lagrange Multipliers

$$L = \underbrace{\|w\|^2/2}_{\text{Minimizing}} + \underbrace{C \sum_{m=1, \dots, M} \xi_m}_{\text{Total margin violation error}} - \sum_{m=1, \dots, M} \underbrace{\lambda_m \{y^m \cdot (w \cdot X^m + b) - 1 + \xi_m\}}_{\lambda_m, \mu_m \geq 0 \text{ Lagrange Multipliers}} - \sum_{m=1, \dots, M} \underbrace{\mu_m \xi_m}_{\text{Tradeoff (i.e., relative importance) between error and margin}}$$

↔  
Maximizing the margin

**Objective function**  
to be minimized w.r.t.  $w$ ,  $b$  &  $\xi_m$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

46

$$L = \|w\|^2/2 + C \sum_{m=1, \dots, M} \xi_m - \sum_{m=1, \dots, M} \lambda_m \{y^m \cdot (w \cdot X^m + b) - 1 + \xi_m\} - \sum_{m=1, \dots, M} \mu_m \xi_m$$

Setting the **gradient of  $L$  w.r.t.  $w$ ,  $b$  and  $\xi_m$  to zero**, we have :

$$\left. \begin{aligned} \frac{\partial L}{\partial w} &= w - \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m = 0 \\ \frac{\partial L}{\partial b} &= - \sum_{m=1, \dots, M} \lambda_m y^m = 0 \\ \frac{\partial L}{\partial \xi_m} &= C - \lambda_m - \mu_m = 0 \end{aligned} \right\} \begin{array}{l} \text{Same solution} \\ \text{as for the H-SVM} \end{array}$$

$\mu_m \geq 0$   
 $\lambda_m \geq 0$   $\Rightarrow$   $\lambda_m \leq C$   
 (Otherwise,  $C - \lambda_m \leq 0$ ,  
 thus no  $\mu_m$  will verify:  
 $C - \lambda_m - \mu_m = 0$ )



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

47



$$w^* = \sum_{m=1, \dots, M} \lambda_m y^m \cdot X^m$$

$$\sum_{m=1, \dots, M} \lambda_m y^m = 0$$

with the only difference  
w.r.t (H-SVM):  
 $0 \leq \lambda_m \leq C$



**Decision rule :** *Given a new data instance  $X$*   
Then, the classifying function is

$$f(X) = \sum_{m \in \text{Support Vectors}} \lambda_m y^m \cdot X^m \cdot X + b^*$$



If  $f(X) > 0$  Then **decide**  $C_1$ ;  
Else **decide**  $C_2$



Machine Learning Approaches ----- B. Solaiman

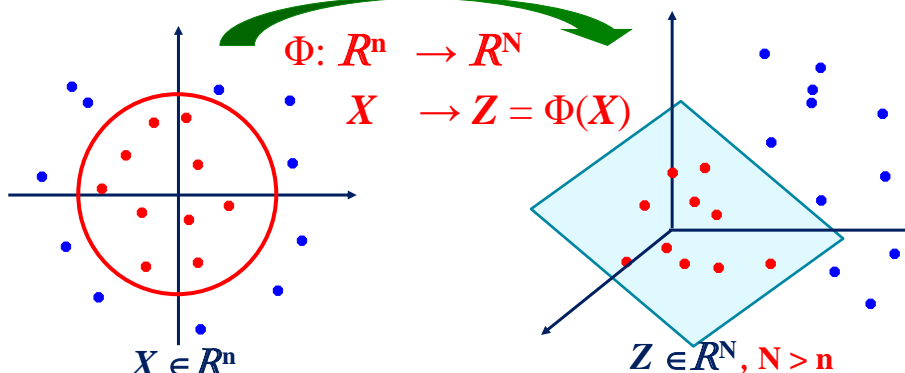


## Support Vector Machines (SVM)

48

### Kernel-based Support Vector Machines

**General idea:** the original feature space is mapped to some higher-dimensional feature space where the training set is linearly separable



Machine Learning Approaches ----- B. Solaiman

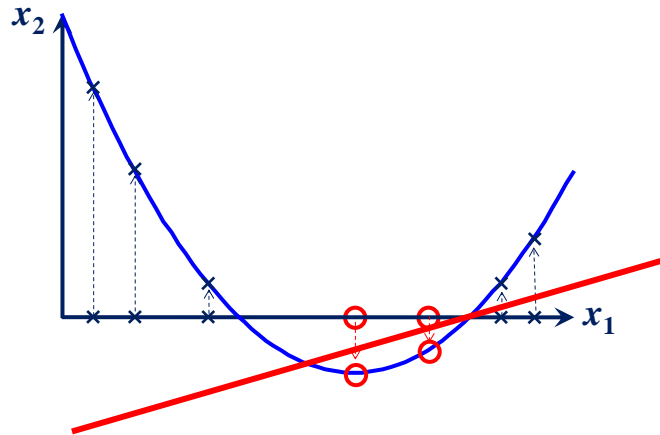




## Support Vector Machines (SVM)

49

### Kernel-based Support Vector Machines



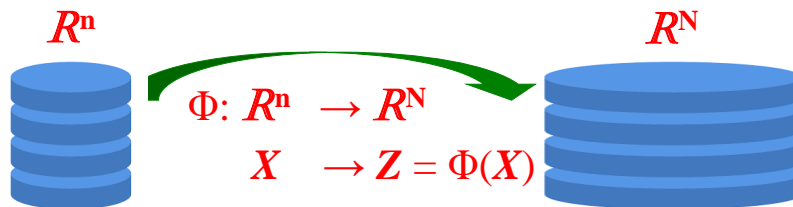
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

50

### Kernel-based Support Vector Machines



*Nonlinearly separable dataset*

$$\{(X^m, y^m), X^m \in R^n, y^m \in \{+1, -1\}\}$$

*Linearly separable dataset*

$$\{(Z^m = \Phi(X^m), y^m), Z^m \in R^N, y^m \in \{+1, -1\}\}$$



*Do we need to know the explicit  
expression of  $Z = \Phi(X)$  ?*

**SVM Machinery**



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

51

### The Kernel trick

*Lagrangian to maximize*

$$L = - (1/2) \cdot \sum_{m=1, \dots, M} \sum_{k=1, \dots, M} \lambda_m \lambda_k y^m y^k \cdot \boxed{Z^m \cdot Z^k} + \sum_{m=1, \dots, M} \lambda_m$$

*Constraints*

$$\sum_{m=1, \dots, M} \lambda_m y^m = 0 \quad \lambda_m \geq 0, m = 1, \dots, M$$

*Decision Rule*

$$f(X) = \sum_{m \in \text{Support Vectors}} \lambda_m y^m \cdot \boxed{Z^m \cdot Z} + b^*$$

*Bias term*

$$b^* = y^k - \sum_{m=1, \dots, M} \lambda_m y^m \boxed{Z^m \cdot Z^k}$$



*We only deal with Z as far as Z-dot product is concerned*



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

52

### The Kernel trick

**SVM Machinery**  
(Quadratic Programming,  
Decision rule) needs

*Z - Transformation*

*Z-dot product*

$$\begin{matrix} X^m \\ X^{m'} \end{matrix} \rightarrow \boxed{K(X^m, X^{m'}) = Z^m \cdot Z^{m'}} \rightarrow \text{Z-dot product}$$

*without explicit computation of  $\Phi(X^m)$  and  $\Phi(X^{m'})$*

**$K(.,.)$  : is called the Transformation Kernel**



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

53

### Example 1

$$X = (x_1, x_2) \dots\dots\dots : X \in \mathbb{R}^2$$

$$\Phi(X) = (1, x_1, x_2, (x_1)^2, (x_2)^2, x_1 \cdot x_2) :$$

Full  $2^d$  order polynomial transformation



$$K(X, X') = Z \cdot Z'$$

$$= 1 + x_1 x'_1 + x_2 x'_2 + (x_1)^2 (x'_1)^2 + (x_2)^2 (x'_2)^2 + x_1 \cdot x'_1 \cdot x_2 x'_2$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

54

### Example 2 $X = (x_1, x_2) \dots\dots\dots : X \in \mathbb{R}^2$

Computing the Kernel without transforming  $X$  and  $X'$  ?

$$K(X, X') = (1 + X \cdot X')^2$$

$$K(X, X') = (1 + X \cdot X')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2$$

$$= 1 + (x_1)^2 (x'_1)^2 + (x_2)^2 (x'_2)^2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2$$

This is the inner product of  $Z = \Phi(X)$  and  $Z' = \Phi(X')$  where

$$\Phi(X) = [1, (x_1)^2, (x_2)^2, \sqrt{2} x_1, \sqrt{2} x_2, \sqrt{2} x_1 x_2]$$

$$\Phi(X') = [1, (x'_1)^2, (x'_2)^2, \sqrt{2} x'_1, \sqrt{2} x'_2, \sqrt{2} x'_1 x'_2]$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

55

### Example 2

*Think about the computation complexity if :*

- $X = (x_1, x_2, \dots, x_d)$  ..... :  $X \in \mathbb{R}^d$ ,  $d \gg 2$  and
- $K(X, X') = (1 + X \cdot X')^Q$ ,  $Q \gg 2$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

56

### Example 3 (Radial Basis Function Kernel)

Show that  $K(X, X') = e^{-\alpha \|X - X'\|^2}$   
corresponds to an “inner product” in an infinite  
dimensional Z-Space (i.e., transforming  $X$  instances)

### Hint

For  $d = 1$  ( $X = x$ : i.e., scalar),  $\alpha = 1$ :

- Expand  $e^{-\alpha \|X - X'\|^2}$
- Expand cross exponential term using Taylor series



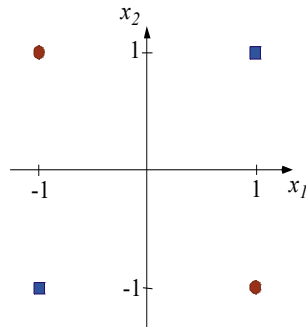
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

57

### Example 4 XOR



Index $i$	$X$	$y$
1	(1,1)	1
2	(1,-1)	-1
3	(-1,-1)	1
4	(-1,1)	-1

$$K(X, X') = (1 + X \cdot X')^2$$



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

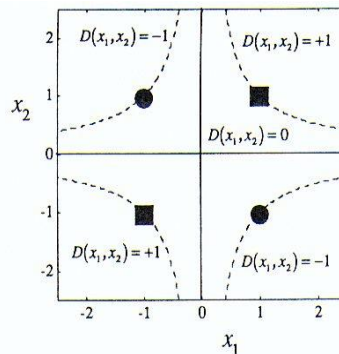
58

### Example 4 XOR

Optimal Lagrange multipliers:  $\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = 1/8$



The 4 instance data are SV



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

59

### Example 5 Face Detection application

#### Training images database:

266 pictures: 150 faces + 116 non-faces



#### *Preprocessing*

- Gray scale transformation
- Histogram equalization
- Adjust resolution to 30x40 pixel

#### *Training the SVM*

*based on the 266 training instances, a polynomial kernel*



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

60

### Example 5 Face Detection application

#### Test phase:

**Given an input image:**

- - Moving 30x40 pixel sub window over the input image
- - Histogram equalization of a sub window
- - Classification by SVM
- - Removing intersections



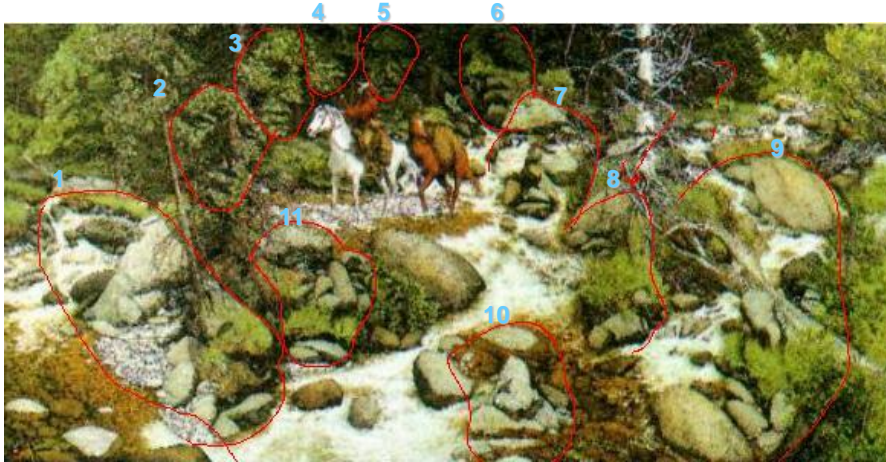
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

61

*How many faces ?*



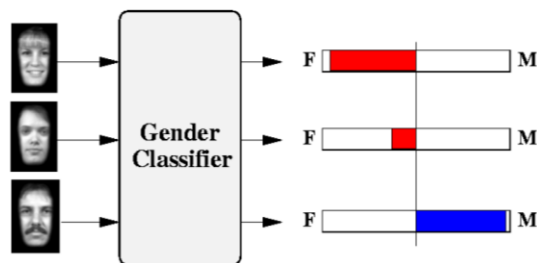
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

62

Example 6 Learning gender with SVMs



*Processed  
faces*



Machine Learning Approaches ----- B. Solaiman



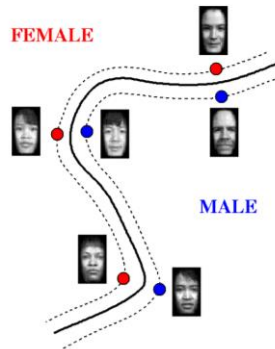
## Support Vector Machines (SVM)

63

### Example 6 Learning gender with SVMs

**Training dataset:** 1044 males, 713 females

Experiment with various kernels, select Gaussian RBF



Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

64

### Example 6 Learning gender with SVMs

Classifier	Error Rate		
	Overall	Male	Female
<b>SVM with RBF kernel</b>	<b>3.38%</b>	<b>2.05%</b>	<b>4.79%</b>
<b>SVM with cubic polynomial kernel</b>	<b>4.88%</b>	<b>4.21%</b>	<b>5.59%</b>
Large Ensemble of RBF	5.54%	4.59%	6.55%
Classical RBF	7.79%	6.89%	8.75%
Quadratic classifier	10.63%	9.44%	11.88%
Fisher linear discriminant	13.03%	12.31%	13.78%
Nearest neighbor	27.16%	26.53%	28.04%
Linear classifier	58.95%	58.47%	59.45%



Machine Learning Approaches ----- B. Solaiman





## Support Vector Machines (SVM)

65

### Example 7 Two spirals benchmark problem

#### Carnegie Mellon AI Repository

$$\alpha = \frac{i\pi}{104\gamma} \quad i = 1, 2, \dots, n \quad n : \text{num of patterns}$$

$$\gamma : \text{density}$$

$$x = R \cos \alpha \quad R : \text{radius}$$

$$y = R \sin \alpha$$

$$x^- = -x^+ \quad y^- = -y^+$$

#### Data generation :

$$(x, y) \quad \text{with} \quad R = 3.5 \quad \gamma = 1.0 \quad n^- = 100 \quad n^+ = 100$$



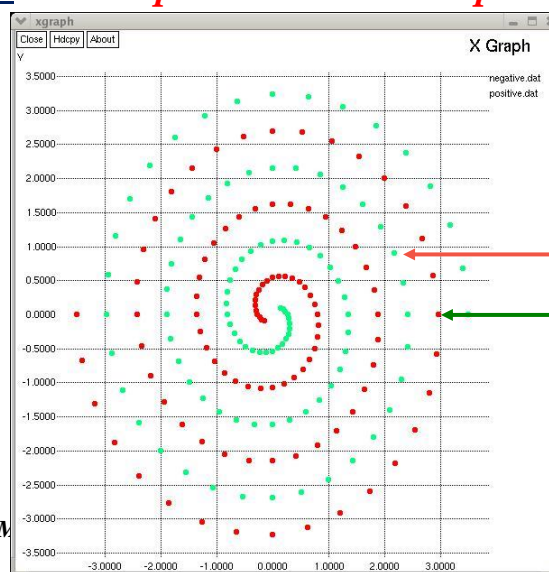
Machine Learning Approaches ----- B. Solaiman



## Support Vector Machines (SVM)

66

### Example 7 Two spirals benchmark problem



Class  
-1

Class  
+1





## Support Vector Machines (SVM)

67

### Example 7 *Two spirals benchmark problem*

