



# **The Introduction To Artificial Intelligence**

**Yuni Zeng [yunizeng@zstu.edu.cn](mailto:yunizeng@zstu.edu.cn)  
2024-2025-1**

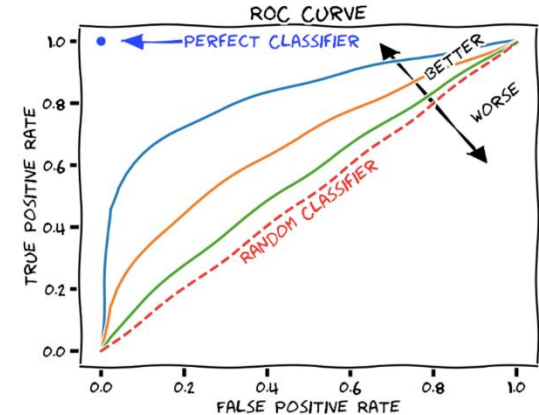
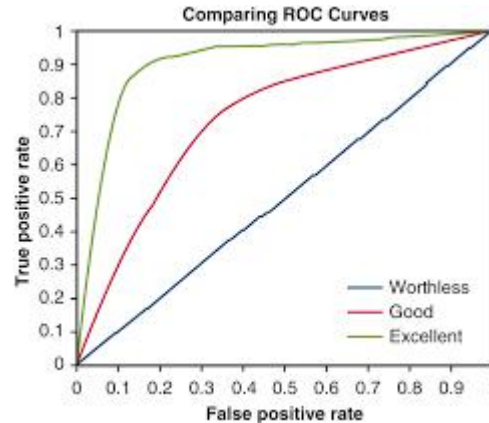
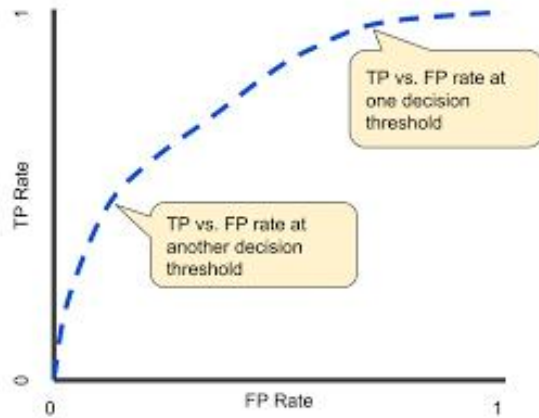
# A brief review

---

- ❑ How to make a model convincing?
  - Error, Training error, Generalization error
  - Overfitting and Underfitting
  - Evaluation Methods: Hold-out method, Cross Validation, Bootstrapping
- ❑ How to evaluate a model?
  - Measure metrics: ACC, Recall, F1, AUC...

# 1.3 Performance Measure

## ❑ ROC Curve (Receiver Operating Characteristic)



- An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds.
- TPR – FPR:
  - TPR: True positive rate
  - FPR: False positive rate

# Test

## ROC Curve (Receiver Operating Characteristic)

样本编号 (No.)	真实标签 (True label)	模型输出概率 (output probability)	样本编号 (No.)	真实标签 (True label)	模型输出概率 (output probability)
1	p	0.9	11	p	0.4
2	p	0.8	12	n	0.39
3	n	0.7	13	p	0.38
4	p	0.6	14	n	0.37
5	p	0.55	15	n	0.36
6	p	0.54	16	n	0.35
7	n	0.53	17	p	0.34
8	n	0.52	18	n	0.33
9	p	0.51	19	p	0.30
10	n	0.505	20	n	0.10

• p : positive sample, n: negative sample

# Test

## □ ROC Curve (Receiver Operating Characteristic)

- P: number of positive samples; TP: number of true positive samples
- N: number of negative samples; FP: number of false positive samples

Thresholds	0.9	0.8	0.7	0.6	0.55	0.54	0.53	0.52	0.51	0.505
TPR										
FPR										

Thresholds	0.4	0.39	0.38	0.37	0.36	0.35	0.34	0.33	0.30	0.10
TPR										
FPR										

# Test

## ❑ ROC Curve (Receiver Operating Characteristic)

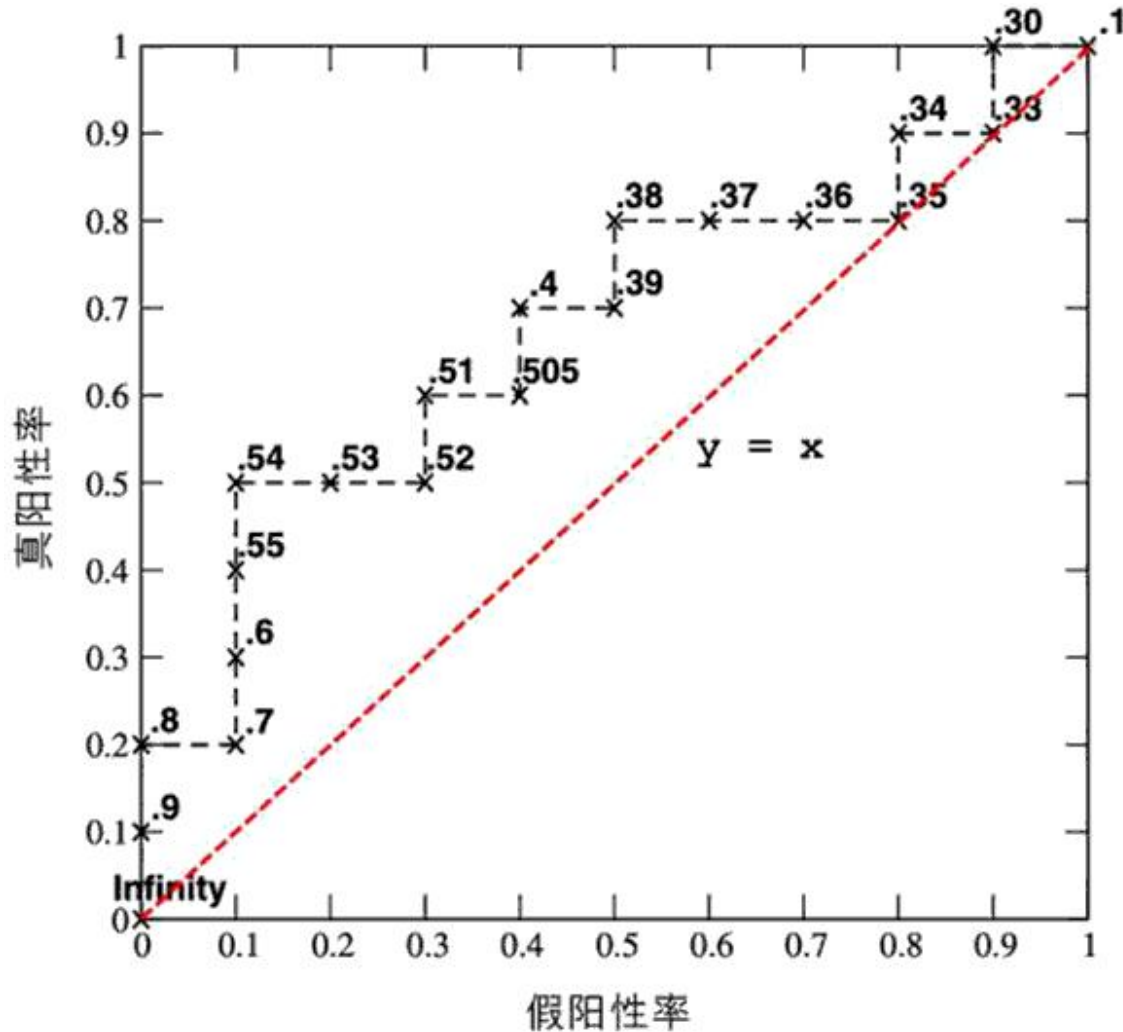
- P: number of positive samples; TP: number of true positive samples
- N: number of negative samples; FP: number of false positive samples

Thresholds	0.9	0.8	0.7	0.6	0.55	0.54	0.53	0.52	0.51	0.505
TPR = TP/P	0.1	0.2	0.2	0.3	0.4	0.5	0.5	0.5	0.6	0.6
FPR = FP/N	0	0	0.1	0.1	0.1	0.1	0.2	0.3	0.3	0.4

Thresholds	0.4	0.39	0.38	0.37	0.36	0.35	0.34	0.33	0.30	0.10
TPR = TP/P	0.7	0.7	0.8	0.8	0.8	0.8	0.9	0.9	1.0	1.0
FPR = FP/N	0.4	0.5	0.5	0.6	0.7	0.8	0.8	0.9	0.9	1.0

# Test

## ROC Curve (Receiver Operating Characteristic)



# The Introduction to Artificial Intelligence

- Part I Brief Introduction to AI & Different AI tribes
- Part II Knowledge Representation & Reasoning
- Part III AI GAMES and Searching
- Part IV Model Evaluation and Selection

✚ Part V Machine Learning



# Machine Learning



Supervised  
learning

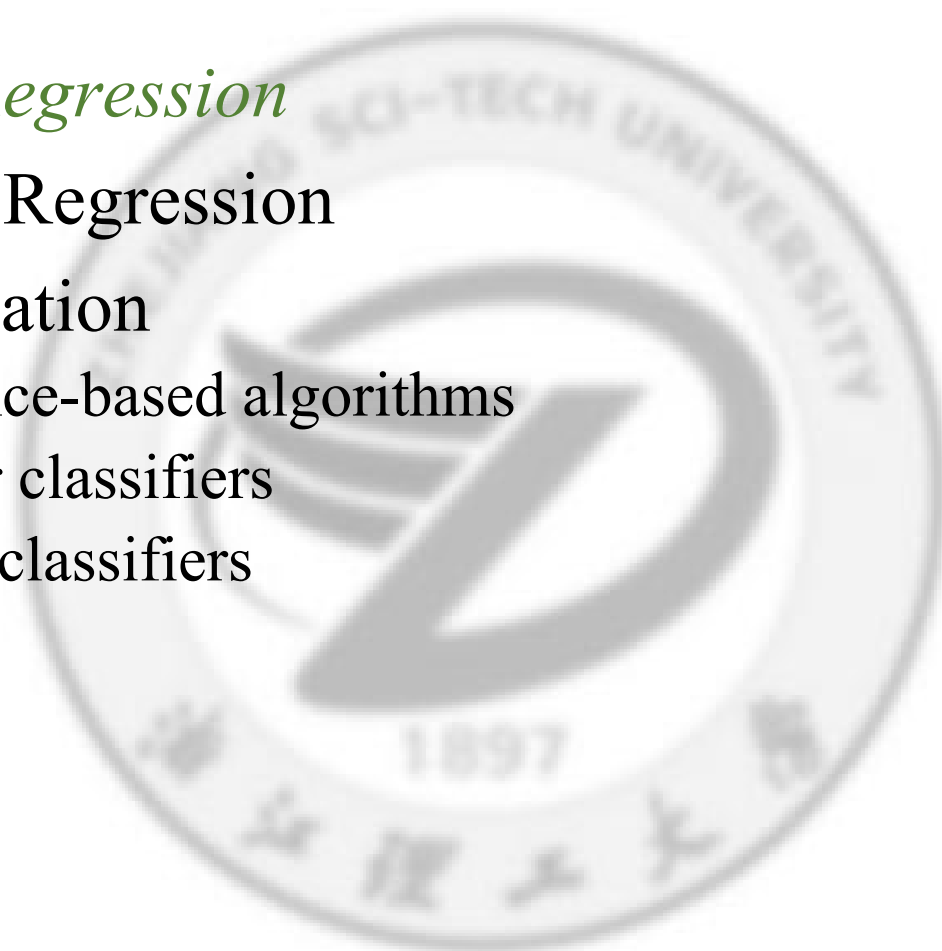
Unsupervised  
learning

Reinforcement  
learning

# Supervised learning

---

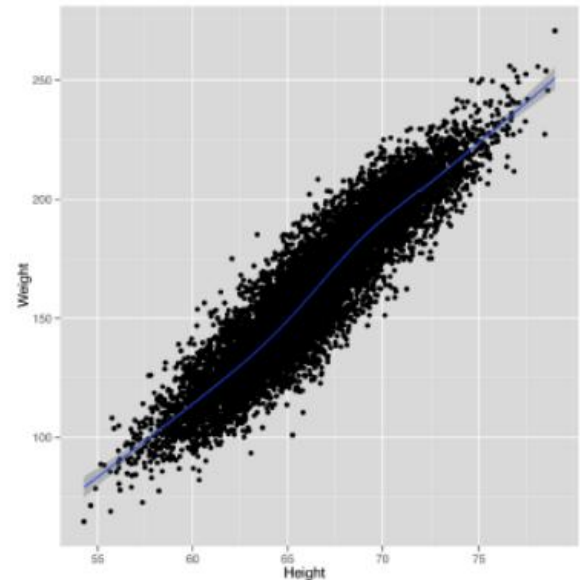
- *Linear Regression*
- Logistic Regression
- Classification
  - Distance-based algorithms
  - Linear classifiers
  - Other classifiers
- .....



# Linear Regression

## □ What is regression?

Regression is to relate **input variables** to the **output variable**, to either **predict** outputs for new inputs and/or to **interpret** the effect of the input on the output.



Height is correlated with weight.

# Linear Regression

---

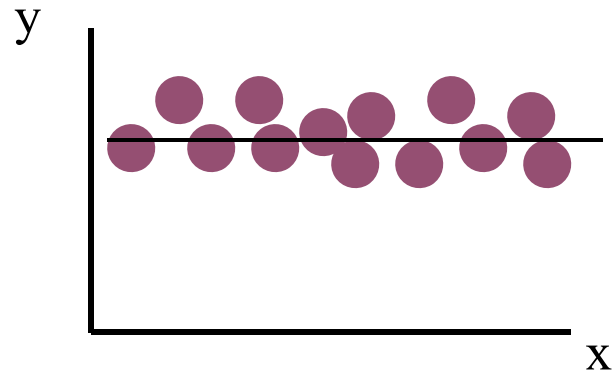
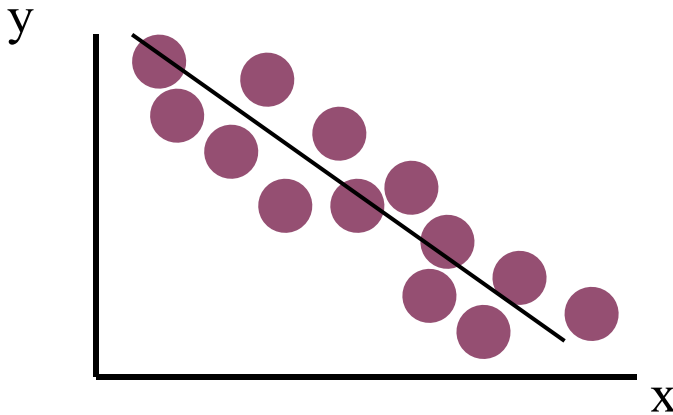
## □ Linear Regression Model

- Only **one independent variable**,  $x$
- Relationship between  $x$  and  $y$  is described by a **linear function**
- Changes in  $y$  are assumed to be related to changes in  $x$

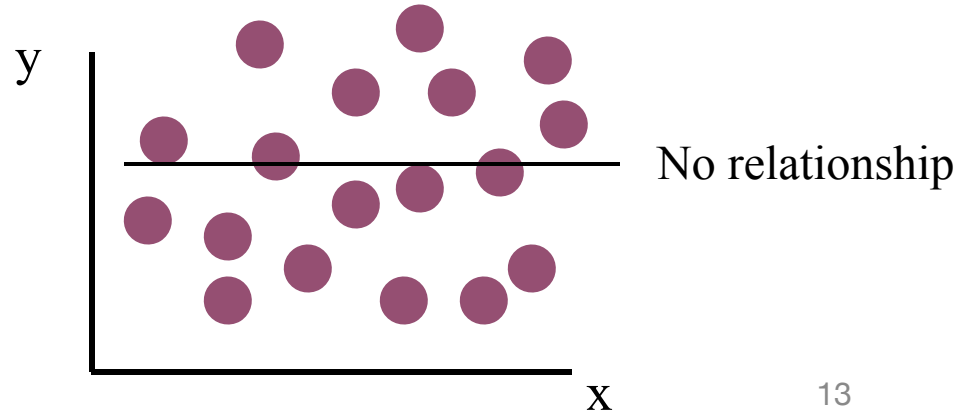
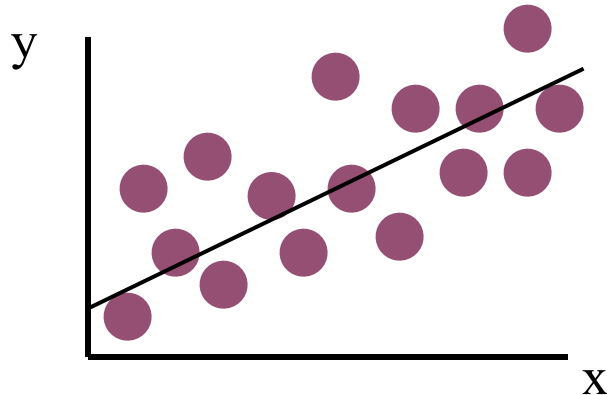
# Linear Regression

## □ Linear Regression Model

Linear relationships



Question: How to describe the linear relationships?



# Linear Regression

## □ Linear Regression Model

The diagram illustrates the Linear Regression Model equation,  $y_i = b_0 + b_1 x_i + \epsilon_i$ , with various components labeled and grouped.

Labels and arrows pointing to the equation:

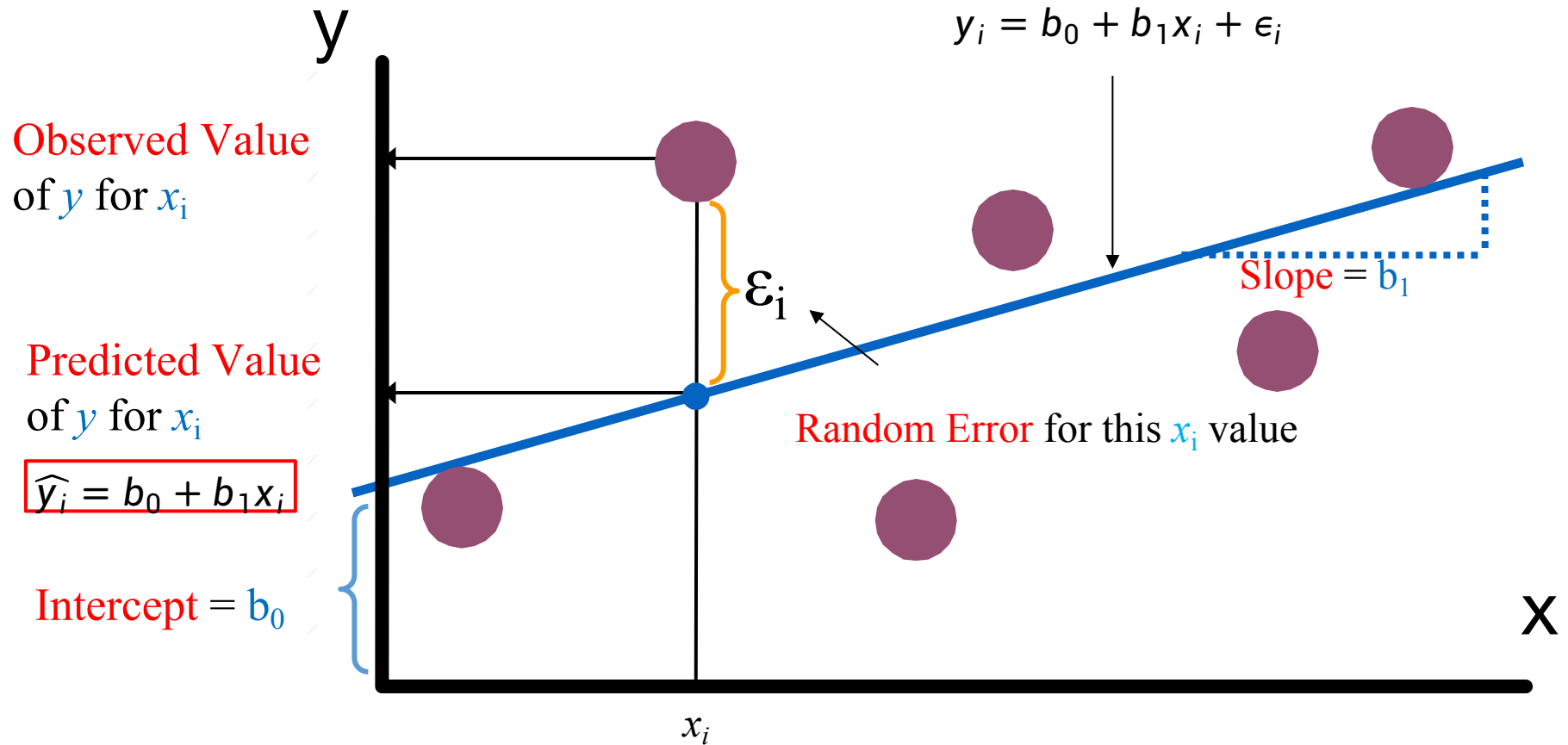
- Dependent Variable (points to  $y_i$ )
- intercept (points to  $b_0$ )
- Slope Coefficient (points to  $b_1$ )
- Independent Variable (points to  $x_i$ )
- Random Error term (points to  $\epsilon_i$ )

Groupings below the equation:

- Linear component (bracketed under  $b_0 + b_1 x_i$ )
- Random Error component (bracketed under  $\epsilon_i$ )

# Linear Regression

## □ Linear Regression Model



Question: How to obtain the best line?

# Linear Regression

## □ The Least Squares Method

$b_0$  and  $b_1$  are obtained by finding the values of that minimize the **sum** of the squared **differences** between  $y_i$  and  $\hat{y}_i$  **for all  $i$** :

$$\min \sum (y_i - \hat{y}_i)^2$$



$$\hat{y}_i = b_0 + b_1 x_i$$

$$\min \sum (y_i - (b_0 + b_1 x_i))^2 \longrightarrow \text{Objective function}$$

**Question:** How to calculate  $b_0$  and  $b_1$ ?

$$\text{derivative}[\sum (y_i - (b_0 + b_1 x_i))^2] = 0 \quad \rightarrow \quad \text{solve for } b_0, b_1$$



# Linear Regression

## □ The Least Squares Method

- Considering the objective function:

$$J = \sum (y_i - (b_0 + b_1 x_i))^2$$

- Rewrite it in matrix form as:

$$J = \|Y - \theta^T X\|_2^2$$

where  $Y = [y_1, \dots, y_n]$ ,  $X = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}$ , and  $\theta = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$

$$\frac{\partial J}{\partial \theta} = -2(Y - \theta^T X)X^T = 0$$

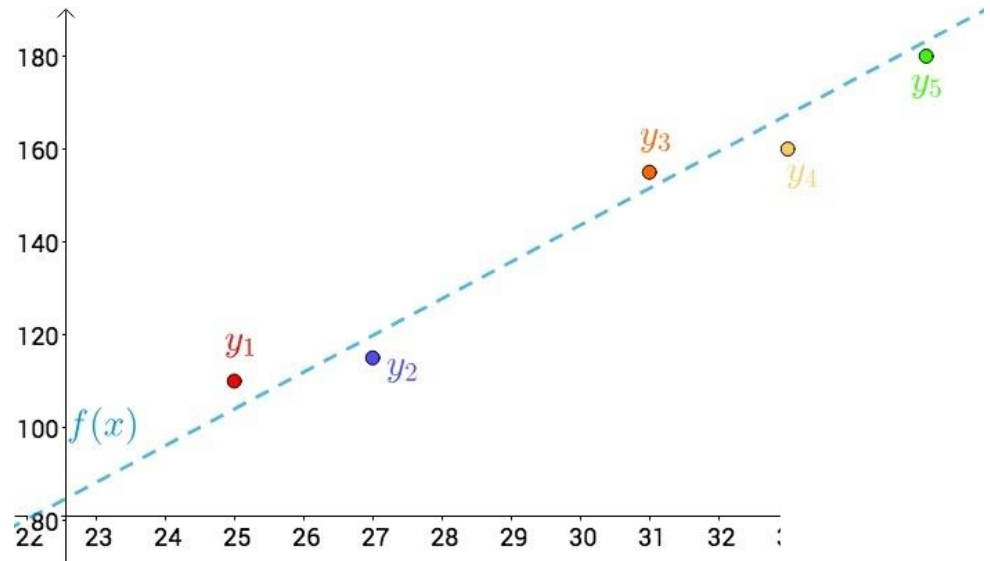
$$\theta^* = (XX^T)^{-1}XY^T$$

# Linear Regression

## □ An Example

- between **temperature** and **ice cream sales**:

Temperature	Sales
25°	110
27°	115
31°	155
33°	160
35°	180



Seems like a linear relationship

# Linear Regression

## □ An Example

- between temperature and ice cream sales:
- Set:  $y = ax + b$

Temperature	Sales
25°	110
27°	115
31°	155
33°	160
35°	180



$i$	$x$	$y$
1	25	110
2	27	115
3	31	155
4	33	160
5	35	180

# Linear Regression

## □ An Example

- between temperature and ice cream sales:
- **Set:**  $y = ax + b$
- $J = \sum (f(x_i) - y_i)^2 = \sum (ax_i + b - y_i)^2$
- $\begin{cases} \frac{\partial}{\partial a} J = 2 \sum (ax_i + b - y_i)x_i = 0 \\ \frac{\partial}{\partial b} J = 2 \sum (ax_i + b - y_i) = 0 \end{cases}$
- $\begin{cases} a \approx 7.2 \\ b \approx -73 \end{cases}$

$i$	$x$	$y$
1	25	110
2	27	115
3	31	155
4	33	160
5	35	180

# Linear Regression

## □ Another Example

- A real estate agent wishes to examine the relationship between **the selling price of a houses** and **its size** (measured in square feet)
- A random sample of 10 houses is selected
  - **Dependent variable ( $y$ ) = house price in \$1000s**
  - **Independent variable ( $x$ ) = square feet**



# Linear Regression

## □ An Example

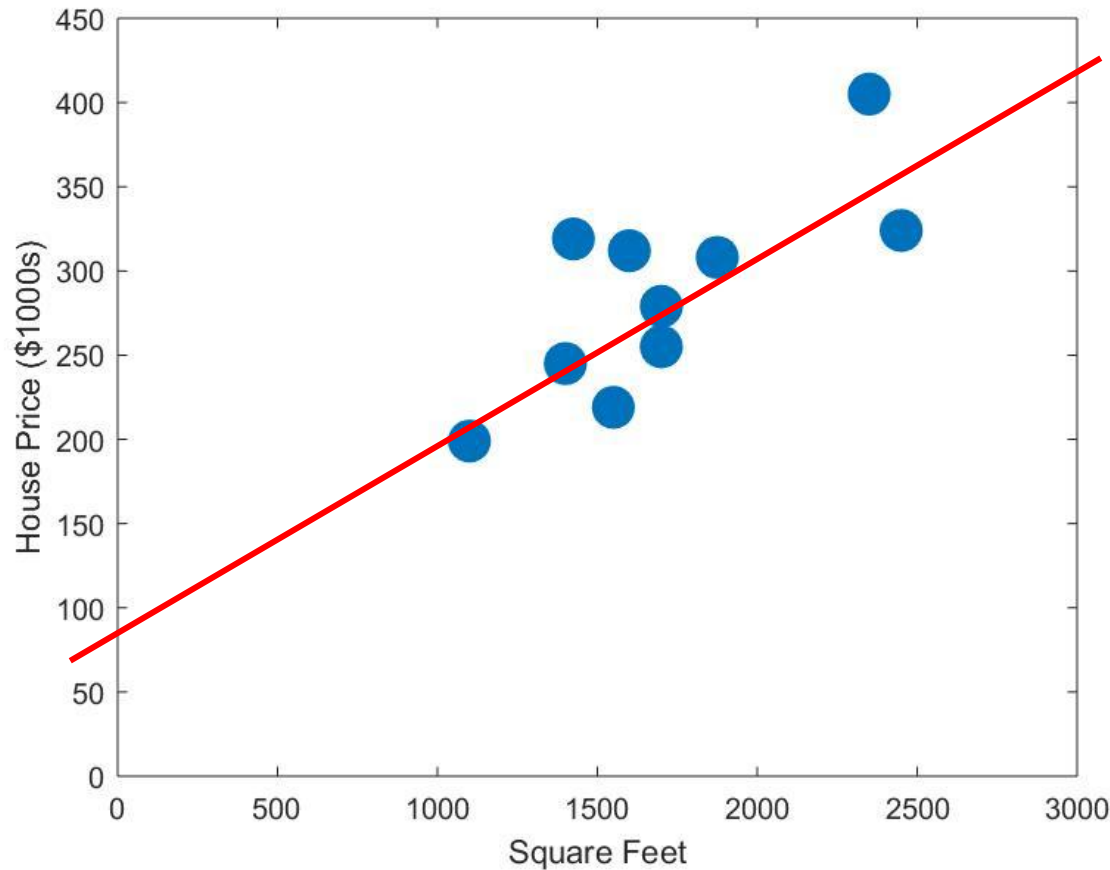
House Price ( $y$ ) in \$1000s	Square Feet ( $x$ )
245	1400
312	1600
279	1700
308	1875
199	1100
219	1550
405	2350
324	2450
319	1425
255	1700

# Linear Regression

## □ An Example

$$\theta^* = (XX^T)^{-1}XY^T$$

Scatter Plot



```
>> theta = inv(X*X')*X*Y'
```

```
theta =
```

```
98.2483
```

```
0.1098
```

```
>> [epsilon, b1, b0] = regression(X, Y)
```

```
epsilon =
```

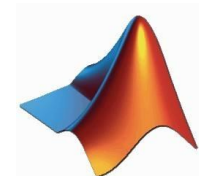
```
0.7621
```

```
b1 =
```

```
0.1098
```

```
b0 =
```

```
98.2483
```



# Linear Regression

- Conclusion: Linear Regression
- Uses least squares estimation to estimate parameters
  - Finds the line that minimizes total squared error around the line:
  - Sum of Squared Error (SSE) =  $\sum (y_i - (b_0 + b_1 x))^2$
  - Minimize the squared error function:  
derivative  $[\sum (y_i - (b_0 + b_1 x))^2] = 0 \rightarrow$  solve for  $b_0, b_1$



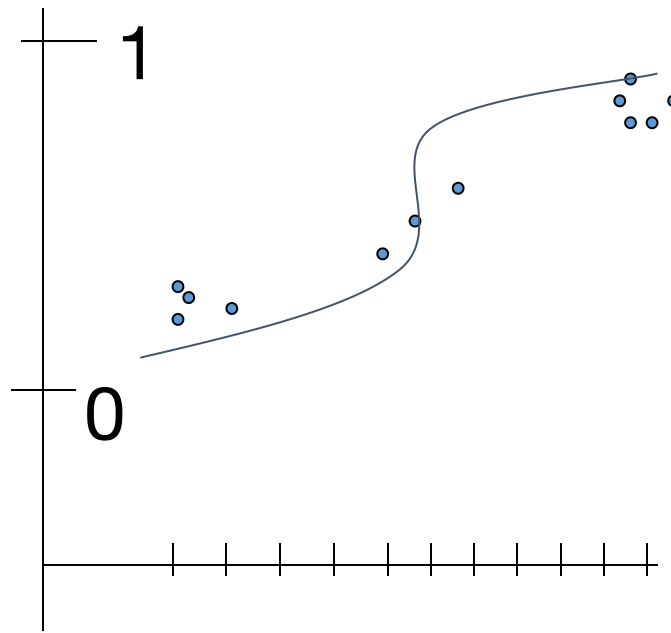
# Linear Regression

## □ Thinking...

Could model **probability** of lung cancer...

$$P \leftarrow b_0 + b_1 x_i$$

The probability of lung cancer (p)



*But why might this not be best modeled as linear?*

Smoking  
(cigarettes/day)

# Logistic Regression

---

## □ Logistic Regression Model

- In medical research, it is often necessary to analyze which **factors** are related to the outcome of a certain outcome.
- How do we find out which factors have a **significant impact** on the outcome?
- Logistic regression analysis can solve these problems better.

# Logistic Regression

---

- Linear regression is written as:

$$y = b_0 + b_1X \quad -\infty \leq y \leq +\infty$$

- If we define y as disease or normal, it can not be modeled by the above equation.
- How about apply the probability to represent it?

$$p \leftarrow b_0, b_1, X$$

# Logistic Regression

## □ Logistic Regression Model

Think about the probability...

probability of disease :  $p$   $0 \leq p \leq 1$

probability of no-disease :  $1-p$   $0 \leq p \leq 1$

odds:  $\frac{p}{1-p}$   $0 \leq \frac{p}{1-p} < +\infty$

$\ln\left(\frac{p}{1-p}\right)$   $-\infty < \ln\left(\frac{p}{1-p}\right) < +\infty$

# Logistic Regression

## □ Logistic Regression Model

Define logistic model as

$$\ln \frac{p}{1-p} = b_0 + b_1 X$$

We obtained that,

$$p = \frac{1}{1 + e^{-(b_0 + b_1 X)}}$$

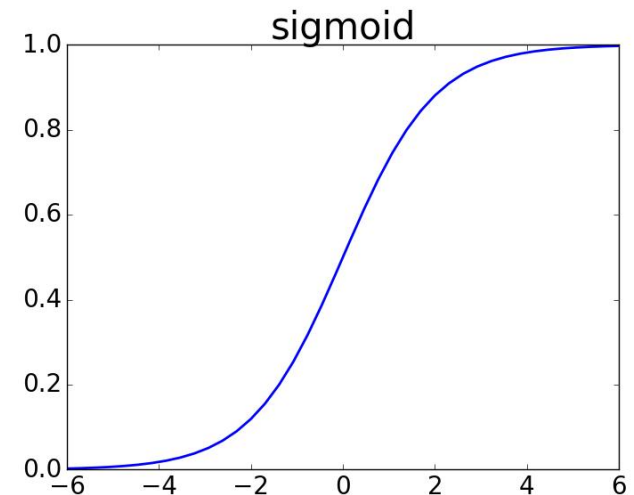
$$h_{\theta}(X) = \frac{1}{1 + e^{-\theta^T X}}$$

Therefore,

$$P(\text{class} = 1|x; \theta) = h_{\theta}(X)$$

$$P(\text{class} = 0|x; \theta) = 1 - h_{\theta}(X)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$



The output of sigmoid function could be used to indicate the probability.

# Logistic Regression

## □ Logistic Regression Model

$$P(\text{class} = 1|x; \theta) = h_{\theta}(X)$$

$$P(\text{class} = 0|x; \theta) = 1 - h_{\theta}(X)$$



$$P(\text{class} = y|x; \theta) = h_{\theta}(X)^y (1 - h_{\theta}(X))^{1-y}$$

Considering all the given data (training set):

$$X = [x_1, \dots, x_n], \quad Y = [y_1, \dots, y_n],$$

$$L(\theta) = \prod_i^n h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

$$\text{The cost function : } J = -\frac{1}{n} \log (L(\theta))$$

# Logistic Regression

## □ Conclusion

### ■ Logistic regression

- Uses sigmoid and log function and to estimate the parameters
- According to the **Maximum Likelihood Estimate**, construct the loss function:

$$J = -\frac{1}{m} \log (L(\theta))$$

where,

$$L(\theta) = \prod_i^n h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

- Minimize the cost:

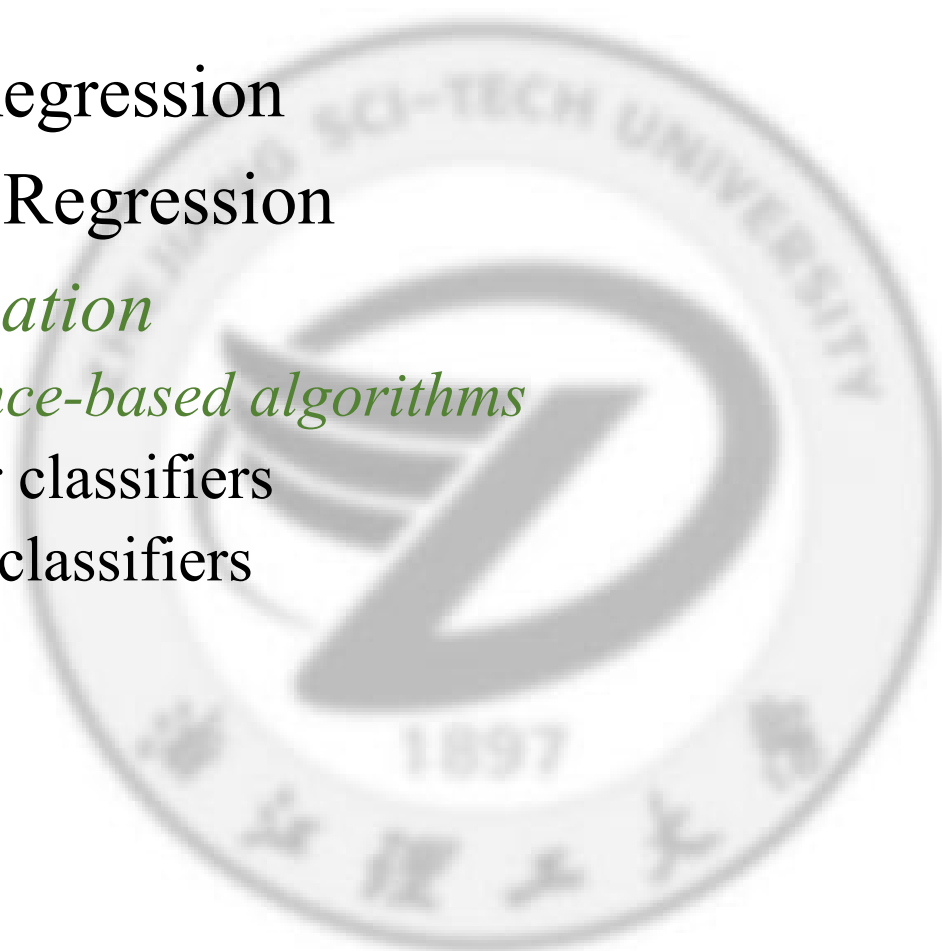
$$\frac{\partial J}{\partial \theta} = 0 \quad \rightarrow \quad \text{solve for } \theta \quad \text{HOW?}$$

Try to solve it by yourself.

# Supervised learning

---

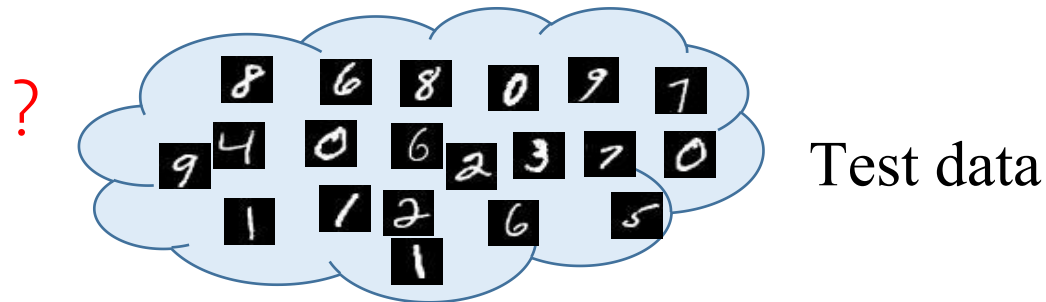
- Linear Regression
- Logistic Regression
- *Classification*
  - *Distance-based algorithms*
  - Linear classifiers
  - Other classifiers
- .....





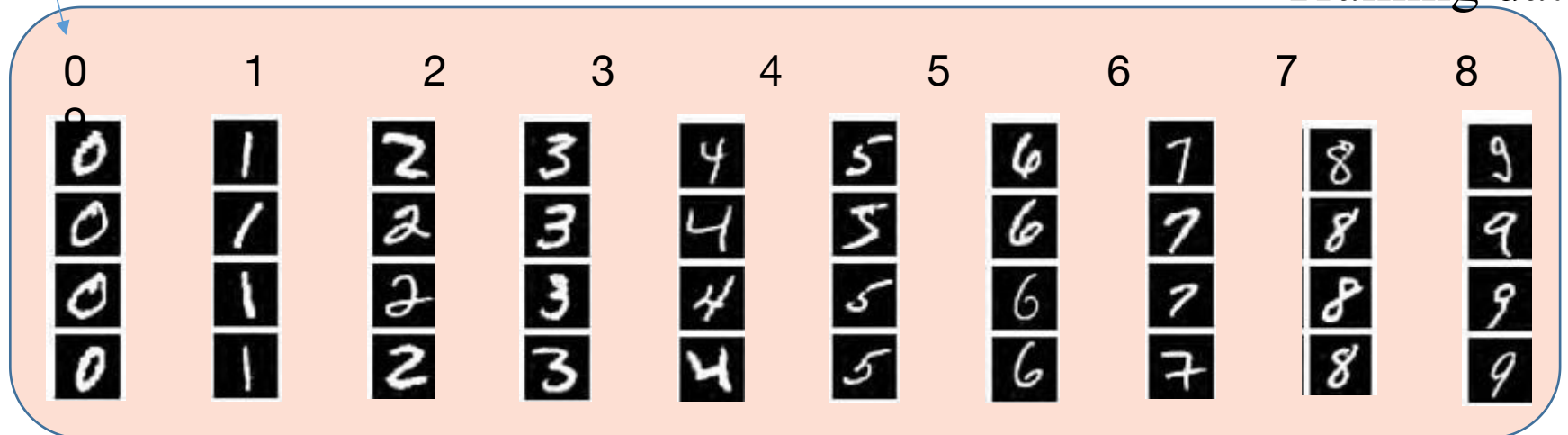
# Classification

Multi-class classification assigns test samples to a certain class.

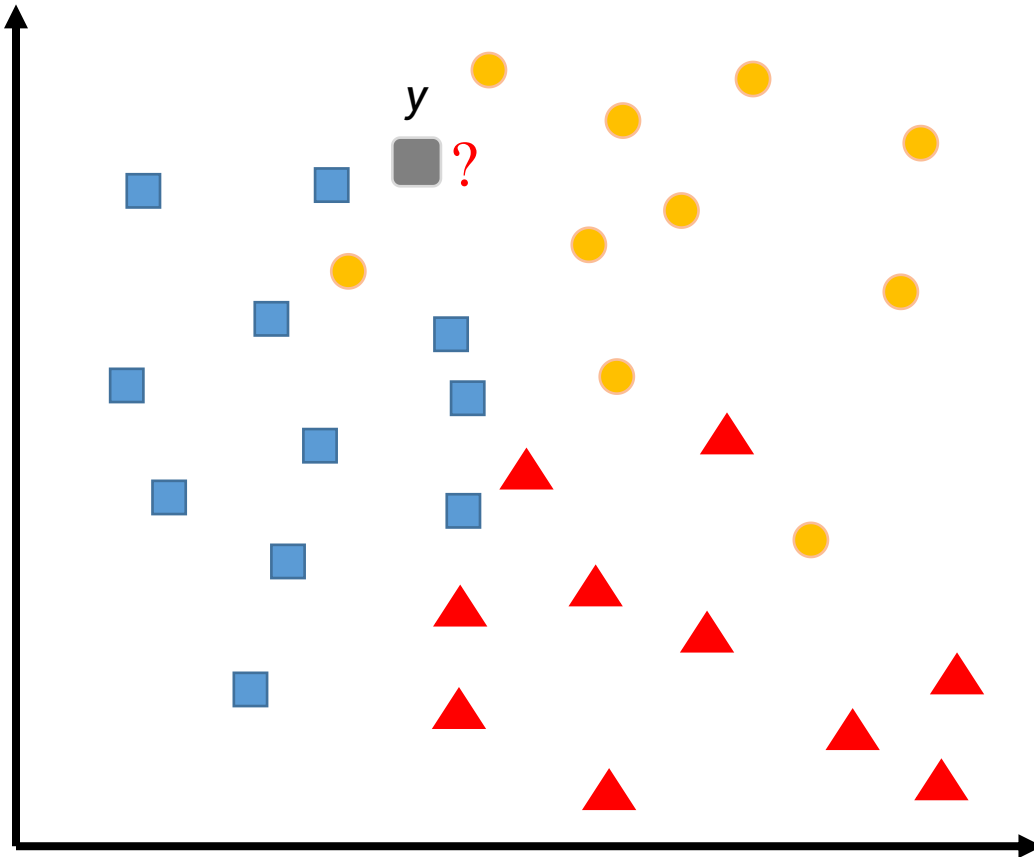


Labels

Training data



# Classification



Training data:

$$X = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$$

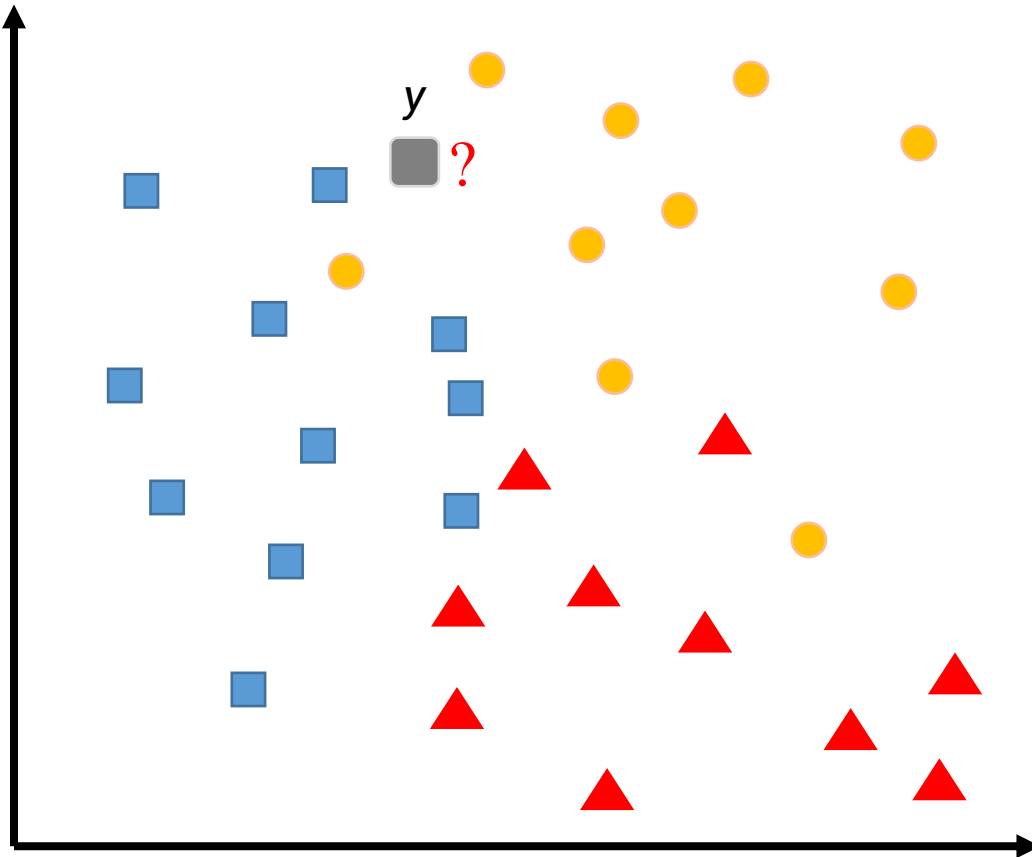
and  
training labels:

$$L = \{l^{(1)}, l^{(2)}, \dots, l^{(N)}\}$$

$N$ : the number of training data

# Classification

## □ Nearest neighbor



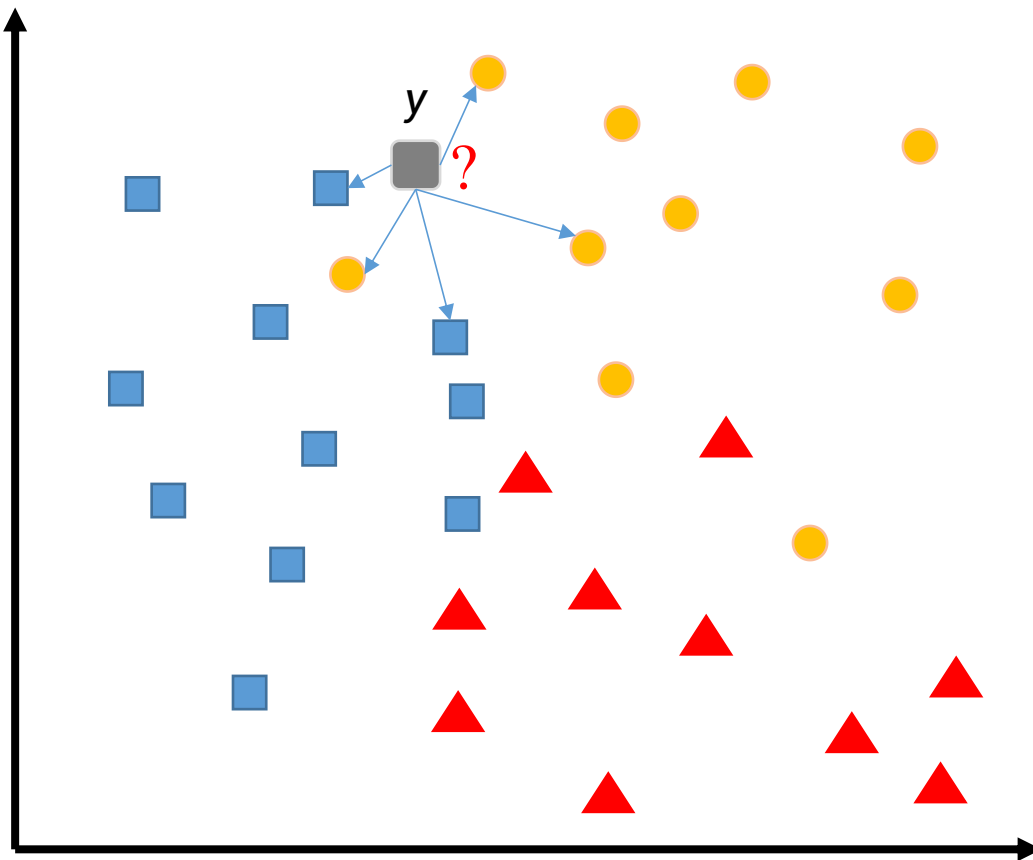
How to decide which is the nearest one?

*The distance  $d(\mathbf{x}, \mathbf{y})$  between two points  $\mathbf{x} \in R^n$  and  $\mathbf{y} \in R^n$  can for example be measured by the Euclidean distance.*

$$d(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sqrt{\sum_{i=1}^n (x_i^{(1)} - x_i^{(2)})^2}$$

# Classification

## □ Nearest neighbor



How to decide which is the nearest

$$d^j(x^{(y)}, y) = \sqrt{\sum_{i=1}^n (x_i^{(j)} - y)^2}$$

Calculate all the distances from the training data to the test data  $y$ , and we obtain:

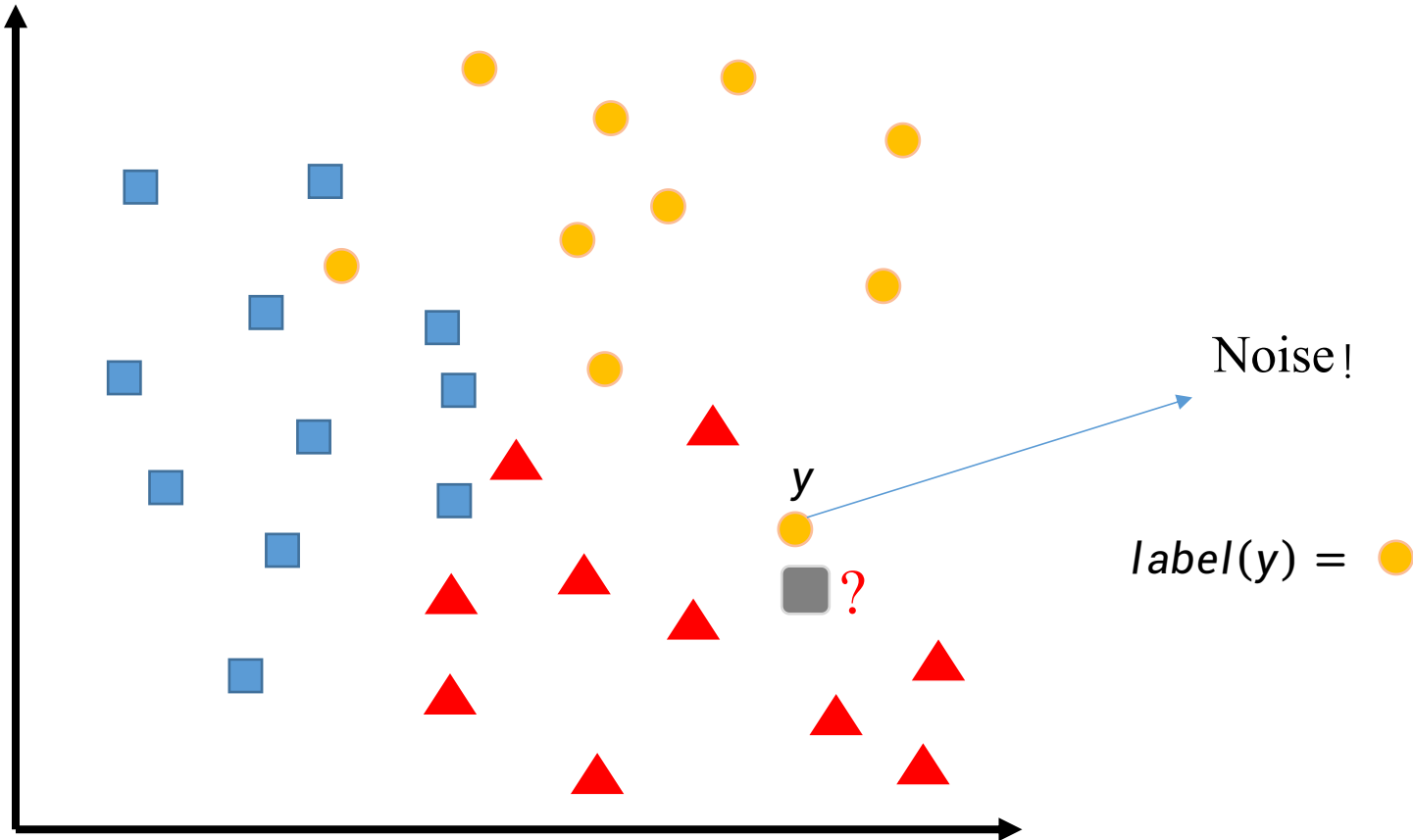
$$D = [d^{(1)}, d^{(2)}, \dots, d^{(N)}]$$

$$s = \operatorname{argmin}_i d^{(i)}$$

$$\operatorname{label}(y) = \operatorname{label}(x^{(s)}) = \text{blue square}$$

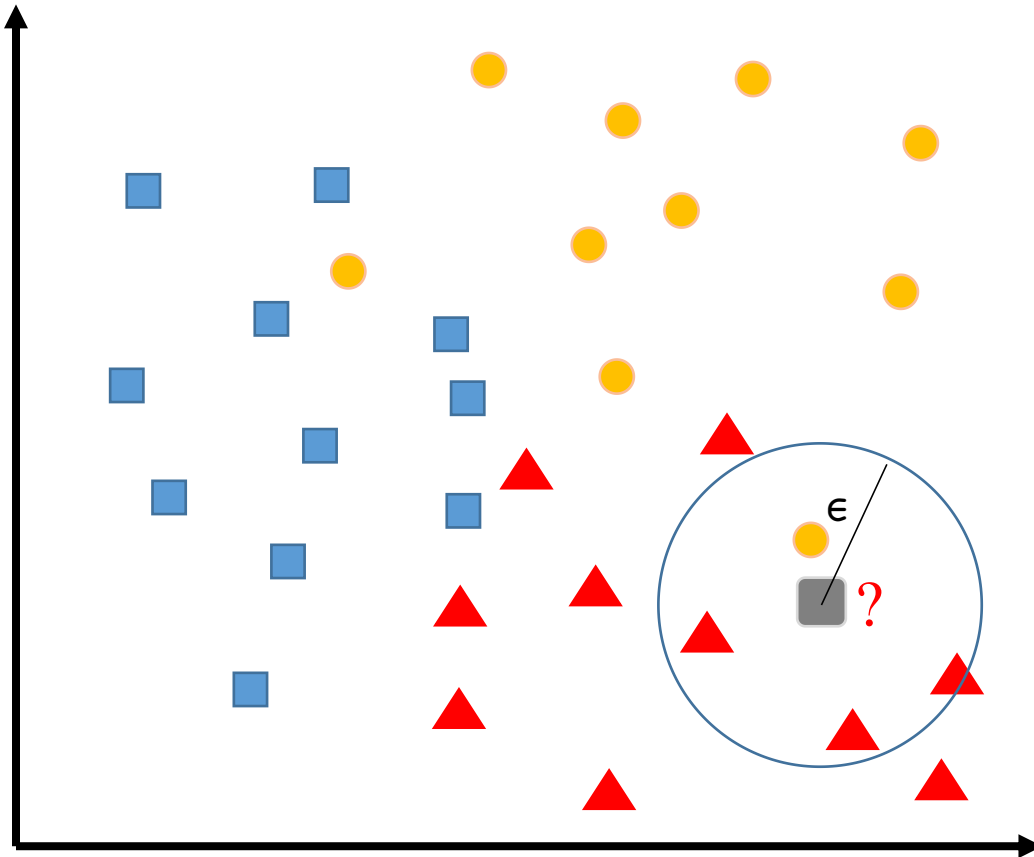
# Classification

## □ Nearest neighbor



# Classification

## □ $\epsilon$ -ball Nearest neighbor



Select a value  $\epsilon$ , then draw a ball in  $\mathbb{R}^n$  with  $y$  as the center and  $\epsilon$  as the radius.

The label of  $y$  is decided by majority labels of points in this ball.

In this ball:

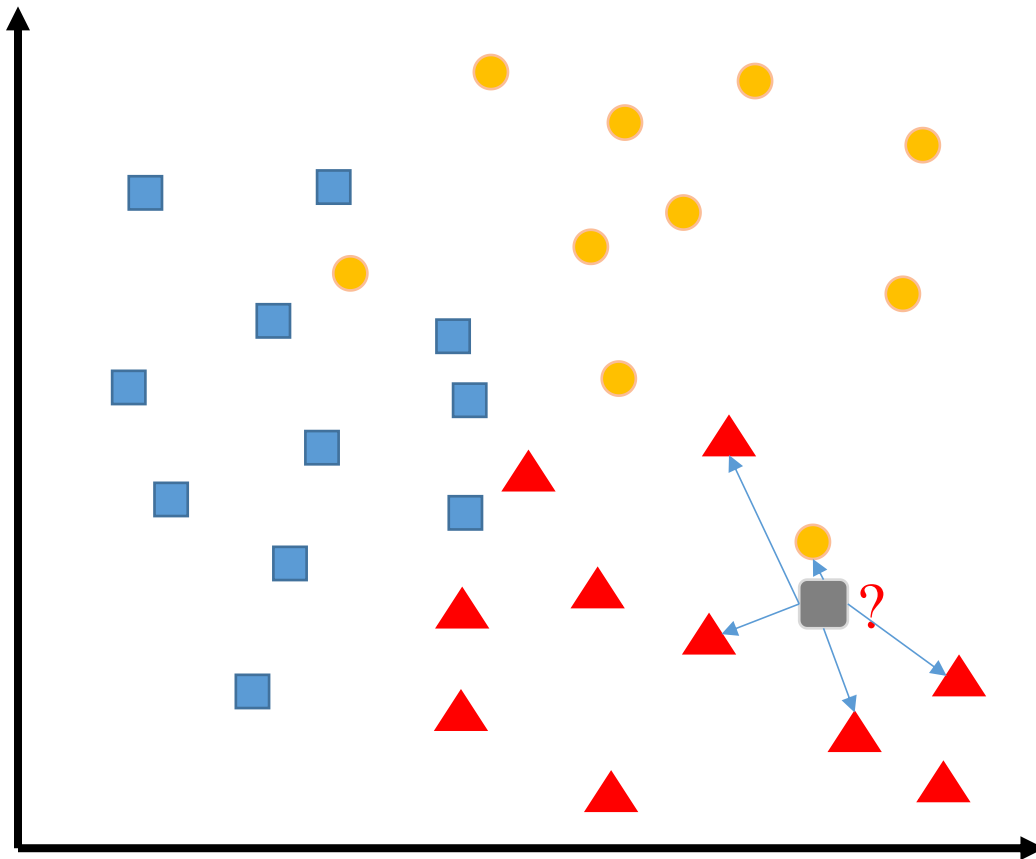
▲ : 3

● : 1

■ belongs to ▲

# Classification

## □ K Nearest neighbor



Select a value  $k$ , then find  $y$ 's  $k$  nearest neighbor.

The label of  $y$  is decided by majority labels of  $y$ 's  $k$  neighbors.

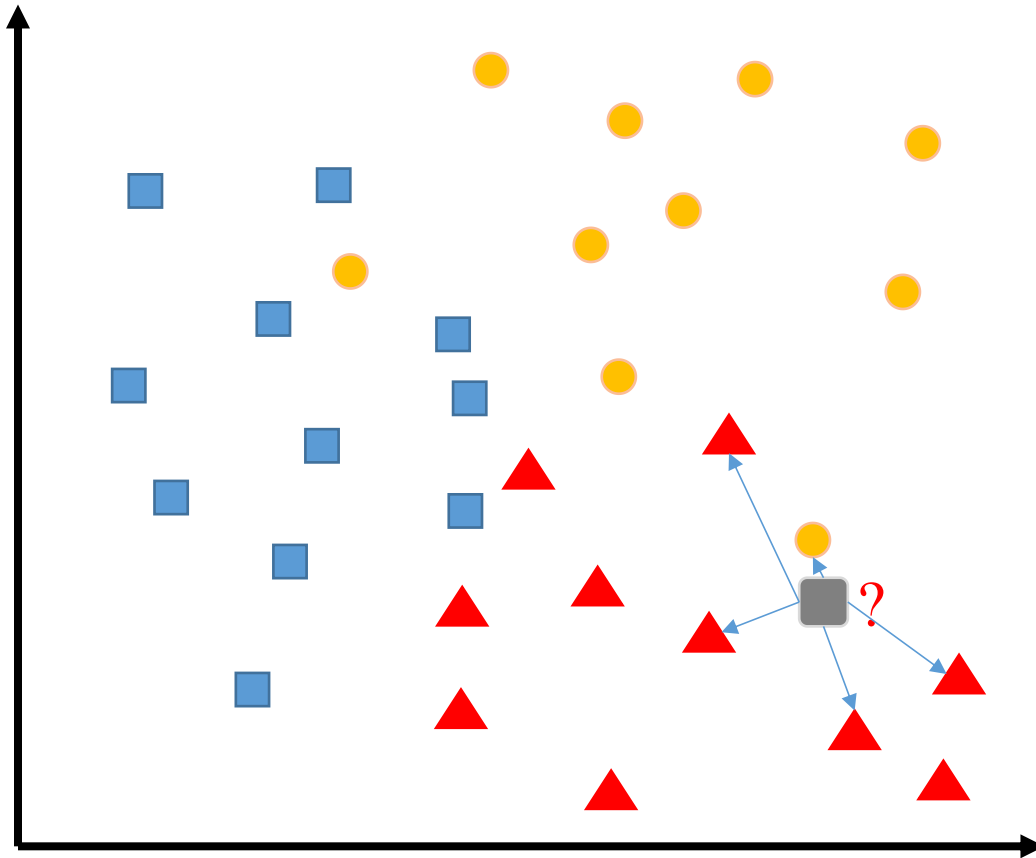
Let  $k$  be 5,

▲ : 5      ● : 1

■ belongs to ▲

# Classification

## □ K Nearest neighbor



Question:

How to decide  $k$ ?

Which algorithm achieve better performance?

▲ : 5

● : 1

■ belongs to ▲



# Classification

## □ Distance Metrics

- Euclidean distance

- $d_e(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

- Sum of squared distance

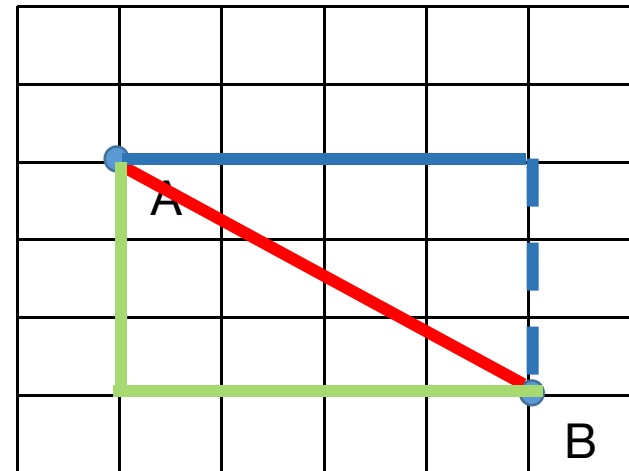
- $d_q(x, y) = \sum_{i=1}^n (x_i - y_i)^2$

- Manhattan distance

- $d_m(x, y) = \sum_{i=1}^n |x_i - y_i|$

- Chebyshev distance

- $d_c(x, y) = \max_{i=1, \dots, n} |x_i - y_i|$



# Classification



## □ Nearest neighbor classifier

Problem:





















- Need to determine value of parameter  $K$
- Distance based learning is not clear which **type of distance** to use and which attribute to use to produce the best results.
- Computation cost is quite high because we need to compute distance of each query instance to all training samples.

# Classification

## □ Example

- Each image is represented by a vector of dimension 784.

The matrix indicates the pairwise distances.

										
	0	2.8735	2.1766	2.6559	2.2201	2.2500	2.0893	2.4795	2.8443	2.1202
	2.8735	0	2.5055	2.8681	2.9475	2.6062	2.8493	2.8330	2.9434	3.1619
	2.1766	2.5055	0	2.9024	2.3556	0.7858	2.3561	2.2060	2.5274	2.4331
	2.6559	2.8681	2.9024	0	2.7428	2.9531	3.0539	2.8362	2.8488	2.6425
	2.2201	2.9475	2.3556	2.7428	0	2.5284	2.1733	2.4262	2.3432	2.5895
	2.2500	2.6062	0.7858	2.9531	2.5284	0	2.4679	2.2906	2.5549	2.3900
	2.0893	2.8493	2.3561	3.0539	2.1733	2.4679	0	2.5580	2.7456	2.3759
	2.4795	2.8330	2.2060	2.8362	2.4262	2.2906	2.5580	0	2.8885	2.5823
	2.8443	2.9434	2.5274	2.8488	2.3432	2.5549	2.7456	2.8885	0	2.9773
	2.1202	3.1619	2.4331	2.6425	2.5895	2.3900	2.3759	2.5823	2.9773	0

The distance between the data is inconsistent with similarity of the content of the image .

# Supervised learning

---

- Linear Regression
- Logistic Regression
- Classification
  - Distance-based algorithms
  - *Linear classifiers*
  - Other classifiers
- .....



# Classification

## □ Linearly separable

Apple = [diameter, color, shape, spots, place of production]



$$A_1 = \begin{bmatrix} 7.8 \\ 0.2 \end{bmatrix}$$



$$A_2 = \begin{bmatrix} 7.4 \\ 0.2 \end{bmatrix}$$



$$A_3 = \begin{bmatrix} 7.1 \\ 0.1 \end{bmatrix}$$



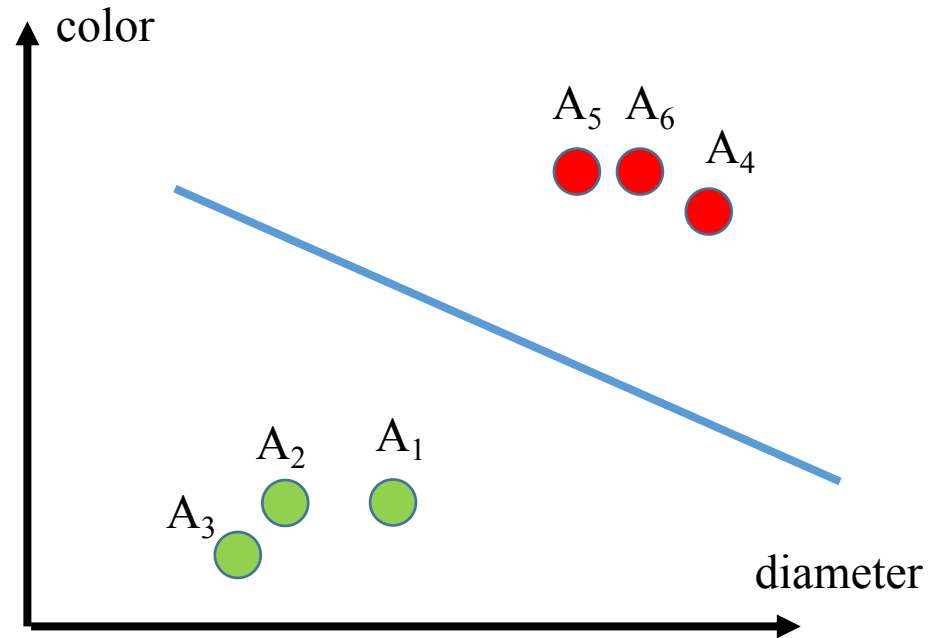
$$A_4 = \begin{bmatrix} 8.5 \\ 0.7 \end{bmatrix}$$



$$A_5 = \begin{bmatrix} 8.1 \\ 0.8 \end{bmatrix}$$



$$A_6 = \begin{bmatrix} 8.3 \\ 0.8 \end{bmatrix}$$



These training data are *linearly separable*

# Classification

## □ Linearly separable

Apple = [diameter, color, shape, spots, place of production]



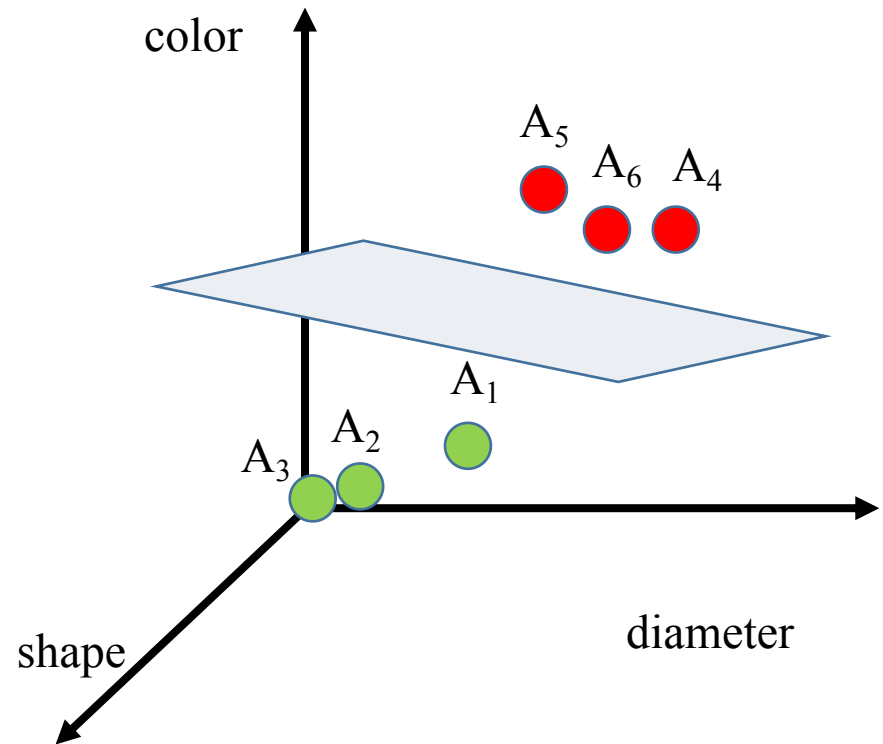
$$A_1 = \begin{bmatrix} 7.8 \\ 0.2 \\ 0.6 \end{bmatrix}$$



$$A_2 = \begin{bmatrix} 7.4 \\ 0.2 \\ 0.7 \end{bmatrix}$$



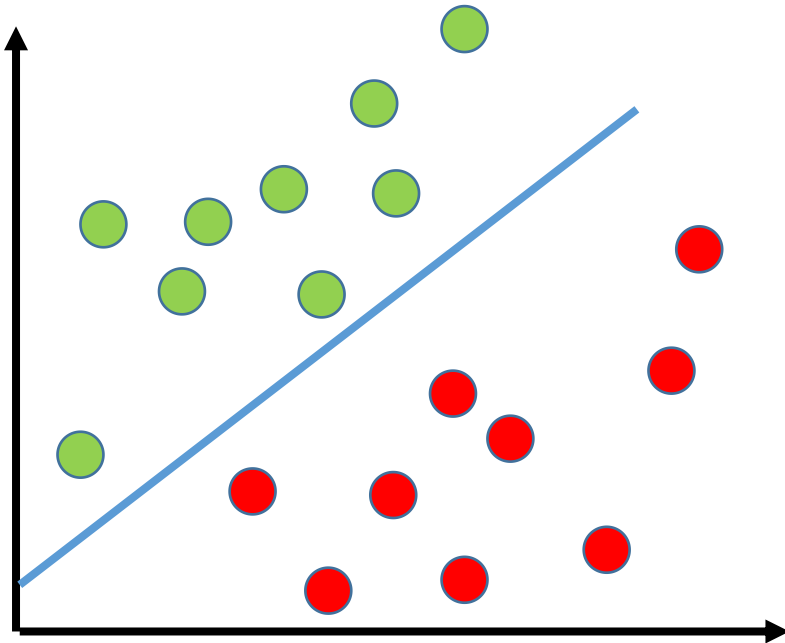
$$A_3 = \begin{bmatrix} 7.1 \\ 0.1 \\ 0.6 \end{bmatrix}$$



In  $n$  dimensions a hyperplane is needed for the separation.

# Classification

## □ Linearly separable

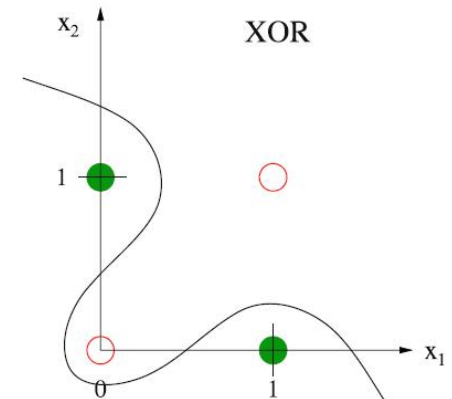
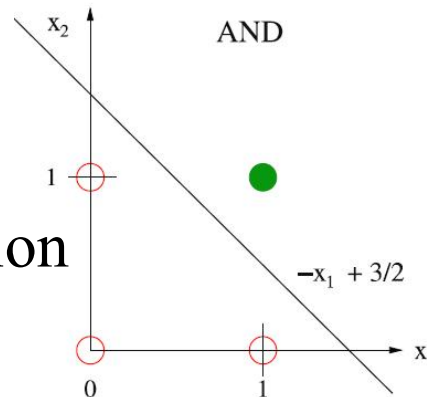


A linearly separable two dimensional data set. The equation for the dividing straight line is

$$w_1x_1 + w_2x_2 = 1$$

Every  $(n - 1)$ -dimensional hyperplane in  $R^n$  can be described by an equation

$$\sum_{i=1}^n w_i x_i + b = 0$$



The boolean function AND is linearly separable, but XOR is not (●  $\hat{=}$  true, ○  $\hat{=}$  false)

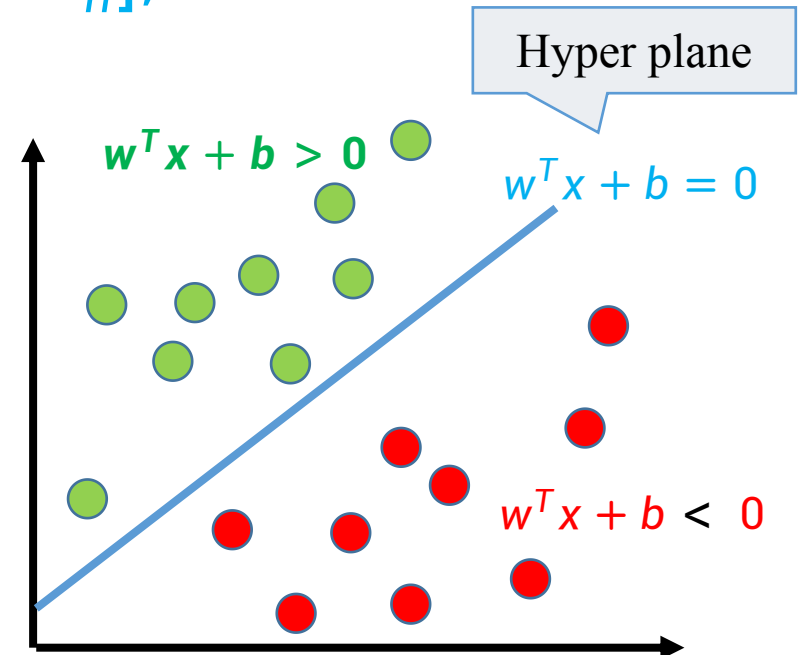
# Classification

## □ Linearly separable

- **Definition** Two sets  $M_1 \subset \mathbb{R}^n$  and  $M_2 \subset \mathbb{R}^n$  are called *linearly separable*.
- if real vector  $\mathbf{w}=[w_1, w_2, \dots, w_n], b$  exist with

$$\begin{aligned} \mathbf{w}^T \mathbf{x} + b &> 0 \text{ for all } \mathbf{x} \in M_1 \\ \text{and} \\ \mathbf{w}^T \mathbf{x} + b &\leq 0 \text{ for all } \mathbf{x} \in M_2 \end{aligned}$$

$$\text{classify}(\mathbf{x}; \mathbf{w}, b) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$





# Classification

## □ Linearly separable

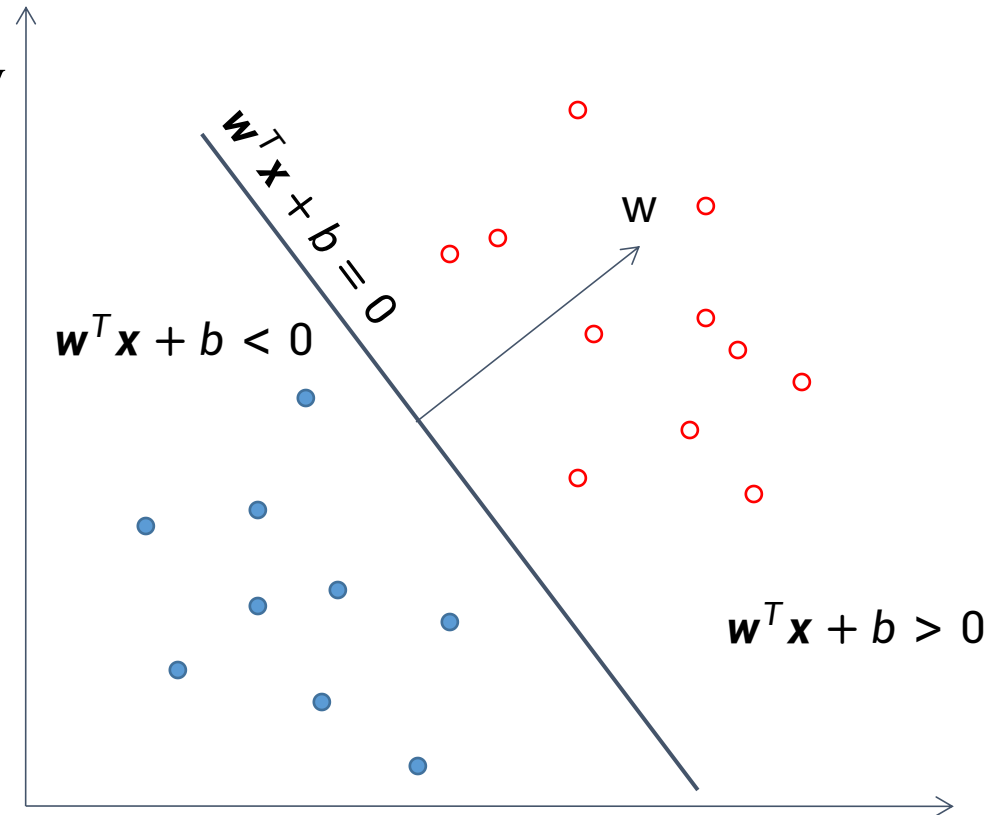
- Given a training set which is linearly separable:

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\},$$

$$\mathbf{x}_i \in X = \mathbb{R}^d,$$

$$y_i \in Y = \{-1, +1\}, \quad i = 1, 2, \dots, m$$

- Goal: solve a separating hyperplane  
 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$

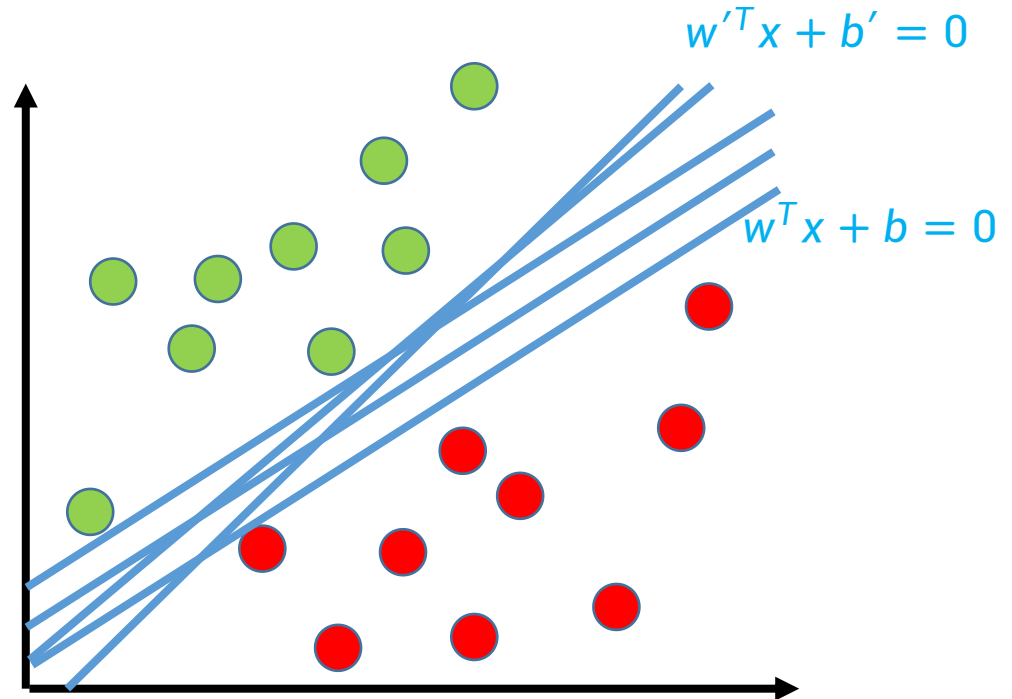


# Classification

## □ Linearly separable

Any of these would be fine.

But which is best?

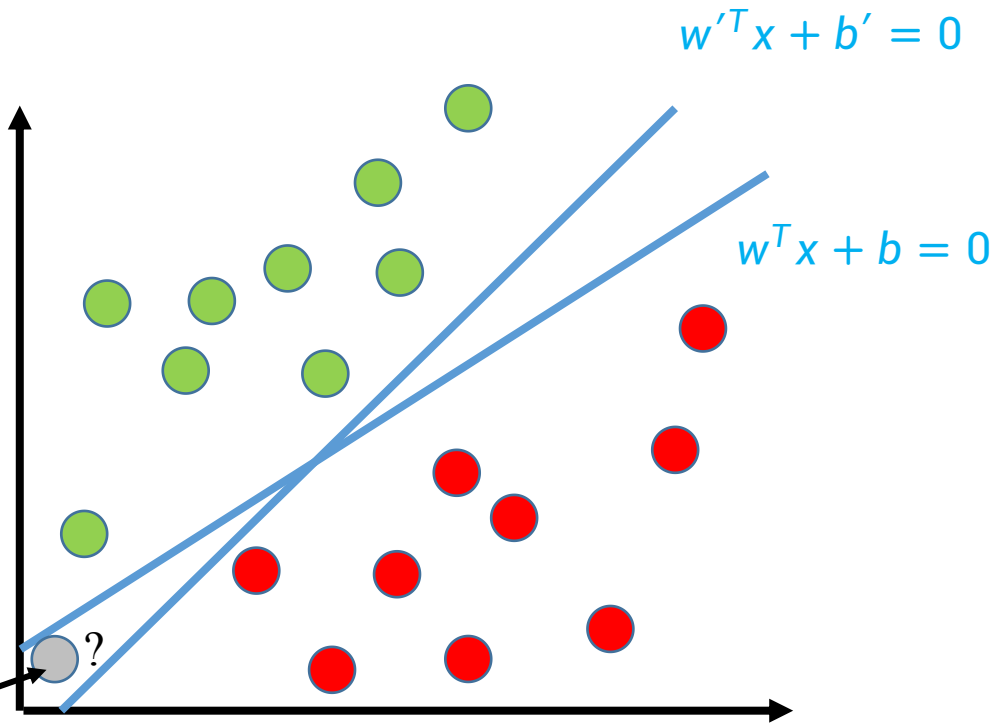


# Classification

## □ Linearly separable

These two lines can separate the two classes of points.

How would you classify this point?

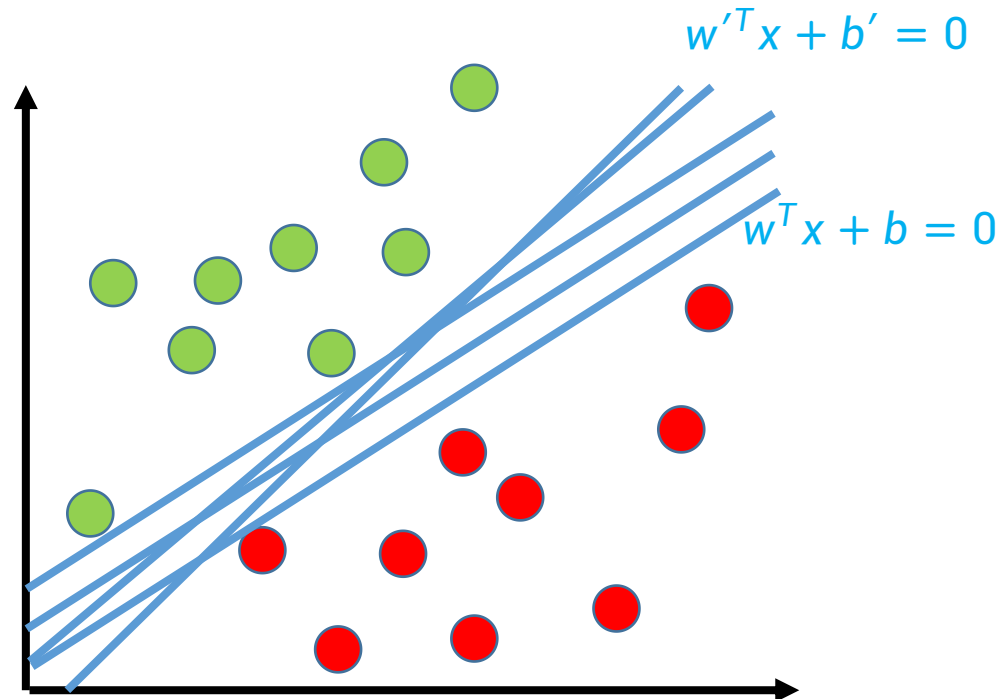


# Classification

## □ Linearly separable

The distance between a sample and hyperplane indicates the **classification confidence**:

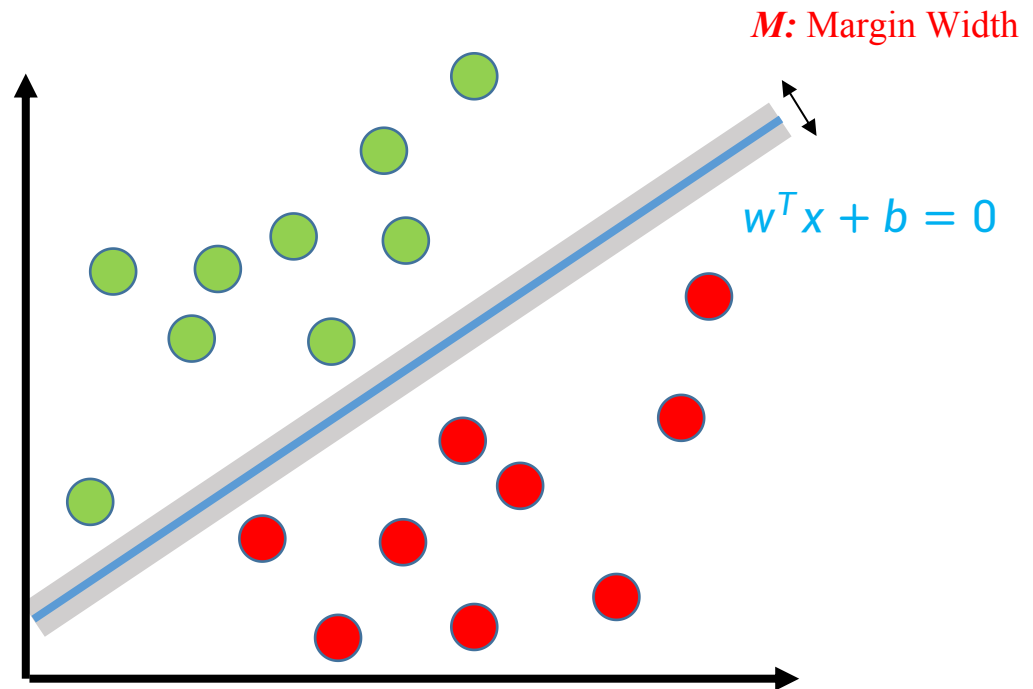
- The farther the sample is from the hyperplane, the higher the confidence that it will be correctly classified.
- The closer the sample is to the hyperplane, the lower the confidence that it will be correctly classified.



# Classification

## □ Linearly separable

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

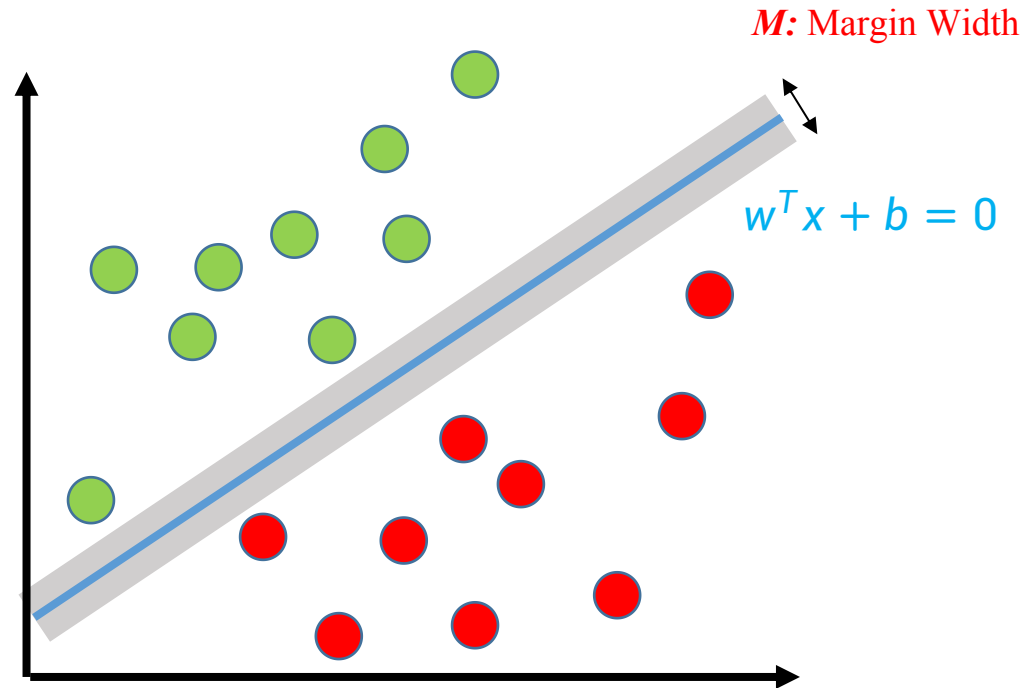


# Classification

## □ Linearly separable

**Hyperplane with maximum margin:** the hyperplane separating samples with maximum margin, which is more robust for classification.

- Two-class samples are separating on corresponding side of the hyperplane;
- The distance from the closest sample point to the hyperplane on both sides to the hyperplane is maximized.

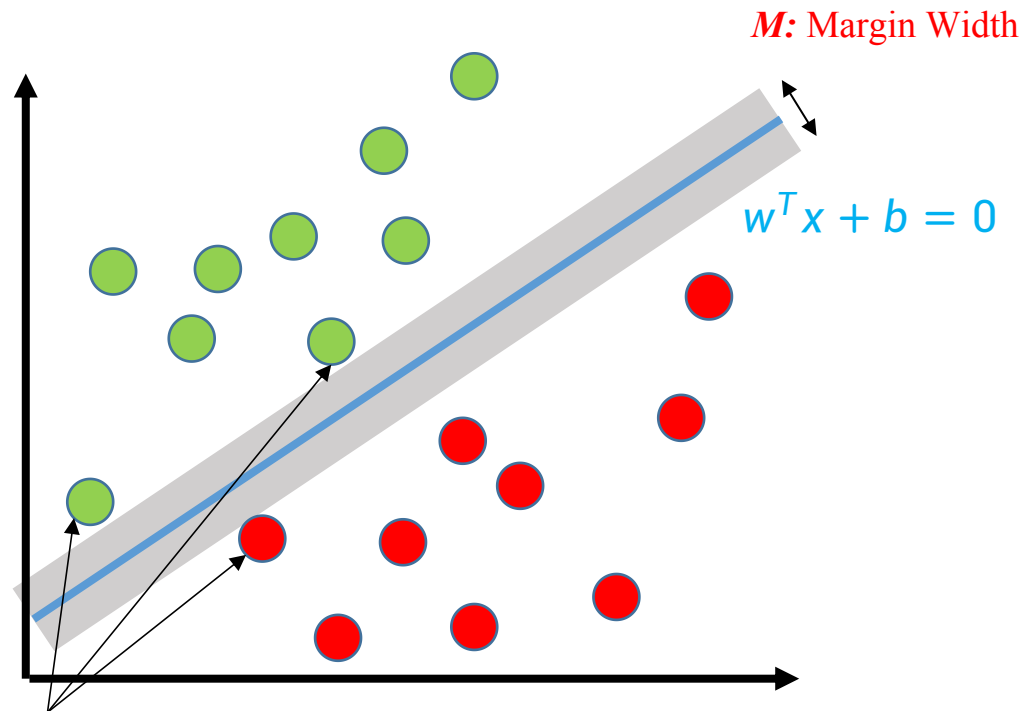


**Support Vector Machine (SVM):** the optimal separating hyperplane when samples are linearly separable.

# Classification

## □ Support vector machine (SVM)

1. Maximizing the margin is good.
2. Implies that only support vectors are important; other training examples are ignorable.
3. Empirically it works very well.



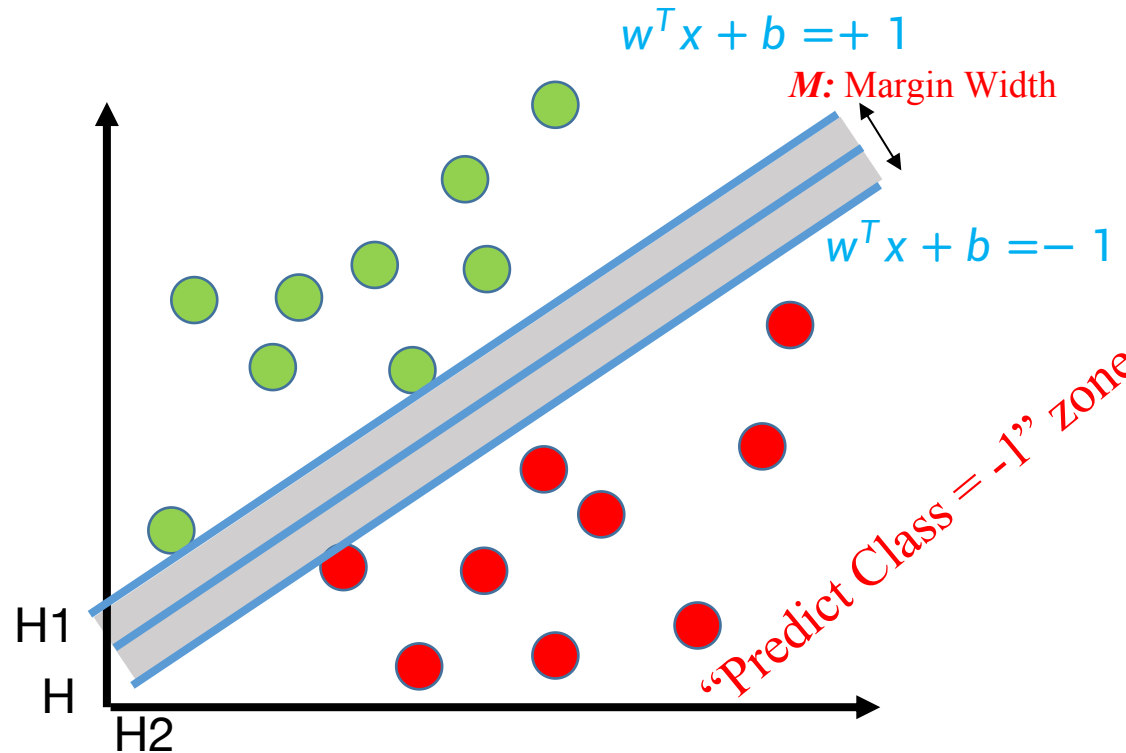
**Support Vectors:** are those data points that the margin pushes up against.

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

# Classification

## □ Linearly separable

- **Margin:**  $H_1$  and  $H_2$  are boundary hyperplanes that pass through the samples closest to the  $H$  and parallel to  $H$ . The distance between  $H_1$  and  $H_2$  is called separating margin.
- **Optimal separating hyperplane:** A hyperplane separating samples correctly, and samples (two-class) closest to the hyperplane also has the maximum distance from the hyperplane.





# Classification

## □ Linearly separable

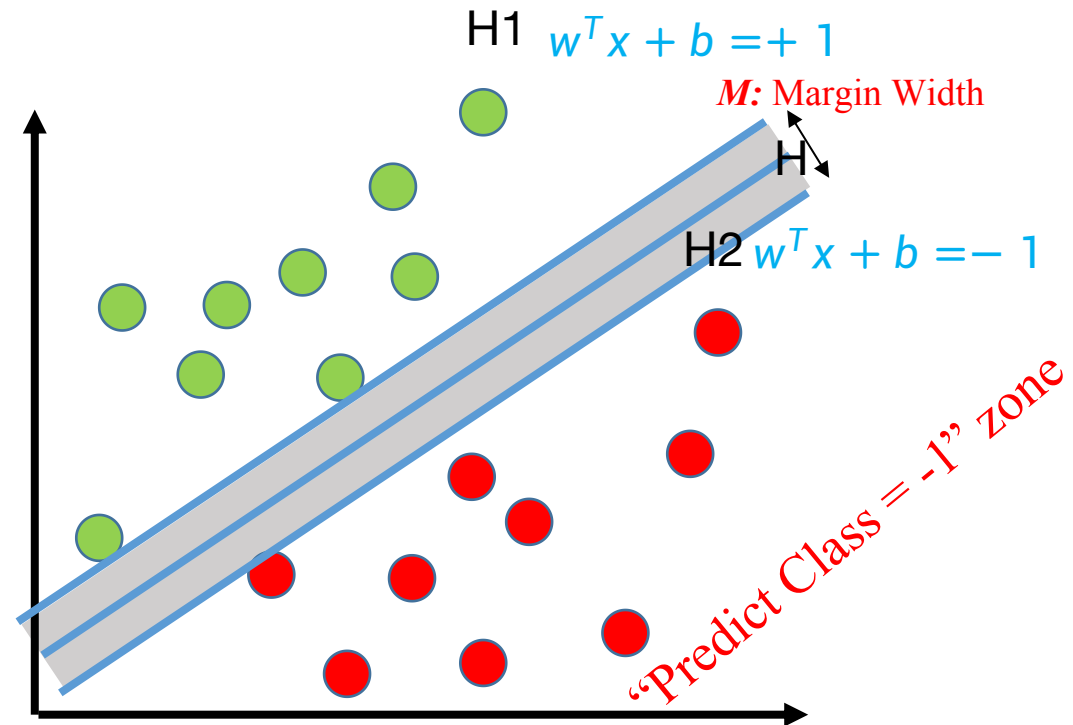
For any line for classification, we can calculate the margin.

Then, which line is the best?

Goal:

- 1) Correctly classify all training data
- 2) Maximize the Margin

How to achieve this goal?



# Classification

## □ SVM

- Given a training set which is linearly separable:  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ,  $\mathbf{x}_i \in X = R^d$ ,  $y_i \in Y = \{-1, +1\}$
- Hyperplane H:  $\mathbf{w}^T \mathbf{x} + b = 0$
- The distance between any sample  $\mathbf{x}$  in feature space to H:  
$$r = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$$

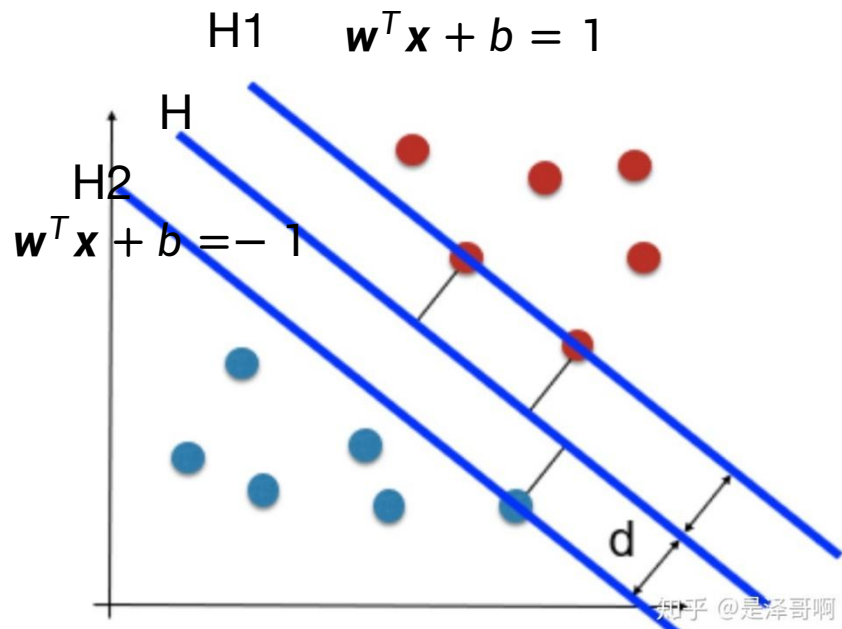
# Classification

## □ SVM-- Goal 1

- Linearly separable
- $$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b > 0, y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b < 0, y_i = -1 \end{cases}$$
- Linearly separable sample with high confidence and accuracy

$$\begin{cases} \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|} \geq d, y_i = +1 \\ \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|} \leq -d, y_i = -1 \end{cases}$$

$$\begin{aligned} \Rightarrow \begin{cases} \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|d} \geq 1, y_i = +1 \\ \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|d} \leq -1, y_i = -1 \end{cases} & \xrightarrow[\text{Normalization}]{\|\mathbf{w}\|d=1} \begin{cases} \mathbf{w}^T \mathbf{x}_i + w_0 \geq 1, y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + w_0 \leq -1, y_i = -1 \end{cases} \end{aligned}$$



$$y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1,$$

Then  $|\mathbf{w}^T \mathbf{x}_i + w_0| = y_i(\mathbf{w}^T \mathbf{x}_i + w_0)$

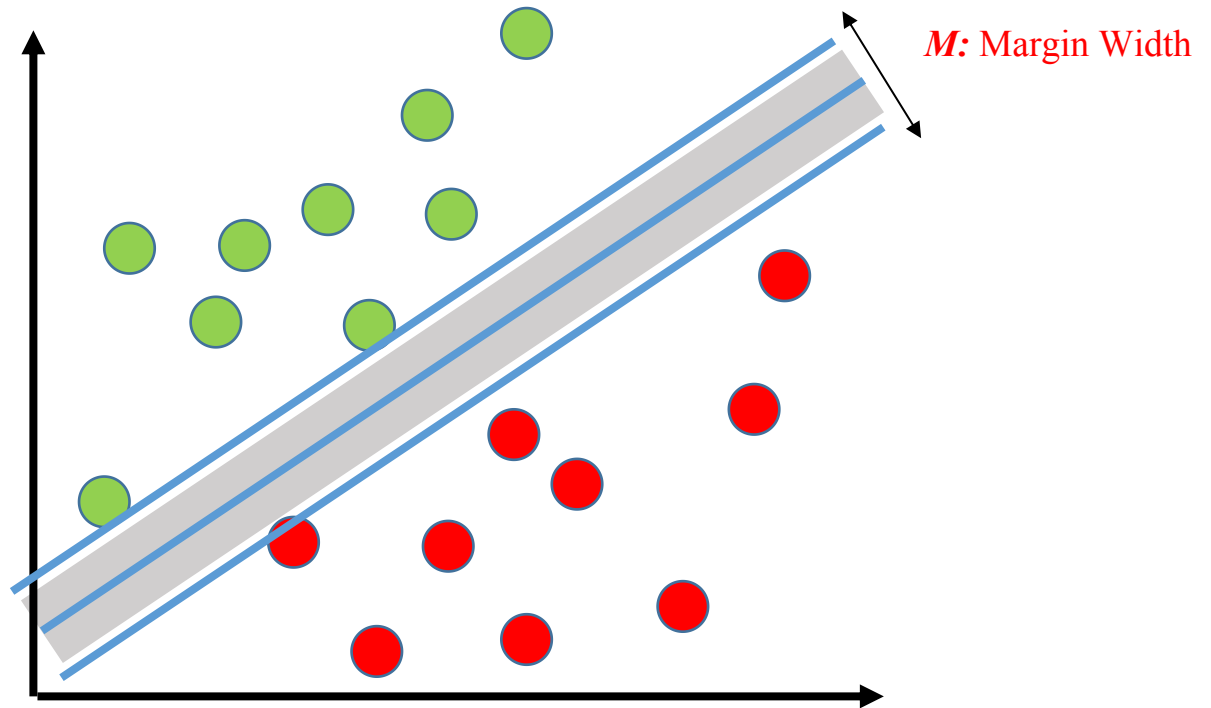
# Classification

## □ SVM

Goal:

1) Correctly classify all training data:

$$y (w^T x + b) \geq 1 \text{ for all } y$$



# Classification

## □ SVM

- Margin:

$$r = 2d = \frac{2}{\|w\|}$$

$$w^T(x_i^+ - x_i^-) = 2$$

$$\begin{aligned} w^T(x_i^+ - x_i^-) &= \|w\| \cdot \|x_i^+ - x_i^-\| \cdot \cos\theta \\ &= w^T(x_i^+ - x_i^-) = \|w\| \cdot \|x_i^+ - x_i^-\| \cdot \frac{r}{\|x_i^+ - x_i^-\|} = \|w\| \cdot r \end{aligned}$$

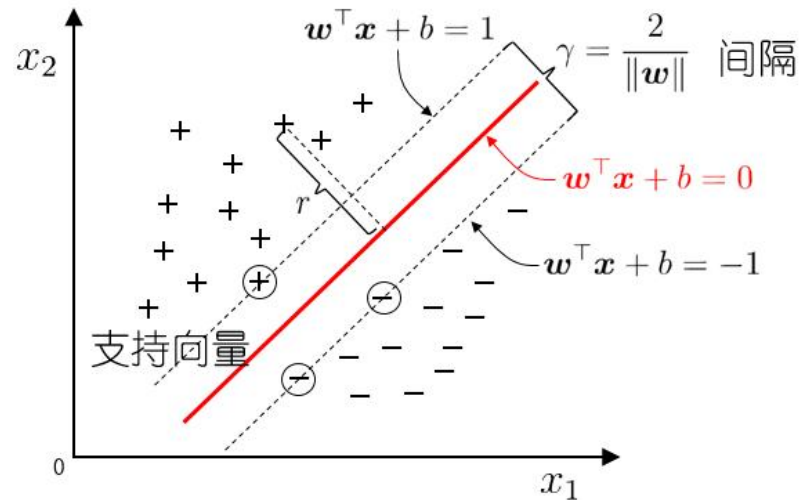
- **Maximum margin:** solve  $w$  and  $b$  to get maximum margin  $\gamma$

$$\max_{w,b} \frac{2}{\|w\|} \text{ s.t. } y_i(w^T x_i + b) \geq 1, i = 1, 2, \dots, m$$



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

$$\text{s.t. } y_i(w^T x_i + b) \geq 1, i = 1, 2, \dots, m$$



# Classification

## □ SVM

Goal:

- 1) Correctly classify all training data:

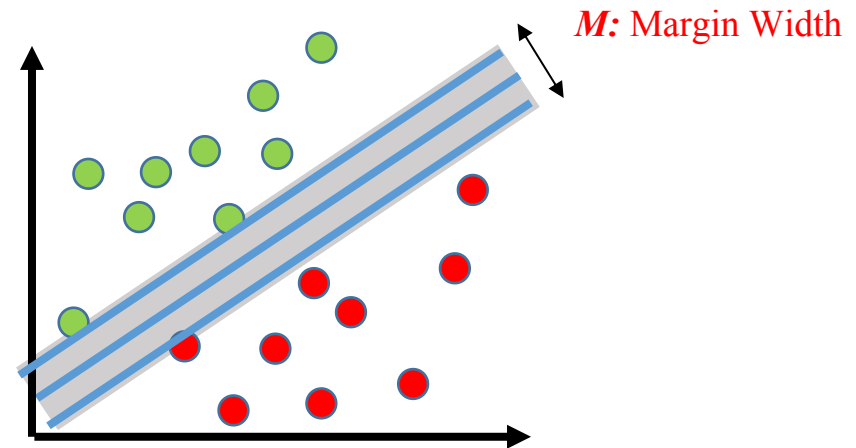
$$y_i (w^T x_i + b) \geq 1 \text{ for all } i$$

- 2) Maximize the margin

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

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

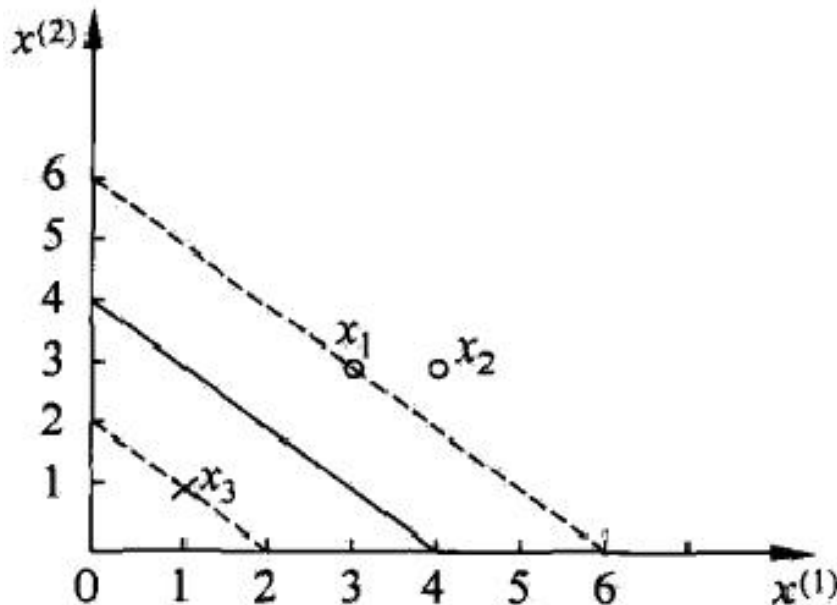
$$\begin{aligned} &\text{Minimize } \Phi(w) = \frac{1}{2} w^T w \\ &\text{subject to } y_i (w x_i + b) \geq 1 \text{ for all } i \end{aligned}$$



The primary problem of SVM: Solve  $d+1$  variables ( $w, b$ ) with  $m$  inequality constraints, which is suitable for low dimensions.

# Example 1——Solve the primary problem of SVM

- Problem: Given a training set size of 3, in which  $(x_1, y_1)$  and  $(x_2, y_2)$  are positive samples, and  $(x_3, y_3)$  is negative sample.  $x_1 = (3; 3)$ ,  $y_1 = +1$ ,  $x_2 = (4; 3)$ ,  $y_2 = +1$ ;  $x_3 = (1; 1)$ ,  $y_3 = -1$ . Solve the optimal separating hyperplane  $H$  with maximum margin.



Solve  $\mathbf{w} = (w_1; w_2)$  and  $w_0$  ( $w_0$  is  $b$ )

# Example 1——Solve the primary problem of SVM

Answer:

- Step 1: Build the primary problem of SVM upon the sample set

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} (w_1^2 + w_2^2) \\ \text{s. t. } & \begin{cases} 3w_1 + 3w_2 + w_0 \geq 1 \\ 4w_1 + 3w_2 + w_0 \geq 1 \\ -w_1 - w_2 + w_0 \geq 1 \end{cases} \end{aligned}$$

$$\begin{aligned} & \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t. } & y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, \quad i = 1, 2, \dots, m \end{aligned}$$



# Example 1——Solve the primary problem of SVM

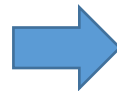
- Step 2: Build Lagrange function by setting Lagrange multiplier  $a_i \geq 0$  for each inequality constraint

$$L(w_1, w_2, w_0, a_1, a_2, a_3)$$

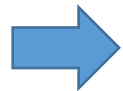
$$= \frac{1}{2}(w_1^2 + w_2^2) - a_1(3w_1 + 3w_2 + w_0 - 1) - a_2(4w_1 + 3w_2 + w_0 - 1) - a_3(-w_1 - w_2 + w_0 - 1)$$

Set the partial as 0

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_1} = w_1 - 3a_1 - 4a_2 + a_3 = 0 \\ \frac{\partial L}{\partial w_2} = w_2 - 3a_1 - 3a_2 + a_3 = 0 \\ \frac{\partial L}{\partial w_0} = -a_1 - a_2 + a_3 = 0 \\ a_1(3w_1 + 3w_2 + w_0 - 1) = 0 \\ a_2(4w_1 + 3w_2 + w_0 - 1) = 0 \\ a_3(-w_1 - w_2 + w_0 - 1) = 0 \\ 3w_1 + 3w_2 + w_0 - 1 \geq 0 \\ 4w_1 + 3w_2 + w_0 - 1 \geq 0 \\ -w_1 - w_2 + w_0 - 1 \geq 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{array} \right.$$



$$\left\{ \begin{array}{l} w_1 = 3a_1 + 4a_2 - a_3 \\ w_2 = 3a_1 + 3a_2 - a_3 \\ a_3 = a_1 + a_2 \\ a_1(3w_1 + 3w_2 + w_0 - 1) = 0 \\ a_2(4w_1 + 3w_2 + w_0 - 1) = 0 \\ a_3(-w_1 - w_2 + w_0 - 1) = 0 \\ 3w_1 + 3w_2 + w_0 - 1 \geq 0 \\ 4w_1 + 3w_2 + w_0 - 1 \geq 0 \\ -w_1 - w_2 + w_0 - 1 \geq 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{array} \right.$$



# Example 1——Solve the primary problem of SVM

- Step 3: only keep  $a_1, a_2, w_0$

$$\left\{ \begin{array}{l} w_1 = 3a_1 + 4a_2 - a_3 = 2a_1 + 3a_2 \\ w_2 = 3a_1 + 3a_2 - a_3 = 2a_1 + 2a_2 \\ a_3 = a_1 + a_2 \\ a_1(12a_1 + 15a_2 + w_0 - 1) = 0 \\ a_2(14a_1 + 18a_2 + w_0 - 1) = 0 \\ (a_1 + a_2)(-4a_1 - 5a_2 - w_0 - 1) = 0 \\ 12a_1 + 15a_2 + w_0 - 1 \geq 0 \\ 14a_1 + 18a_2 + w_0 - 1 \geq 0 \\ -4a_1 - 5a_2 - w_0 - 1 \geq 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{array} \right.$$

(1) If  $a_1 = 0, a_2 = 0$ , then

$$w_1 = w_2 = 0$$

However,  $w_0 \geq 1$  and  $w_0 \leq -1$  are contradictory!!!

(2) If  $a_1 > 0, a_2 = 0$ , then

$$\begin{cases} 12a_1 + w_0 - 1 = 0 \\ -4a_1 - w_0 - 1 = 0 \end{cases}$$

We will have

$$\begin{cases} a_1 = a_3 = \frac{1}{4} \\ a_2 = 0 \\ w_0 = -2 \\ w_1 = w_2 = \frac{1}{2} \end{cases}$$

Which satisfy all inequality constraints!

# Example 1——Solve the primary problem of SVM

- Step 3: only keep  $a_1, a_2, w_0$

$$\left\{ \begin{array}{l} w_1 = 3a_1 + 4a_2 - a_3 = 2a_1 + 3a_2 \\ w_2 = 3a_1 + 3a_2 - a_3 = 2a_1 + 2a_2 \\ a_3 = a_1 + a_2 \\ a_1(12a_1 + 15a_2 + w_0 - 1) = 0 \\ a_2(14a_1 + 18a_2 + w_0 - 1) = 0 \\ (a_1 + a_2)(-4a_1 - 5a_2 - w_0 - 1) = 0 \\ 12a_1 + 15a_2 + w_0 - 1 \geq 0 \\ 14a_1 + 18a_2 + w_0 - 1 \geq 0 \\ -4a_1 - 5a_2 - w_0 - 1 \geq 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{array} \right.$$

(3) If  $a_1 = 0, a_2 > 0$ , then

$$\begin{cases} 18a_2 + w_0 - 1 = 0 \\ -5a_2 - w_0 - 1 = 0 \end{cases}$$

We will have

$$\left\{ \begin{array}{l} a_1 = 0 \\ a_2 = a_3 = \frac{2}{13} \\ w_0 = -\frac{23}{13} \\ w_1 = \frac{6}{13} \\ w_2 = \frac{4}{13} \end{array} \right.$$

Substitute in the inequality constraints to get the following expression

$$12a_1 + 15a_2 + w_0 - 1 = -\frac{6}{13} < 0$$

Which violates the constraint!

# Example 1——Solve the primary problem of SVM

- Step 3: only keep  $a_1, a_2, w_0$

$$\left\{ \begin{array}{l} w_1 = 3a_1 + 4a_2 - a_3 = 2a_1 + 3a_2 \\ w_2 = 3a_1 + 3a_2 - a_3 = 2a_1 + 2a_2 \\ \quad a_3 = a_1 + a_2 \\ a_1(12a_1 + 15a_2 + w_0 - 1) = 0 \\ a_2(14a_1 + 18a_2 + w_0 - 1) = 0 \\ (a_1 + a_2)(-4a_1 - 5a_2 - w_0 - 1) = 0 \\ 12a_1 + 15a_2 + w_0 - 1 \geq 0 \\ 14a_1 + 18a_2 + w_0 - 1 \geq 0 \\ -4a_1 - 5a_2 - w_0 - 1 \geq 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{array} \right.$$

(4) If  $a_1 > 0, a_2 > 0$ , then

$$\left\{ \begin{array}{l} 12a_1 + 15a_2 + w_0 - 1 = 0 \\ 14a_1 + 18a_2 + w_0 - 1 = 0 \\ -4a_1 - 5a_2 - w_0 - 1 = 0 \end{array} \right.$$

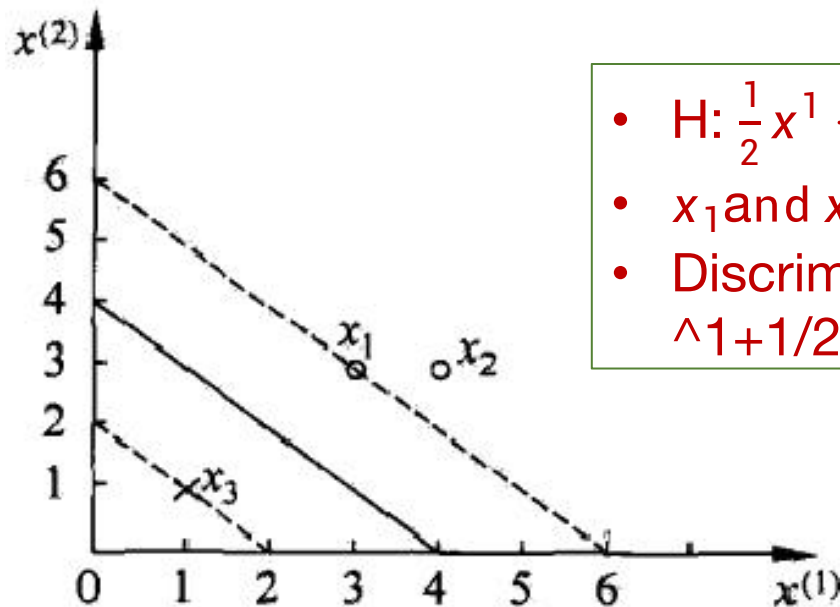
We will have

$$\left\{ \begin{array}{l} a_1 = \frac{3}{2} \\ a_2 = -1 < 0 \\ a_3 = \frac{1}{2} \\ w_0 = -2 \\ w_1 = 0 \\ w_2 = 1 \end{array} \right.$$

Which violates the constraint!

# Example 1——Solve the primary problem of SVM

- Problem: Given a training set size of 3, in which  $(x_1, y_1)$  and  $(x_2, y_2)$  are positive samples, and  $(x_3, y_3)$  is negative sample.  $x_1 = (3; 3)$ ,  $y_1 = +1$ ,  $x_2 = (4; 3)$ ,  $y_2 = +1$ ;  $x_3 = (1; 1)$ ,  $y_3 = -1$ . Solve the optimal separating hyperplane  $H$  with maximum margin.



- $H: \frac{1}{2}x^1 + \frac{1}{2}x^2 - 2 = 0$
- $x_1$  and  $x_3$  are support vectors:  $y_i g(x_i) = 1$
- Discriminant Function:  $g(x) = \text{sign}(\frac{1}{2}x^1 + \frac{1}{2}x^2 - 2)$

- 
- $m$  inequality constraints means  $2^m$  cases!!

Primary Problem → Duel  
Problem

# SVM-Duel Problem

- The primary problem of SVM

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s. t. } y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, i = 1, 2, \dots, m$$

- Lagrange function: Lagrange multiplier  $\mathbf{a} = (a_1, a_2, \dots, a_m)^T$ :

$$L(\mathbf{w}, w_0, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))$$

- **Primary problem → Duel problem** (maxi-mini problem)

$$\max_{\mathbf{a}} \min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$$

- To get the solution of SVM duel problem:

➤ Solve the minimum of  $L(\mathbf{w}, b, \mathbf{a})$  on  $\mathbf{w}, w_0$ :  $\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$

➤ Solve the maximum of  $\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$  on  $\mathbf{a}$ :  $\max_{\mathbf{a}} \min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$

# SVM-Duel Problem

$$\begin{aligned} L(\mathbf{w}, w_0, \mathbf{a}) &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i + w_0)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^m a_i - \sum_{i=1}^m a_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^m a_i y_i w_0 \end{aligned}$$

- (1) Solve  $\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$ , set the partial of  $L(\mathbf{w}, w_0, \mathbf{a})$  on  $\mathbf{w}$  and  $w_0$  as 0

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m a_i y_i \mathbf{x}_i = 0, \text{ then } \mathbf{w} = \sum_{i=1}^m a_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial w_0} = - \sum_{i=1}^m a_i y_i = 0, \text{ then } \sum_{i=1}^m a_i y_i = 0$$

Then,

$$\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a}) = \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$



# SVM-Duel Problem

- (2) Solve the maximum of  $\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$  on  $\mathbf{a}$

$$\max_{\mathbf{a}} \min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a}) = \max_{\mathbf{a}} \left( \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right)$$
$$\text{s. t. } \sum_{i=1}^m a_i y_i = 0, \quad a_i \geq 0, \quad i = 1, 2, \dots, m$$

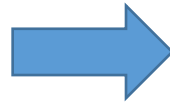
Substitute  $a_i$  after solving to get  $\mathbf{w}$  and  $w_0$ ,  $w_0 = y_j - \sum_{i=1}^m a_i y_i \mathbf{x}_i^T \mathbf{x}_j$

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_{i=1}^m a_i y_i \mathbf{x}_i^T \mathbf{x} + w_0$$

# SVM-Duel Problem——Solution Sparsity

- Final SVM:  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_{i=1}^m a_i y_i \mathbf{x}_i^T \mathbf{x} + w_0$
- SVM satisfies KKT(Karush-Kuhn-Tucker) conditions:

$$\begin{cases} a_i \geq 0 \\ y_i g(\mathbf{x}_i) \geq 1 \\ a_i (y_i g(\mathbf{x}_i) - 1) = 0 \end{cases}$$



For any sample  $(\mathbf{x}_i, y_i)$ , there must exist  $a_i = 0$  or  $y_i g(\mathbf{x}_i) = 1$

- If  $a_i = 0$ , then  $y_i g(\mathbf{x}_i) > 1$ ,  $(\mathbf{x}_i, y_i)$  does not affect SVM  $g(\mathbf{x})$ .
- If  $a_i > 0$ , then  $y_i g(\mathbf{x}_i) = 1$ ,  $(\mathbf{x}_i, y_i)$  is on the boundary hyperplane H1 or H2, which is the support vector.

**Solution Sparsity of SVM:** After training, most of the training samples are not reserved. That is, the final SVM only concerns support vectors which is small amount.

For duel problem ,we only need to solve support vectors and corresponding multiplier  $a$ .

## Example 2——Solve the dual problem of SVM

- Problem: Given a training set size of 3, in which  $(x_1, y_1)$  and  $(x_2, y_2)$  are positive samples, and  $(x_3, y_3)$  is negative sample.  $x_1 = (3; 3)$ ,  $y_1 = +1$ ,  $x_2 = (4; 3)$ ,  $y_2 = +1$ ;  $x_3 = (1; 1)$ ,  $y_3 = -1$ . Solve the linearly separable SVM.

# Example 2——Solve the dual problem of SVM

- SVM Primary Problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} (w_1^2 + w_2^2) \\ \text{s. t.} \quad & \begin{cases} 3w_1 + 3w_2 + w_0 \geq 1 \\ 4w_1 + 3w_2 + w_0 \geq 1 \\ -w_1 - w_2 + w_0 \geq 1 \end{cases} \end{aligned}$$

$$\max_{\mathbf{a}} \min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathbf{a})$$

$$\begin{aligned} &= \max_{\mathbf{a}} \left( \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ \text{s. t.} \quad & \sum_{i=1}^m a_i y_i = 0, \quad a_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned}$$

- Step1: transform to SVM dual problem

$$\begin{aligned} \max_{a_1, a_2, a_3} \quad & \left( \sum_{i=1}^3 a_i - \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 a_i a_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ \text{s. t.} \quad & \begin{cases} a_1 + a_2 - a_3 = 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{cases} \end{aligned}$$



$$\begin{aligned} \min_{a_1, a_2, a_3} \quad & \left( \frac{1}{2} (18a_1^2 + 25a_2^2 + 2a_3^2 + 42a_1a_2 - 12a_1a_3 - 14a_2a_3) - a_1 - a_2 - a_3 \right) \\ \text{s. t.} \quad & \begin{cases} a_1 + a_2 - a_3 = 0 \\ a_1 \geq 0, a_2 \geq 0, a_3 \geq 0 \end{cases} \end{aligned}$$

# Example 2——Solve the dual problem of SVM

- Step 2: Substitute constraints  $a_3 = a_1 + a_2$ , then the objective function is:

$$s(a_1, a_2) = 4a_1^2 + \frac{13}{2}a_2^2 + 10a_1a_2 - 2a_1 - 2a_2$$
$$\text{s. t. } a_1 \geq 0, a_2 \geq 0, a_3 \geq 0$$

Solve the partial of  $s(a_1, a_2)$  on  $a_1, a_2$ , and set as 0:

$$\begin{cases} \frac{\partial s}{\partial a_1} = 8a_1 + 10a_2 - 2 = 0 \\ \frac{\partial s}{\partial a_2} = 13a_2 + 10a_1 - 2 = 0 \end{cases}$$

$$\text{then } \begin{cases} a_1 = \frac{3}{2} \\ a_2 = -1 < 0, \\ a_3 = \frac{1}{2} \end{cases}$$

This violates constraints! We will find the minimum on boundary value of  $a_i$ .

## Example 2——Solve the dual problem of SVM

- Step 3: Solve vector  $\mathbf{w}$  with KKT conditions

(1) When  $a_1 = 0$ ,

$$s(0, a_2) = \frac{13}{2} a_2^2 - 2a_2,$$

$$\text{Set } \frac{\partial s}{\partial a_2} = 13a_2 - 2 = 0, \text{ then } a_2 = \frac{2}{13}, \quad s_{min} = -\frac{2}{13}$$

(2) When  $a_2 = 0$ ,

$$s(a_1, 0) = 4a_1^2 - 2a_1,$$

$$\text{Set } \frac{\partial s}{\partial a_1} = 8a_1 - 2 = 0, \text{ then } a_1 = \frac{1}{4}, \quad s_{min} = -\frac{1}{4}$$

$$\text{Thus, when } a_1 = \frac{1}{4}, a_2 = 0, a_3 = \frac{1}{4}, \quad s_{min} = -\frac{1}{4}$$

$$\text{And } \mathbf{w} = \sum_{i=1}^3 a_i y_i \mathbf{x}_i = a_1 y_1 \mathbf{x}_1 + a_3 y_3 \mathbf{x}_3 = \left(\frac{1}{2}, \frac{1}{2}\right)$$

## Example 2——Solve the dual problem of SVM

- Step 4: Solve  $w_0$  with KKT conditions

Since  $a_1 = \frac{1}{4} > 0$ , the corresponding sample  $\mathbf{x}_1$  is the support vector, then  $y_1 g(\mathbf{x}_1) = 1$  and  $w_0 = -2$ .

- Hyperplane H (SVM)

$$\frac{1}{2}x^1 + \frac{1}{2}x^2 - 2 = 0$$

- For new sample, discriminant function is

$$g(\mathbf{x}) = \text{sign}\left(\frac{1}{2}x^1 + \frac{1}{2}x^2 - 2\right)$$

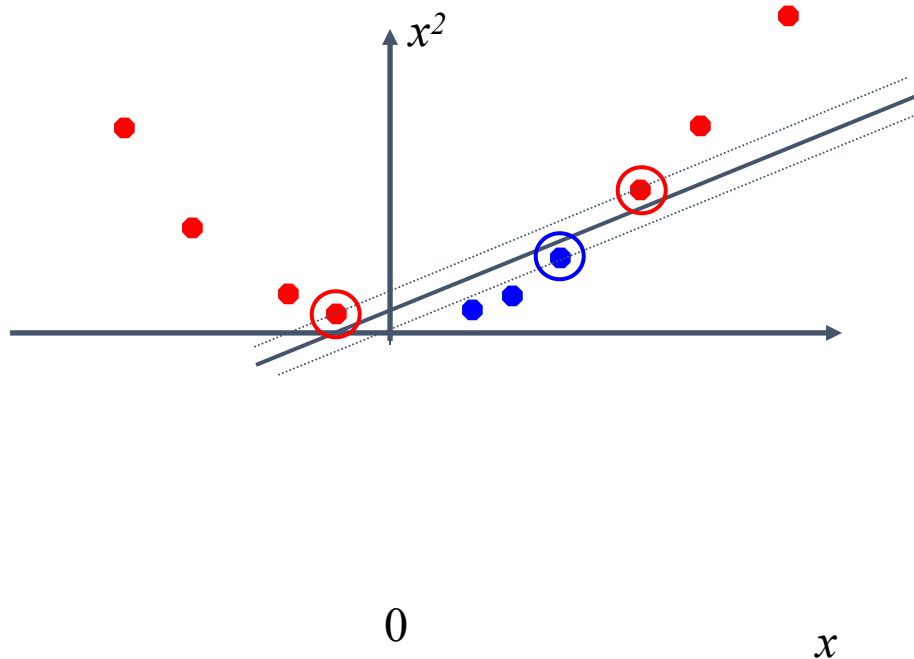
# Classification

## □ NON-linear SVMs

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



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

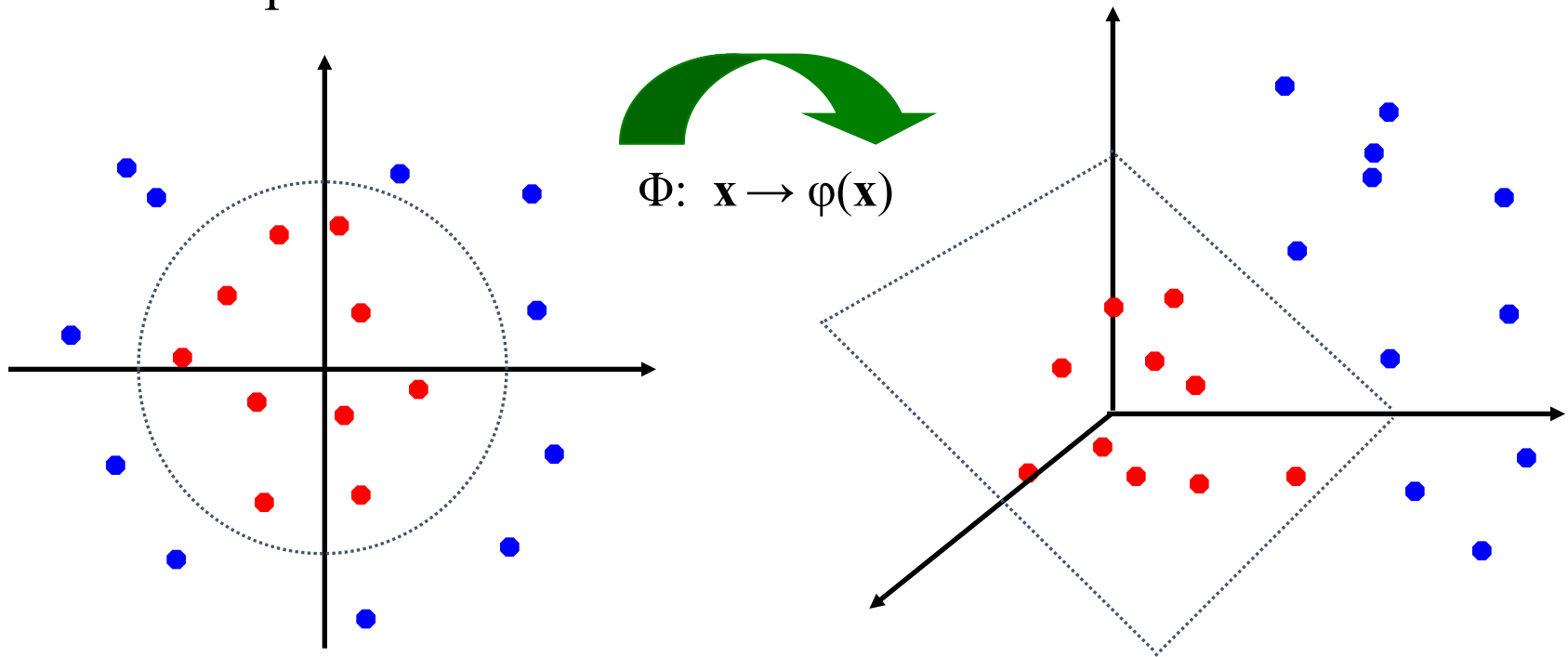




# Classification

## □ NON-linear SVMs

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



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

# Classification

## □ Weakness of SVM

- It is sensitive to noise
- It only considers two classes
  - how to do multi-class classification with SVM?

Suppose there are  $m$  different classes,

### 1) OVA-SVM : ( $m$ SVMs in total)

- $SVM_1$  learns “Output==1” vs “Output != 1”
- $SVM_2$  learns “Output==2” vs “Output != 2”
- $\vdots$
- $SVM_m$  learns “Output== $m$ ” vs “Output !=  $m$ ”

### 2) OVO-SVM: ( $m(m-1)/2$ SVMs in total)

- $SVM_{12}$  learns “Output==1” vs “Output == 2”
- $SVM_{13}$  learns “Output==1” vs “Output == 3”
- $\vdots$
- $SVM_{m(m-1)}$  learns “Output== $m$ ” vs “Output ==  $m-1$ ”

# Classification

---

## □ Other classifiers

- Decision trees

- Sparse representation classifier

- Neural networks

- .....

# Conclusion

---

- Linear Regression
- Logistic Regression
- Classification
  - Distance-based algorithms
  - Linear classifiers
  - Other classifiers
- .....

