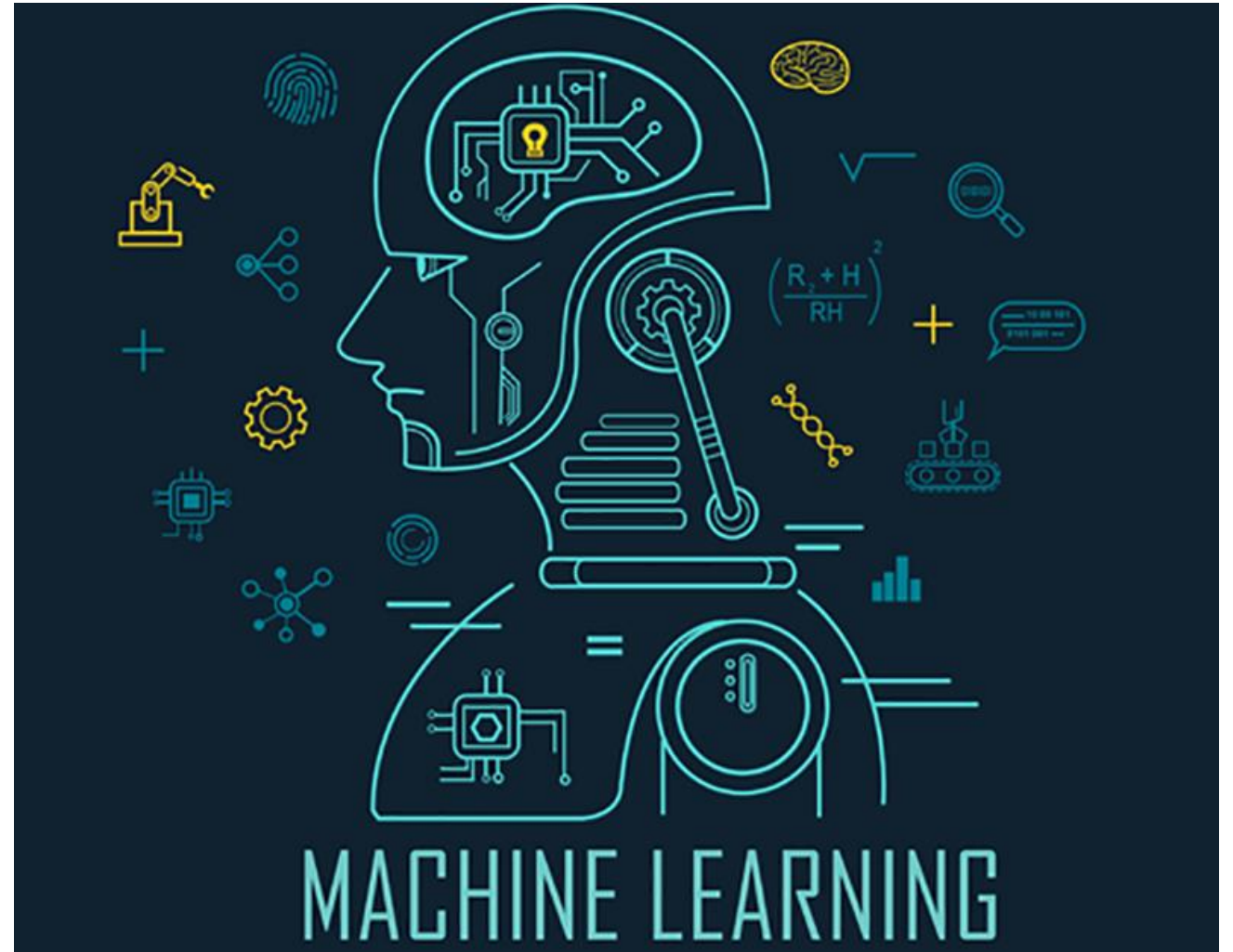


Logistic Regression

Classification

Zahoor Tanoli (PhD)

COMSATS Attock



Classification

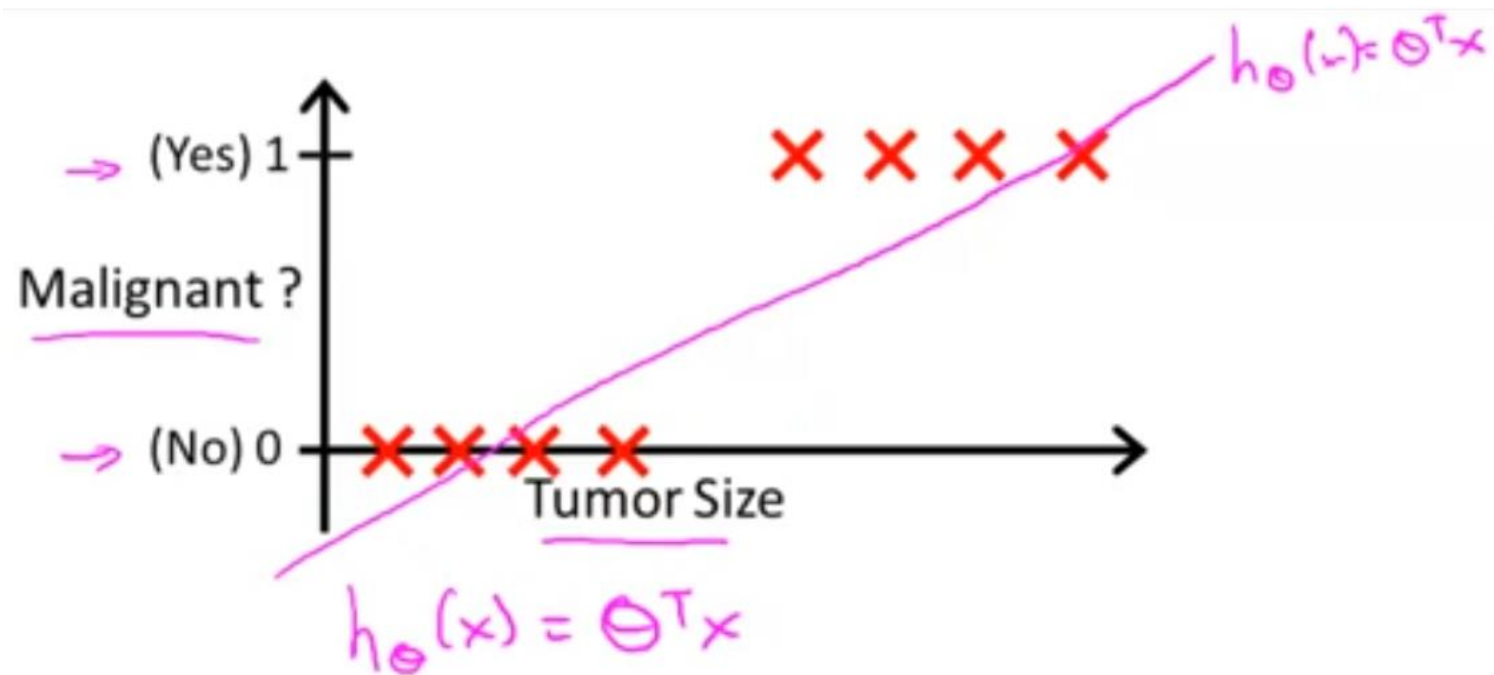
- Email: Spam / Not Spam?
- Online Transactions: Fraudulent (Yes / No)?
- Tumor: Malignant / Benign ?

→ $y \in \{0, 1\}$

0: "Negative Class" (e.g., benign tumor)

1: "Positive Class" (e.g., malignant tumor)

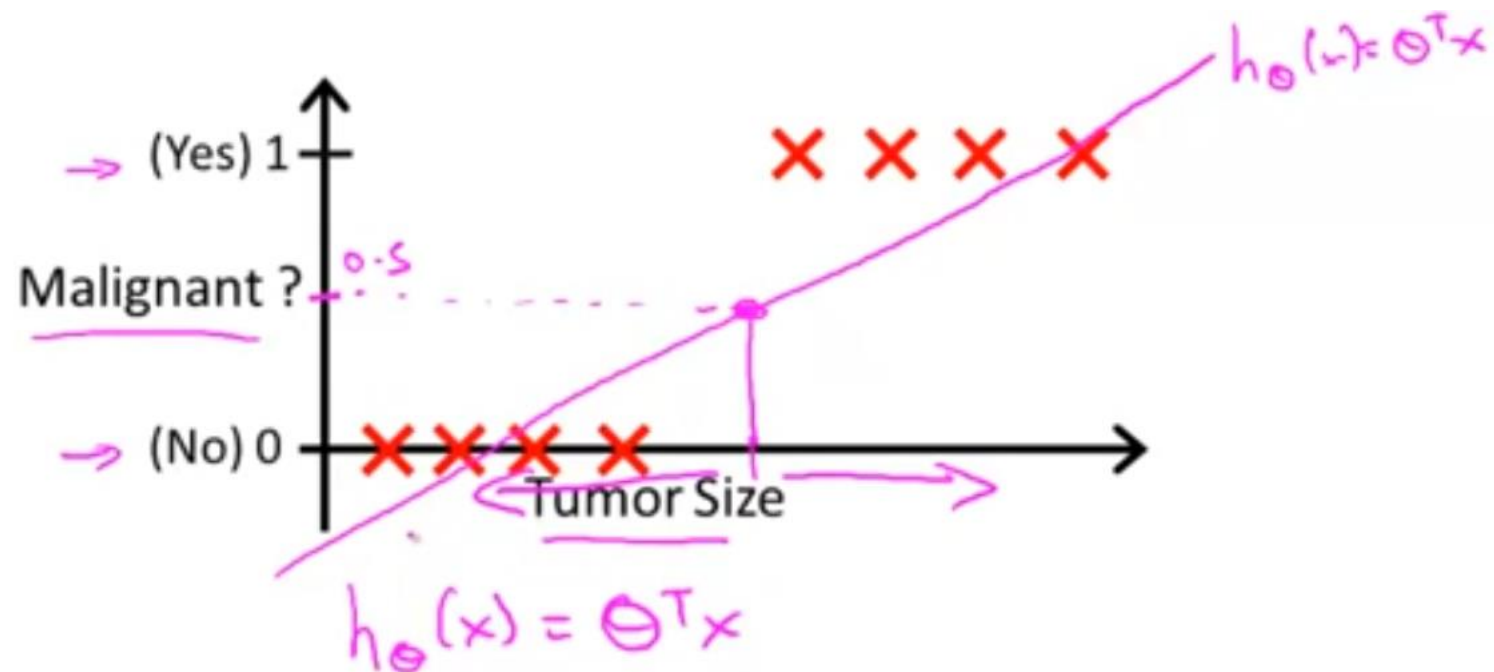
→ $y \in \{0, 1, 2, 3\}$



→ Threshold classifier output $h_{\theta}(x)$ at 0.5:

If $h_{\theta}(x) \geq 0.5$, predict "y = 1"

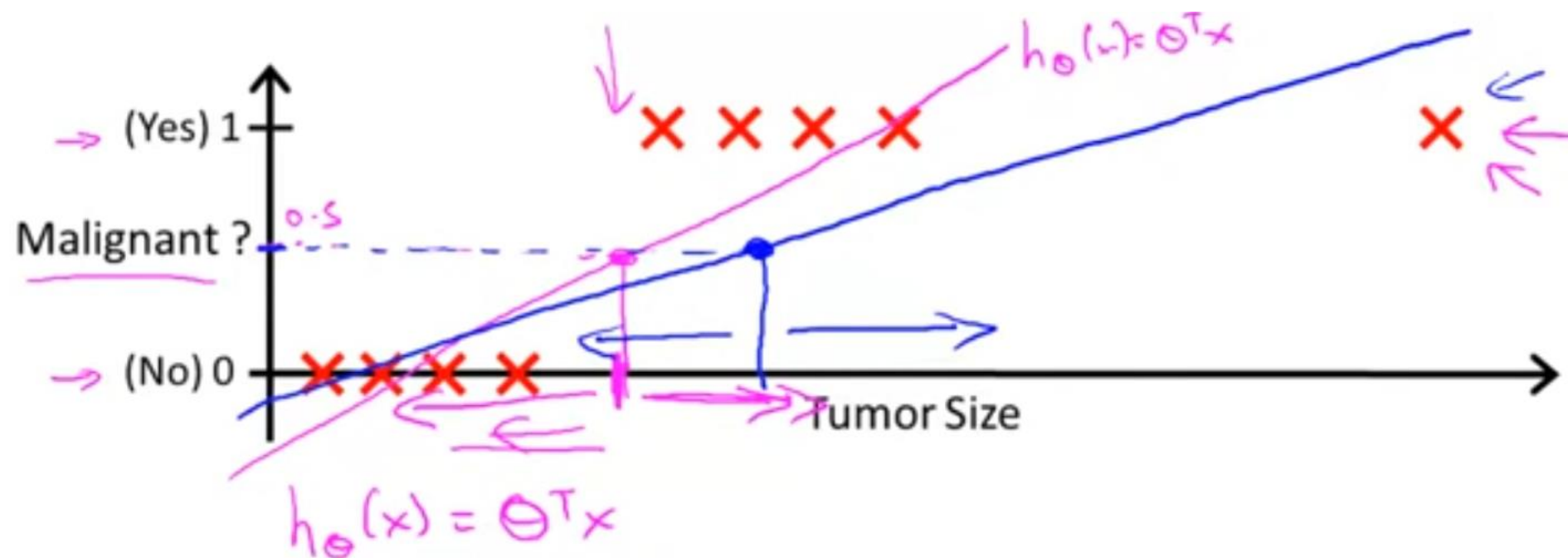
If $h_{\theta}(x) < 0.5$, predict "y = 0"



→ Threshold classifier output $h_{\theta}(x)$ at 0.5:

→ If $h_{\theta}(x) \geq 0.5$, predict "y = 1"

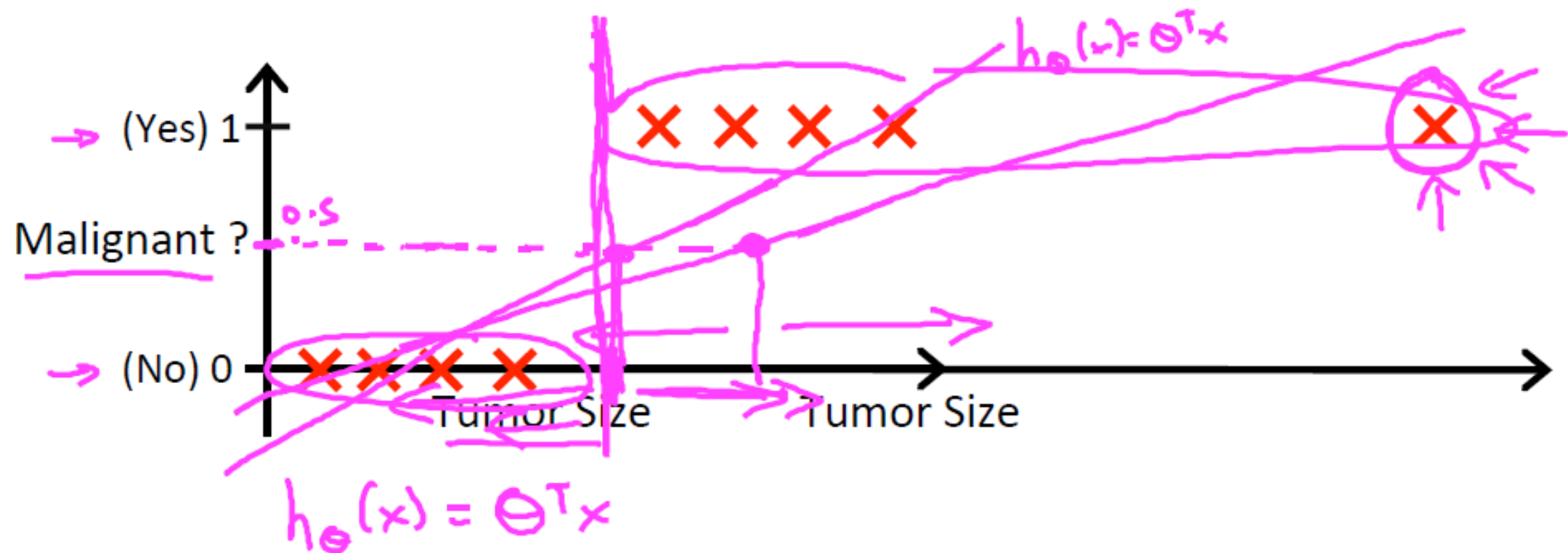
If $h_{\theta}(x) < 0.5$, predict "y = 0"



→ Threshold classifier output $h_{\theta}(x)$ at 0.5:

→ If $h_{\theta}(x) \geq 0.5$, predict "y = 1"

If $h_{\theta}(x) < 0.5$, predict "y = 0"



→ Threshold classifier output $h_{\theta}(x)$ at 0.5:

→ If $h_{\theta}(x) \geq 0.5$, predict "y = 1"

If $h_{\theta}(x) < 0.5$, predict "y = 0"

Classification: $y = 0 \text{ or } 1$

$h_{\theta}(x)$ can be > 1 or < 0

Logistic Regression: $0 \leq h_{\theta}(x) \leq 1$

Classification

$\begin{bmatrix} \square & \square \\ R & L & B & P \\ 3 & 3 & 1 & 0 \end{bmatrix}$
 $\begin{bmatrix} \square & \square \\ R & L & B & P \\ 4 & 7 & 2 & 0 \end{bmatrix}$
 $\begin{bmatrix} \square & \square \\ R & L & B & P \\ 2 & 5 & 3 & 4 \end{bmatrix}$

$\begin{matrix} R & L & B \\ 3 & 6 & 2 \end{matrix}$

$$y = \theta^T X$$

$$y = P = \begin{bmatrix} w_0 & w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} 1 \\ R \\ L \\ B \end{bmatrix} = 7 + 8R + 3L + 2B$$

$$= 7 + 8 \times 3 + 3 \times 6 + 2 \times 2$$

$$= 7 + 24 + 18 + 4$$

$$= 53$$

What is Logistic Regression?

$$y = mx + b - w_0 + w_1 x$$

- ✓ A regression algorithm which does classification
- ✓ Calculates probability of belonging to a particular class
- ✓ If $p > 50\% \rightarrow 1$
- ✓ If $p < 50\% \rightarrow 0$

1000



Regression classification

Spam
Not Spam

	f_1	f_2	f_3	f_4	f_5	L
1	3	4	1	4	3	1
...
1000	0

Logistic Regression

Hypothesis Representation



Logistic Regression Model

Want $0 \leq h_{\theta}(x) \leq 1$

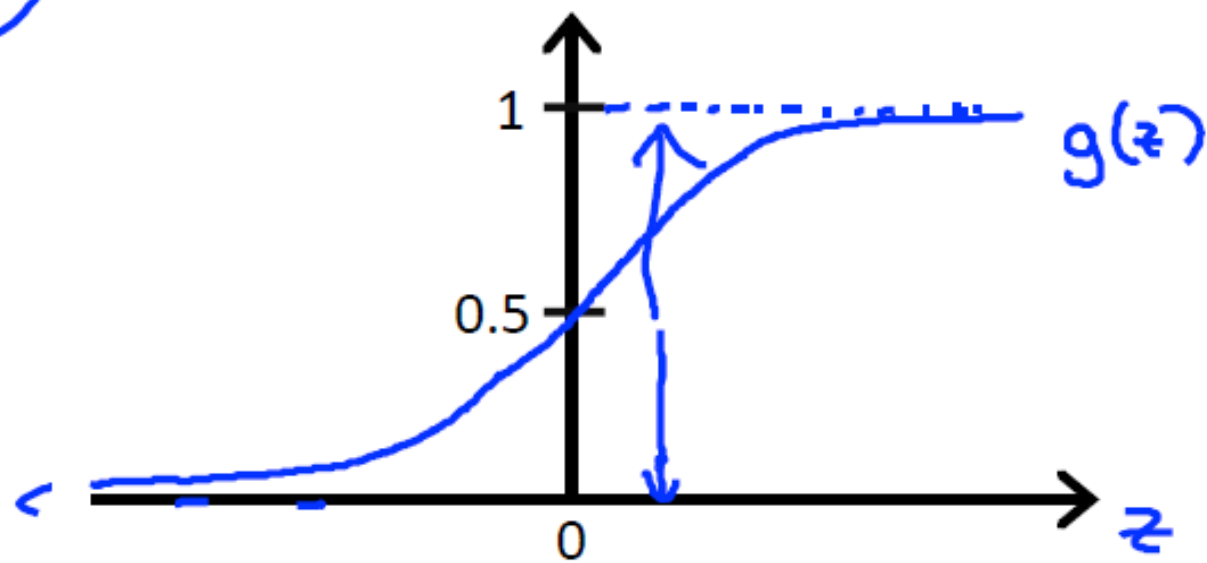
$$h_{\theta}(x) = g(\theta^T x)$$

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

$\theta^T x$

Sigmoid function
Logistic function

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



Parameters θ .

Interpretation of Hypothesis Output

$h_{\theta}(x)$

$h_{\theta}(x)$ = estimated probability that $y = 1$ on input x ←

Example: If x = $\begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$ = $\begin{bmatrix} 1 \leftarrow \\ \text{tumorSize} \end{bmatrix}$ ←

$$\text{↑} \quad \text{y=1}$$

$h_{\theta}(x)$ = 0.7

Tell patient that 70% chance of tumor being malignant

Interpretation of Hypothesis Output

$h_{\theta}(x)$

$h_{\theta}(x)$ = estimated probability that $y = 1$ on input x \leftarrow

Example: If $\underline{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \leftarrow \\ \text{tumorSize} \leftarrow \end{bmatrix}$

$$\underline{h_{\theta}(x)} = \underline{0.7}$$

$y = 1$

Tell patient that 70% chance of tumor being malignant

$$\underline{h_{\theta}(x)} = \underline{P(y=1|x;\theta)}$$

$y = 0 \text{ or } 1$

“probability that $y = 1$, given x ,
parameterized by θ ”

$$\begin{aligned} \rightarrow P(y=0|x;\theta) + P(y=1|x;\theta) &= 1 \\ P(\overline{y=0}|x;\theta) &= 1 - P(y=1|x;\theta) \end{aligned}$$

$$h_{\theta}(x)$$

Interpretation of Hypothesis Output

$h_{\theta}(x)$ = estimated probability that $y = 1$ on input x ←

Example: If x = $\begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$ = $\begin{bmatrix} 1 \leftarrow \\ \text{tumorSize} \leftarrow \end{bmatrix}$

$h_{\theta}(x)$ = 0.7

$y = 1$

Tell patient that 70% chance of tumor being malignant

$h_{\theta}(x) = P(y=1|x;\theta)$

$y = 0 \text{ or } 1$

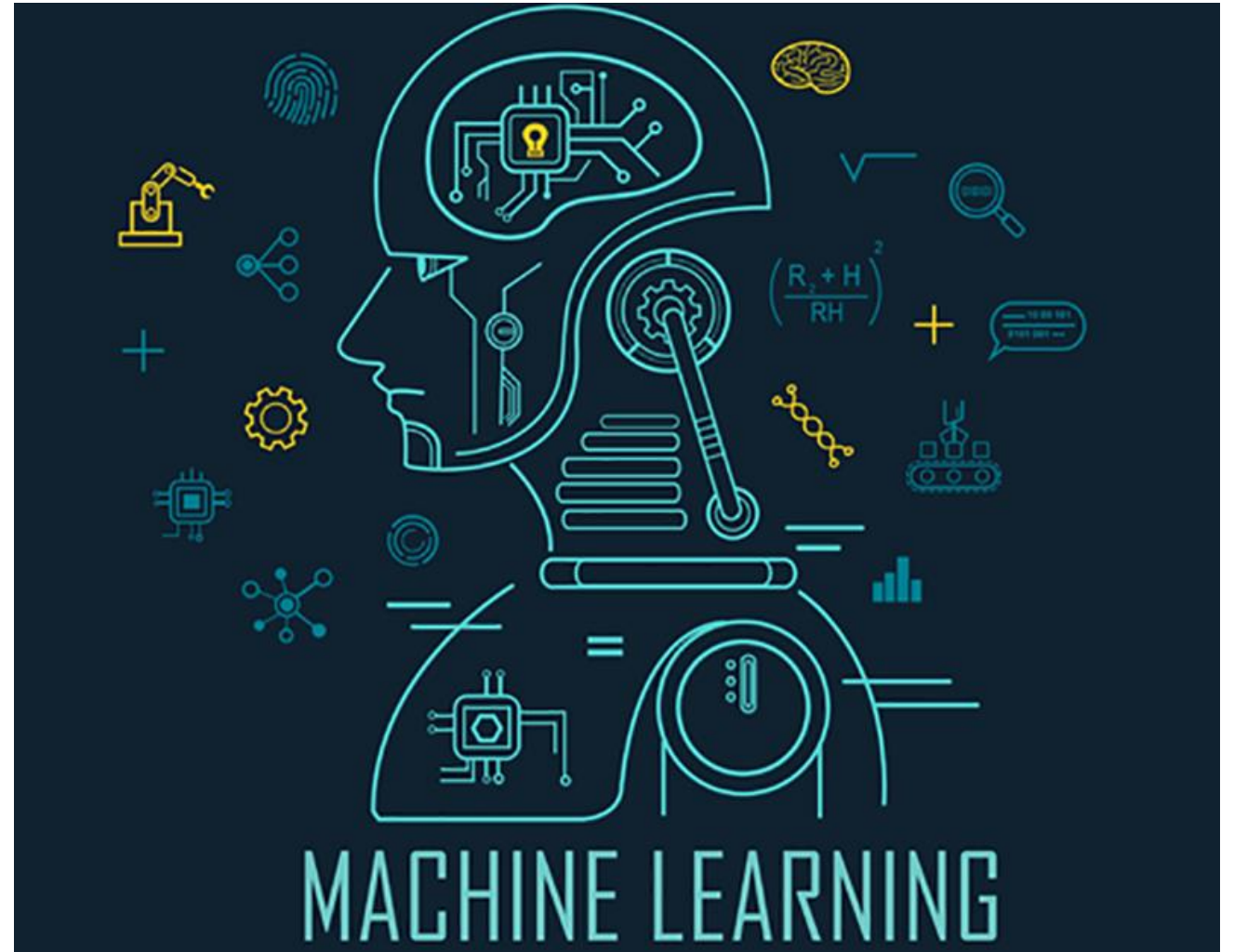
“probability that $y = 1$, given x ,
parameterized by θ ”

→ $P(y = 0|x;\theta) + P(y = 1|x;\theta) = 1$

→ $P(y = 0|x;\theta) = 1 - P(y = 1|x;\theta)$

Logistic Regression

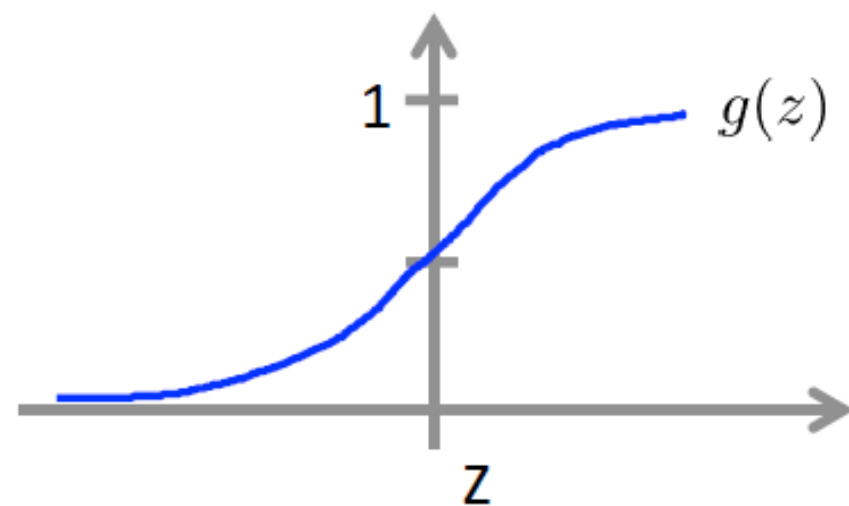
Decision Boundary



Logistic regression

$$h_{\theta}(x) = g(\theta^T x)$$

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



Suppose predict “ $y = 1$ ” if $h_{\theta}(x) \geq 0.5$

$$\theta^T x \geq 0$$

$$g(z) \geq 0.5$$

$$\text{when } z \geq 0$$

$$h_{\theta}(x) = g(\theta^T x)$$

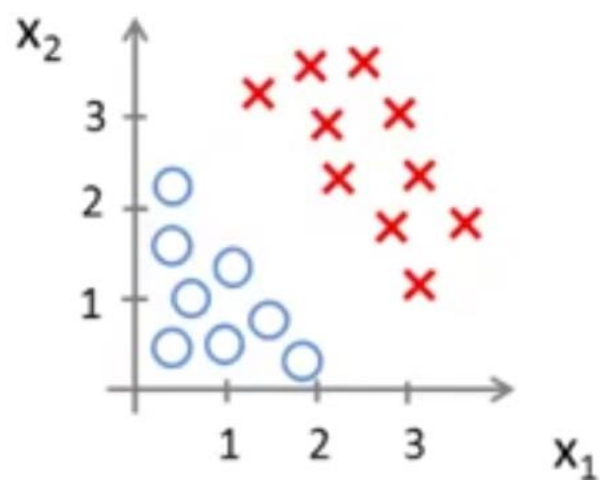
predict “ $y = 0$ ” if $h_{\theta}(x) < 0.5$

$$\theta^T x < 0$$

$$g(z) < 0.5$$

$$\text{when } z < 0$$

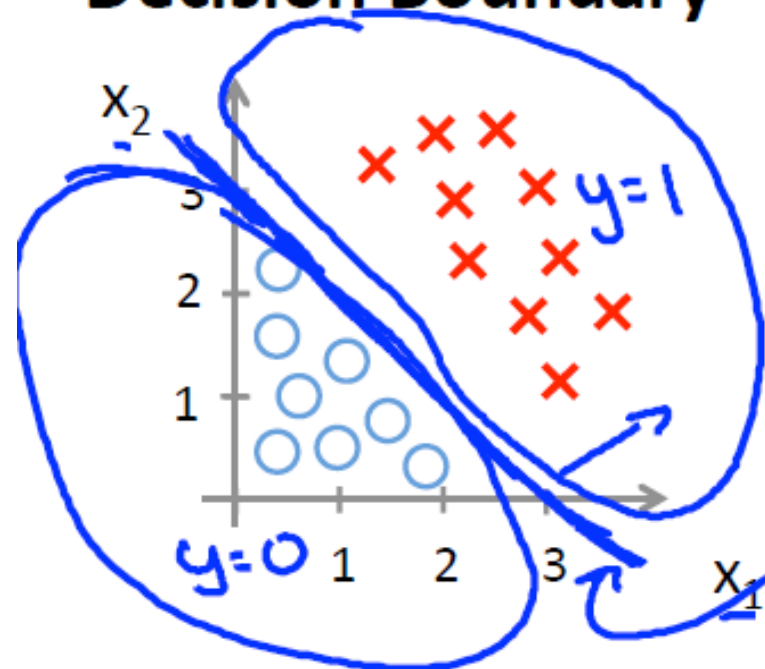
Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$\rightarrow h_{\theta}(x) = g(\underbrace{\theta_0}_{-3} + \underbrace{\theta_1}_{1}x_1 + \underbrace{\theta_2}_{1}x_2)$$

Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

$$h_{\theta}(x) = g(\underbrace{\theta_0}_{-3} + \underbrace{\theta_1}_{1}x_1 + \underbrace{\theta_2}_{1}x_2)$$

Decision boundary

Predict " $y = 1$ " if $-3 + x_1 + x_2 \geq 0$

$$\theta^T x$$

$$\rightarrow \underline{x_1 + x_2 \geq 3}$$

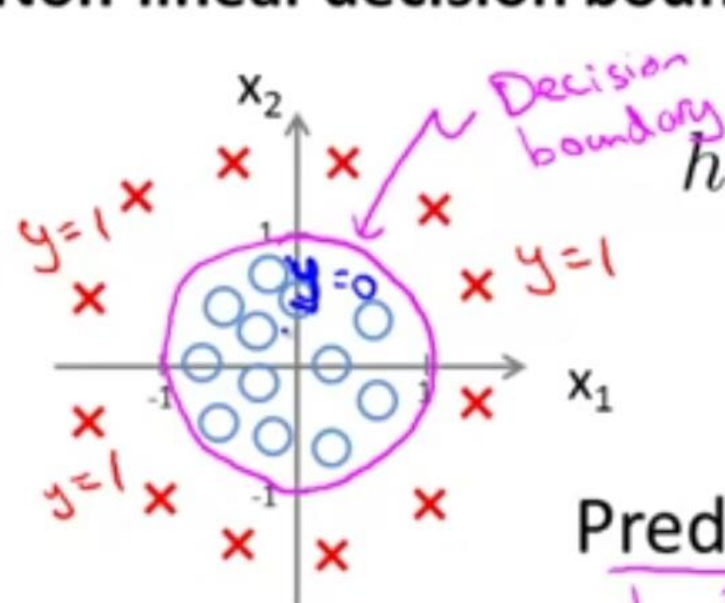
x_1, x_2

$$\rightarrow h_{\theta}(x) = 0.5$$

$$\boxed{x_1 + x_2 = 3}$$

$$\rightarrow \begin{matrix} x_1 + x_2 < 3 \\ y = 0 \end{matrix}$$

Non-linear decision boundaries



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$\theta_0 = -1$ $\theta_1 = 0$ $\theta_2 = 0$ $\theta_3 = 1$ $\theta_4 = 1$

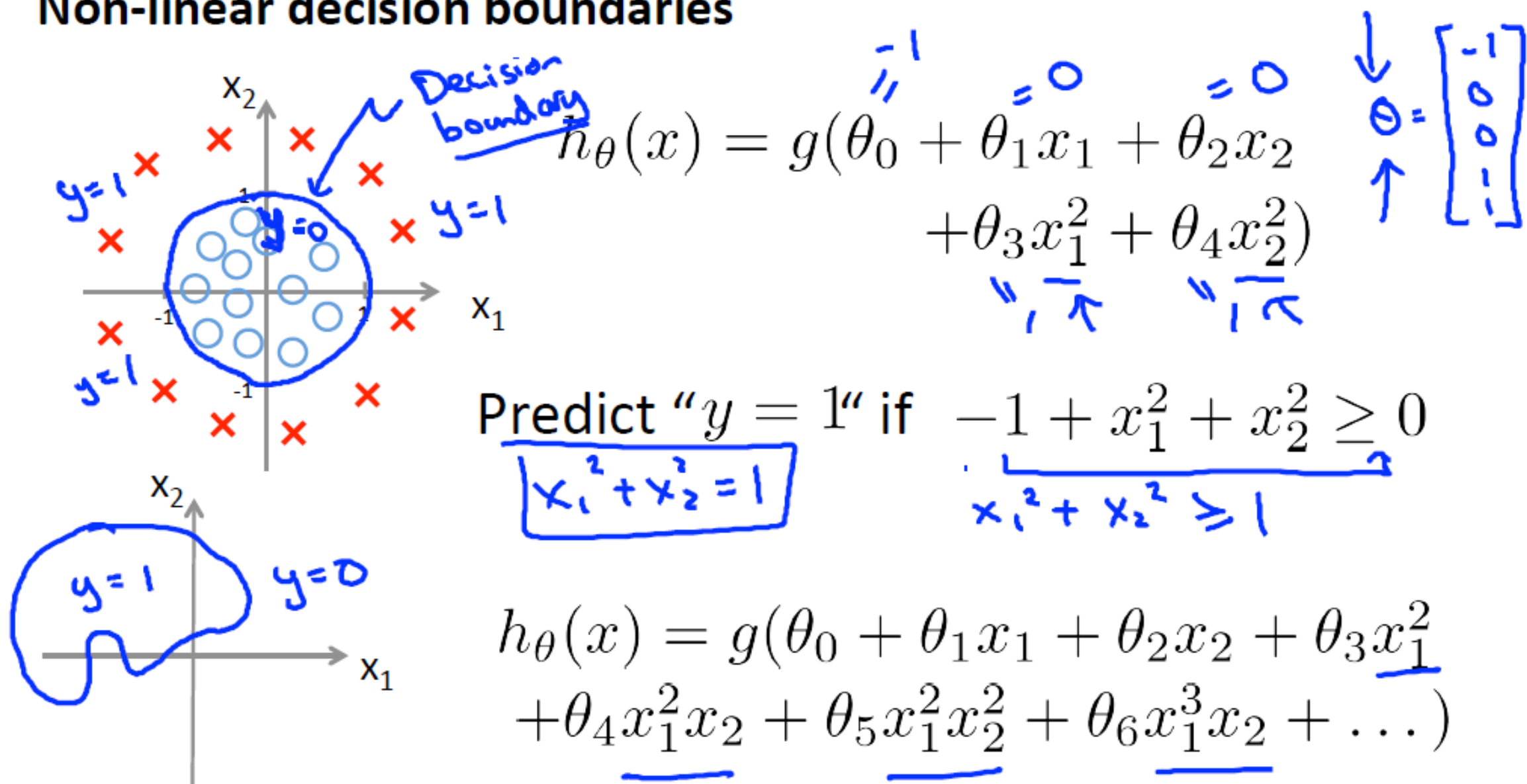
$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Predict " $y = 1$ " if $-1 + x_1^2 + x_2^2 \geq 0$

$\boxed{x_1^2 + x_2^2 = 1}$

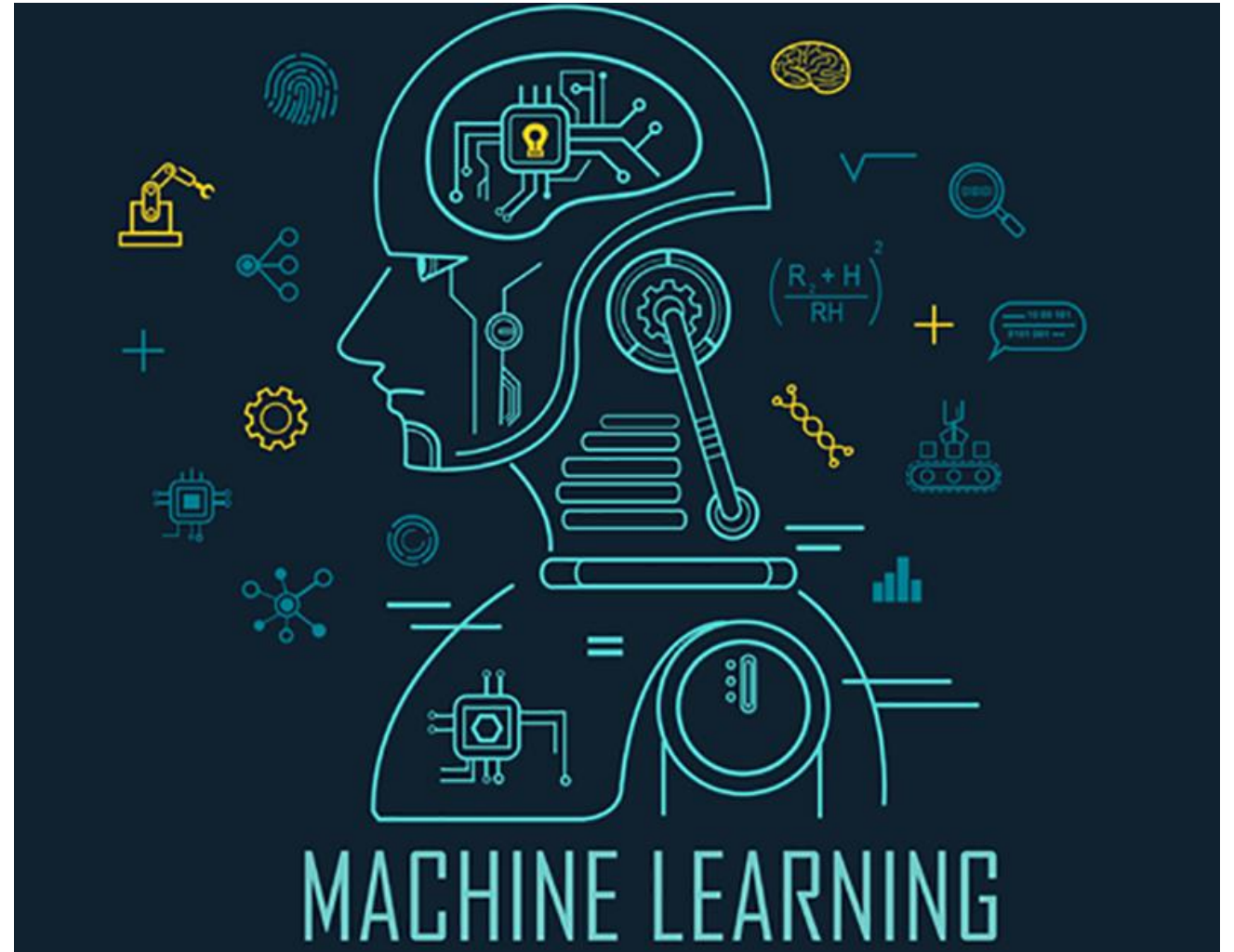
$\underbrace{-1 + x_1^2 + x_2^2}_{x_1^2 + x_2^2 \geq 1} \geq 0$

Non-linear decision boundaries



Logistic Regression

Cost Function



Training
set:

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \mathbb{R}^{n+1}$$

$$\underline{x_0 = 1}, \underline{y \in \{0, 1\}}$$

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

How to choose parameters θ ?

Cost function

→ Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$\text{cost}(h_{\theta}(x^{(i)}), y)$

→ $\text{Cost}(\underbrace{h_{\theta}(x)}_{\text{prediction}}, \underbrace{y}_{\text{target}}) = \frac{1}{2} (h_{\theta}(x) - y)^2$

Cost function

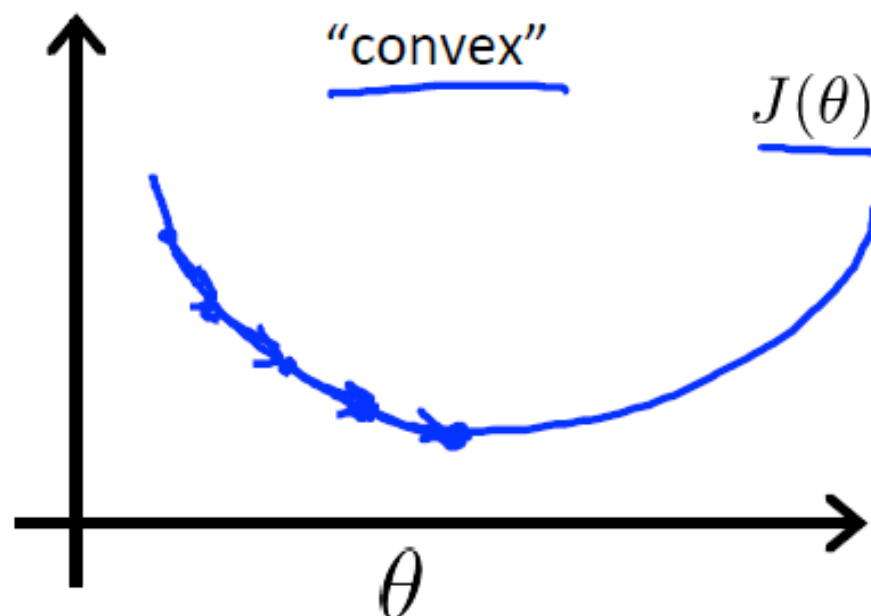
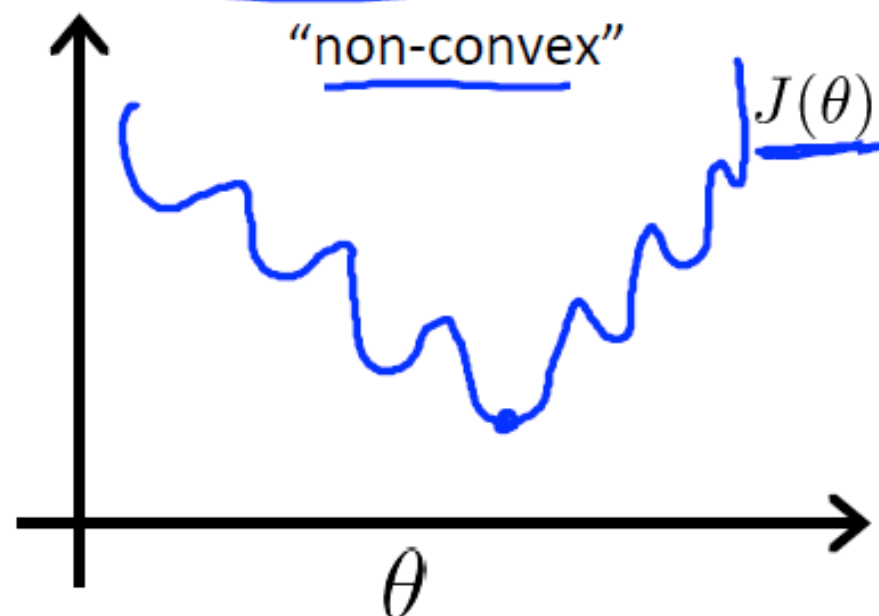
→ ~~Linear~~ regression:
logistic

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

→ $\text{cost}(h_{\theta}(x^{(i)}), y)$

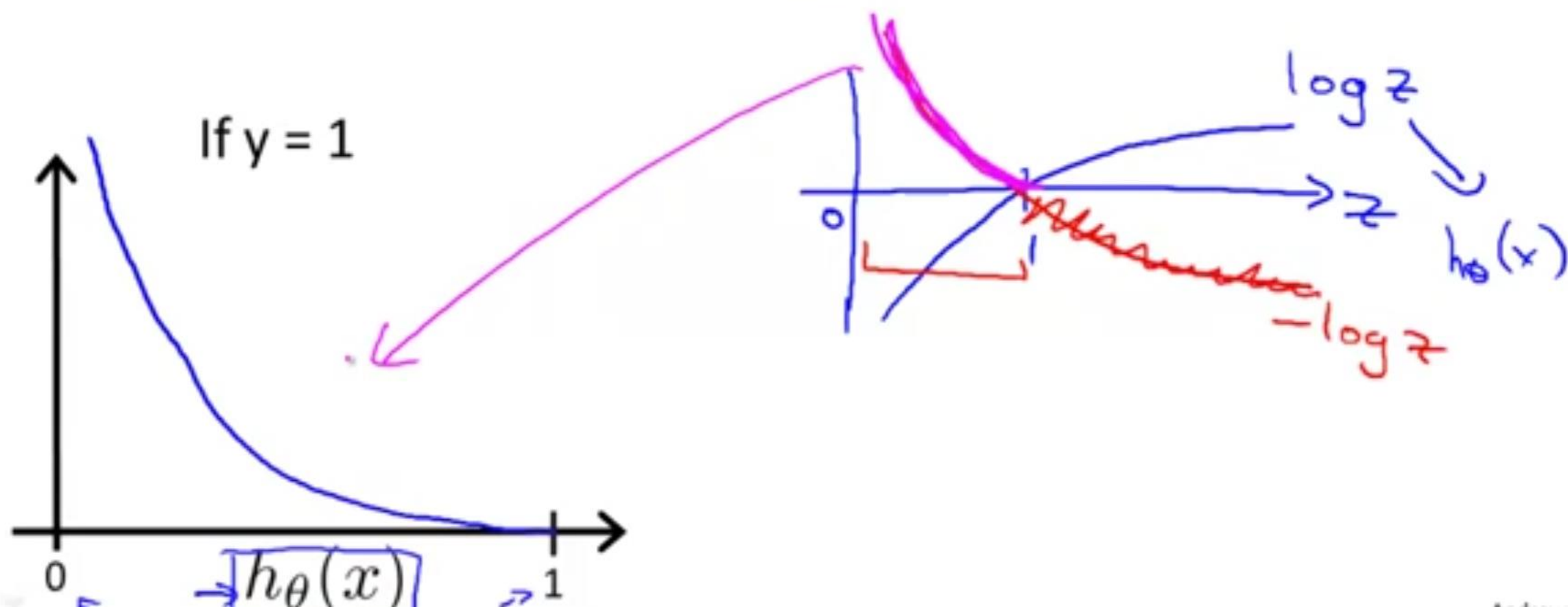
$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

← $\frac{1}{1 + e^{-\theta^T x}}$



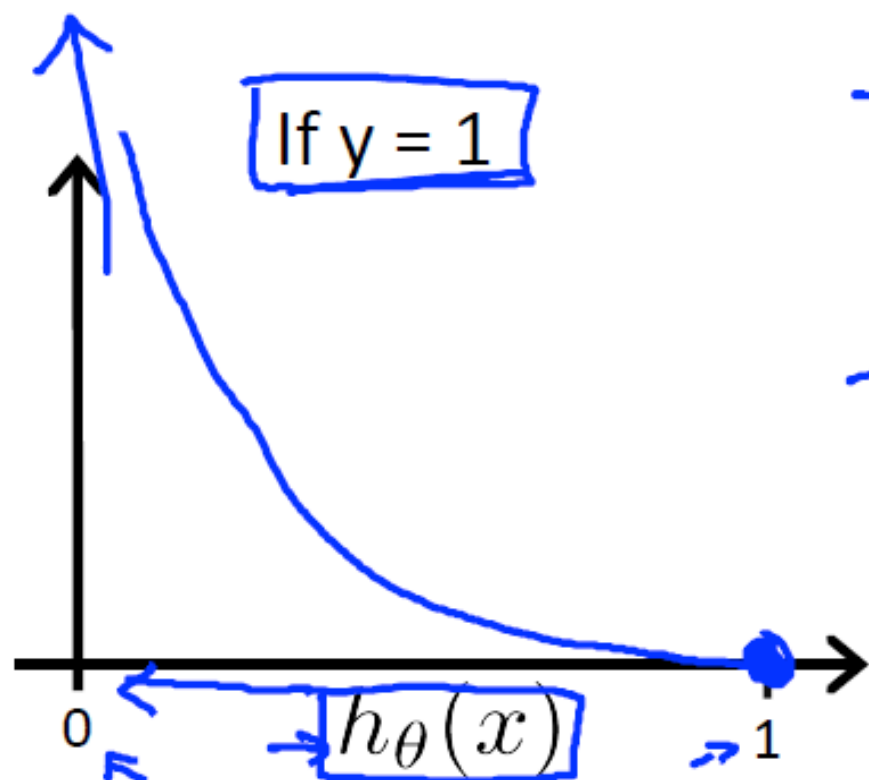
Logistic regression cost function

$$\text{Cost}(\underbrace{h_{\theta}(x)}_{\uparrow}, y) = \begin{cases} \boxed{-\log(h_{\theta}(x))} & \text{if } y = 1 \\ \underline{-\log(1 - h_{\theta}(x))} & \text{if } y = 0 \end{cases}$$



Logistic regression cost function

$$\text{Cost}(\underbrace{h_{\theta}(x)}_{\uparrow}, y) = \begin{cases} \boxed{-\log(h_{\theta}(x))} & \text{if } y = 1 \\ \underline{-\log(1 - h_{\theta}(x))} & \text{if } y = 0 \end{cases}$$

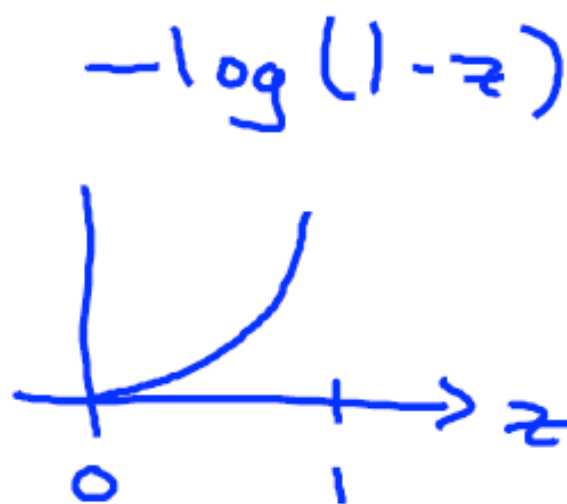
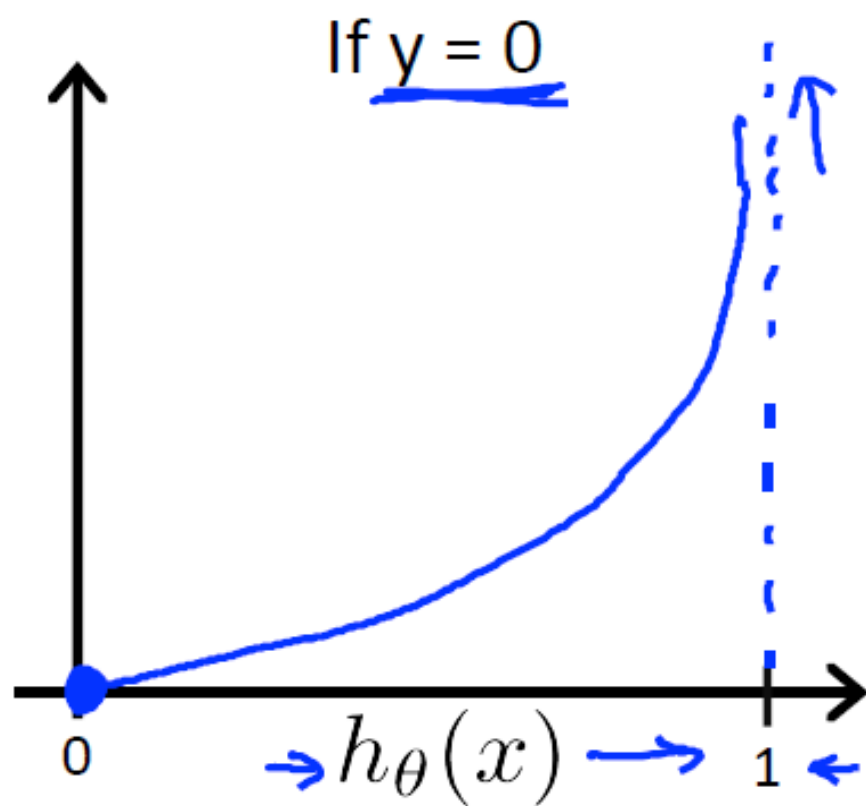


→ Cost = 0 if $y = 1, h_{\theta}(x) = 1$
But as $\frac{h_{\theta}(x) \rightarrow 0}{\text{Cost} \rightarrow \infty}$

→ Captures intuition that if $h_{\theta}(x) = 0$,
(predict $\underline{P(y = 1|x; \theta) = 0}$), but $\underline{y = 1}$,
we'll penalize learning algorithm by a very
large cost.

Logistic regression cost function

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



$-\log t$ very high when $t \rightarrow 0$ ✓

Training a Logistic Regression Model?

- ✓ We need values of parameters in theta
- We need high values of probabilities near 1 for positive instances
- We also want low values of probabilities near 0 for negative instances



$$C(x) = \begin{cases} -\log(\hat{p}) & \text{if } y=1 \\ -\log(1-\hat{p}) & \text{if } y=0 \end{cases}$$

Cost for single training instance =

x_1	x_2	x_3	x_4	x_5	L	predicted	Cost
3	4	1	4	3	1	0.8	0
1	1	0	1	0	0	0.99	

$$C(x) = -[y \log(\hat{p}) + (1-y) \log(1-\hat{p})]$$

Simplified Cost Function and Gradient Descent



Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$\rightarrow \text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

Note: $y = 0$ or 1 always

$$\rightarrow \text{Cost}(h_{\theta}(x), y) = -\underbrace{y}_{=0} \log(h_{\theta}(x)) - \underbrace{(1-y)}_{=1} \log(1 - h_{\theta}(x))$$

If $y=1$: $\text{Cost}(h_{\theta}(x), y) = -\log h_{\theta}(x) \leftarrow$

If $y=0$: $\text{Cost}(h_{\theta}(x), y) = \underline{-\log(1 - h_{\theta}(x))}$

Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$
$$= \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

To fit parameters θ :

$$\min_{\theta} J(\theta) \quad \text{Get } \underline{\Theta}$$

To make a prediction given new x :

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

$$\underline{p(y=1 | x; \theta)}$$

Gradient Descent

$$\rightarrow J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Want $\min_{\theta} J(\theta)$:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

(simultaneously update all θ_j)

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Gradient Descent

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Want $\min_{\theta} J(\theta)$:

Repeat {

$\rightarrow \theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

} (simultaneously update all θ_j)

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$ for $i = 0 \text{ to } n$

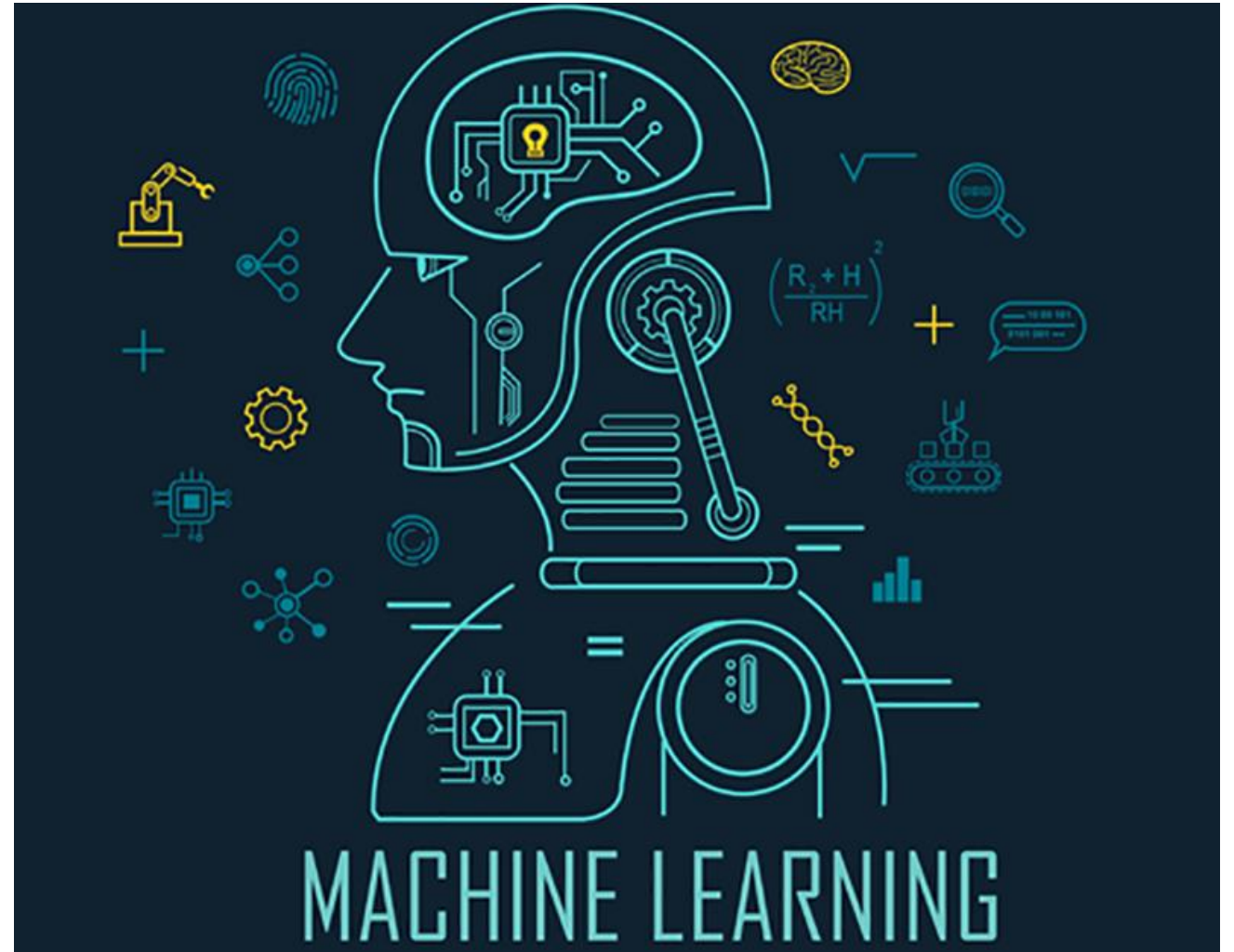
$h_{\theta}(x) = \theta^T x$

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

Algorithm looks identical to linear regression!

Logistic Regression

Advanced Optimization



Optimization algorithm

Cost function $J(\theta)$. Want $\min_{\theta} J(\theta)$.

Given θ , we have code that can compute

$\rightarrow - J(\theta)$
 $\rightarrow - \frac{\partial}{\partial \theta_j} J(\theta)$ (for $j = 0, 1, \dots, n$)

Gradient descent:

Repeat {

$\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$
}

Optimization algorithm

Given θ , we have code that can compute

$$\begin{aligned} & - J(\theta) \\ & - \frac{\partial}{\partial \theta_j} J(\theta) \end{aligned} \quad (\text{for } j = 0, 1, \dots, n)$$

Optimization algorithms:

- - Gradient descent
- Conjugate gradient
- BFGS (Broyden-Fletcher-Goldfarb-Shanno)
- L-BFGS

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

Disadvantages:

- More complex

Example:

- $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

Example: $\min_{\theta} J(\theta)$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_1=5, \theta_2=5.$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

```
function [jVal, gradient]
    = costFunction(theta)
    jVal = (theta(1)-5)^2 + ...
          (theta(2)-5)^2;
    gradient = zeros(2,1);
    gradient(1) = 2*(theta(1)-5);
    gradient(2) = 2*(theta(2)-5);
```

```
options = optimset('GradObj', 'on', 'MaxIter', '100');
```

```
initialTheta = zeros(2,1);
```

```
[optTheta, functionVal, exitFlag] ...
    = fminunc(@costFunction, initialTheta, options);
```

↑

↑

$\theta \in \mathbb{R}^d$ $d \geq 2$.

$$\underline{\text{theta}} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

Annotations: θ_0 is labeled theta(1) with an arrow pointing to it. θ_1 is labeled theta(2). θ_n is labeled theta(n+1).

```
function [jVal gradient] = costFunction(theta)
```

```
    jVal = [code to compute  $J(\theta)$ ];
```

```
    gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];
```

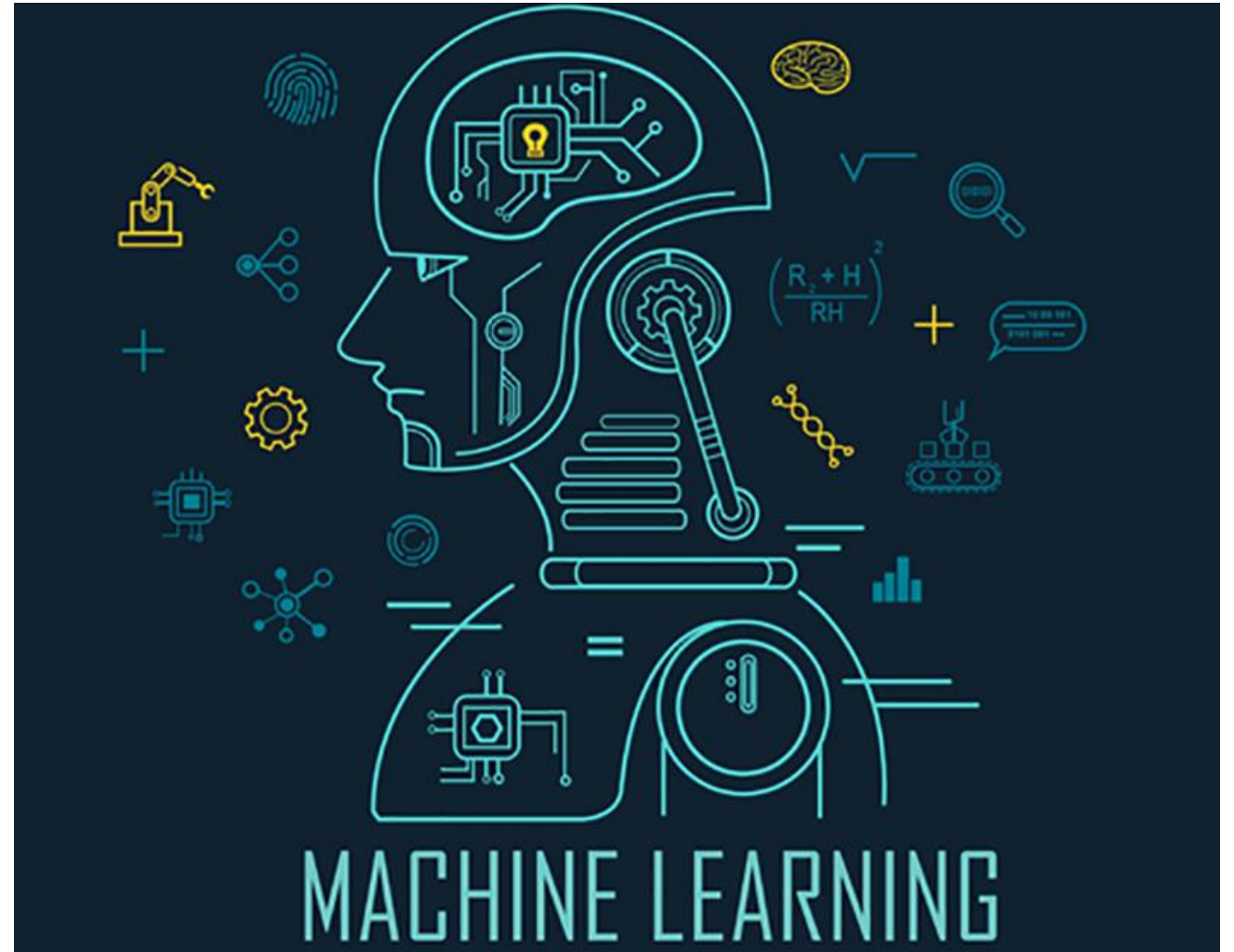
```
    gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];
```

```
    :
```

```
    gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ] ;
```

Logistic Regression

**Multiple-Class
Classification: One VS All**



Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby


$y=1$ $y=2$ $y=3$ $y=4$

Medical diagrams: Not ill, Cold, Flu

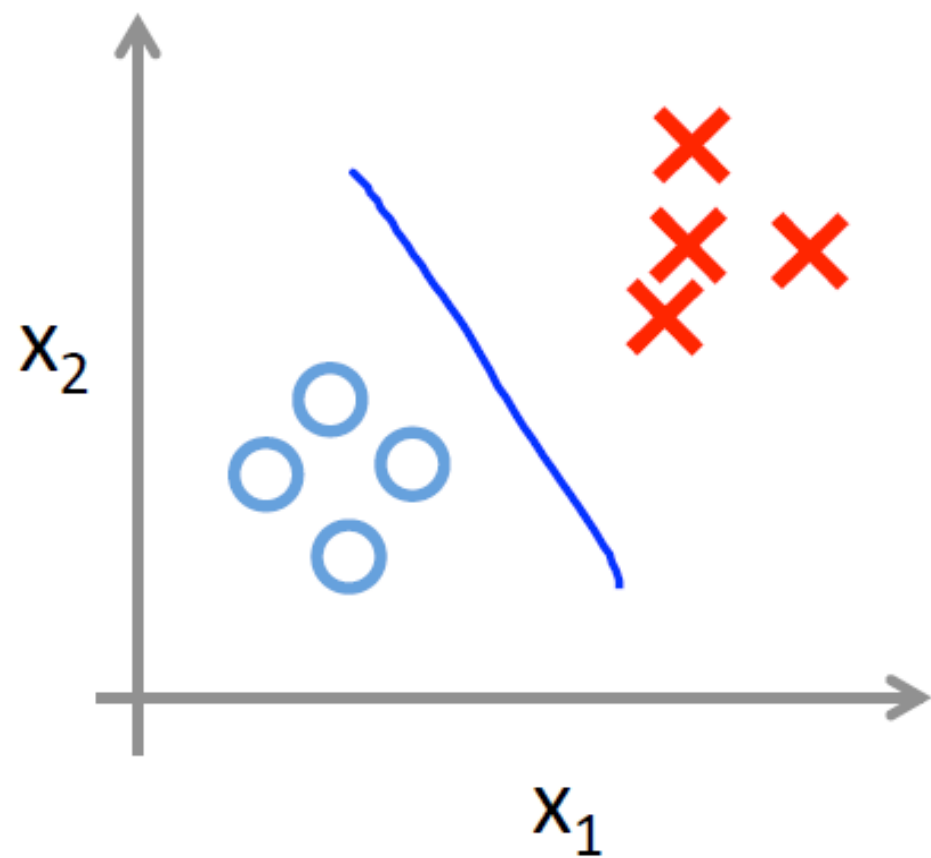
$y=1$ 2 3

Weather: Sunny, Cloudy, Rain, Snow

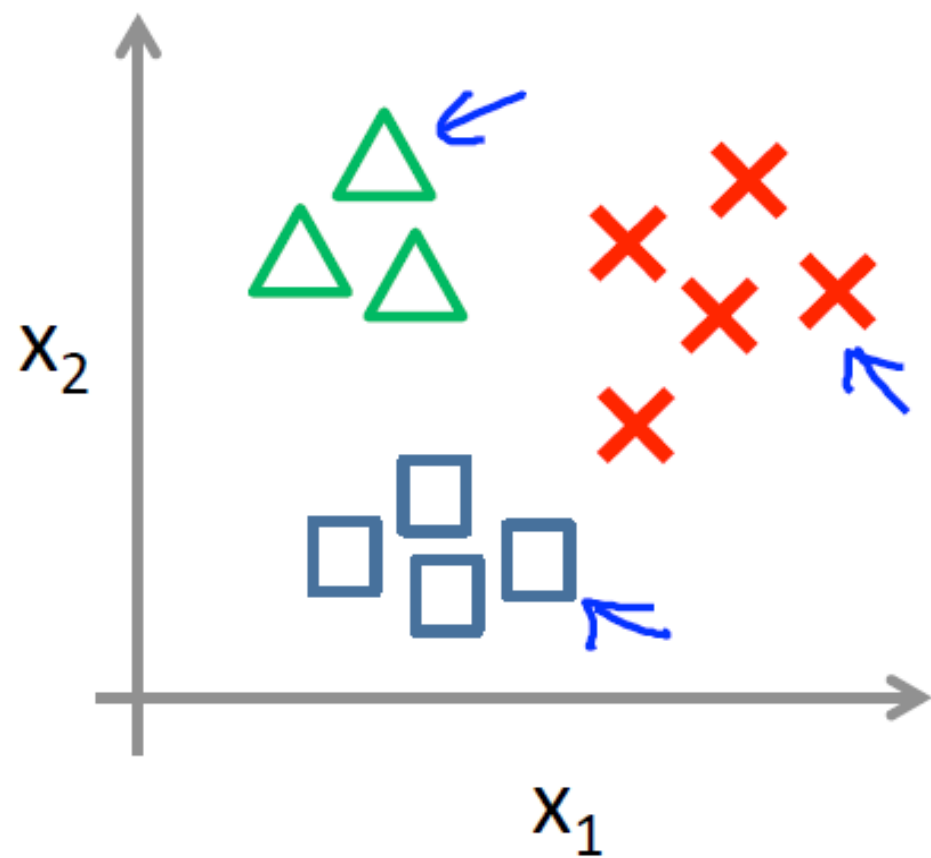
$y=1$ 2 3 4 \leftarrow



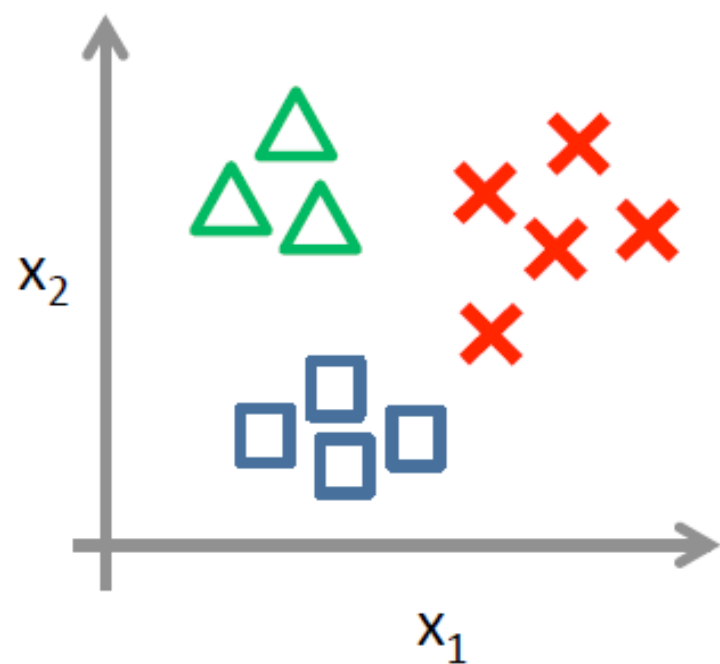
Binary classification:





Multi-class classification:




One-vs-all (one-vs-rest):

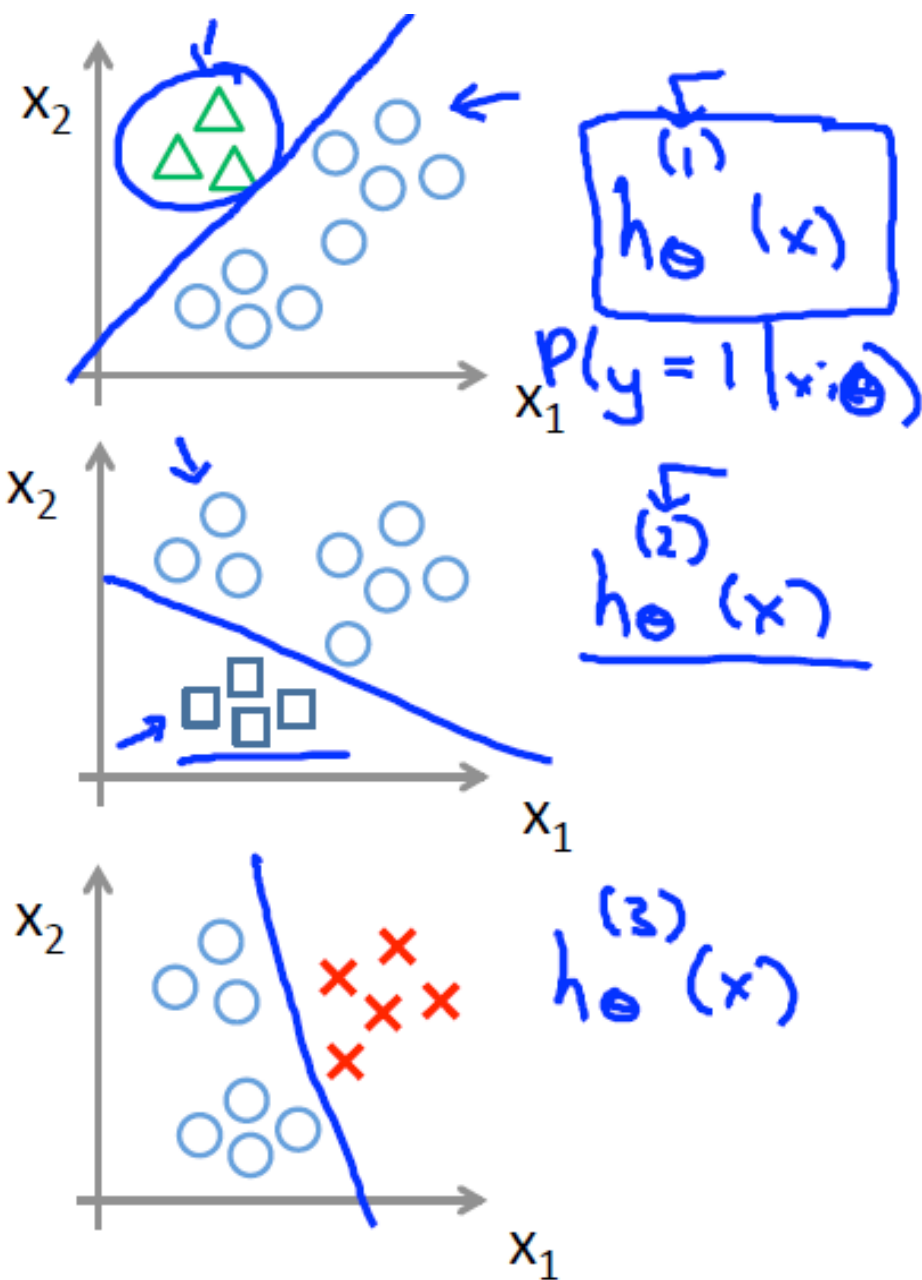


Class 1:  \leftarrow

Class 2:  \leftarrow

Class 3:  \leftarrow

$$h_{\theta}^{(i)}(x) = P(y = i | x; \theta) \quad (i = 1, 2, 3)$$



One-vs-all

Train a logistic regression classifier $\underline{h_{\theta}^{(i)}(x)}$ for each class \underline{i} to predict the probability that $\underline{y = i}$.

On a new input \underline{x} , to make a prediction, pick the class i that maximizes

$$\max_{\underline{i}} \underline{h_{\theta}^{(i)}(x)}$$
