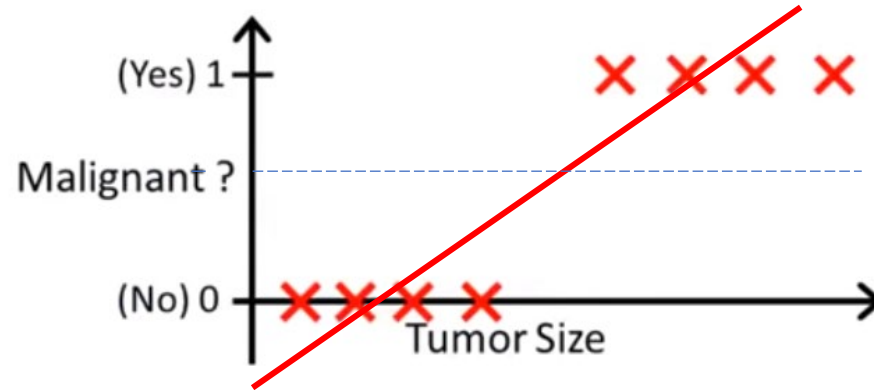


# Logistic Regression

(supervised learning for classification)

# Classification problem



Use linear regression,  $h_{\theta}(x) = x^T \theta$  ?

Set a threshold for classification (at 0.5):

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

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

## Inconvenience of linear regression:

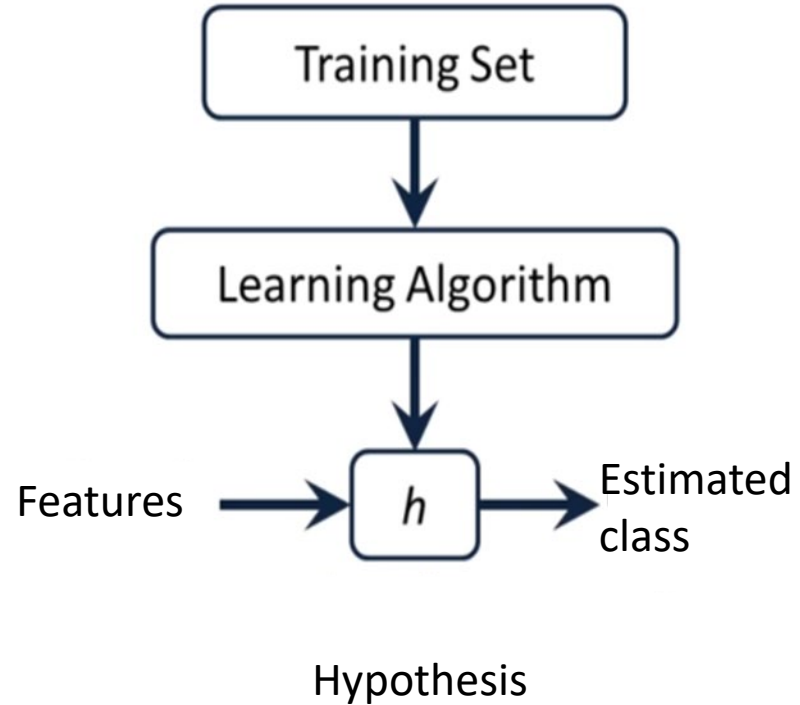
Classification:  $y = 0$  or  $1$

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

## Introduce:

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

## Supervised learning – classification problem



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

$x$ : features

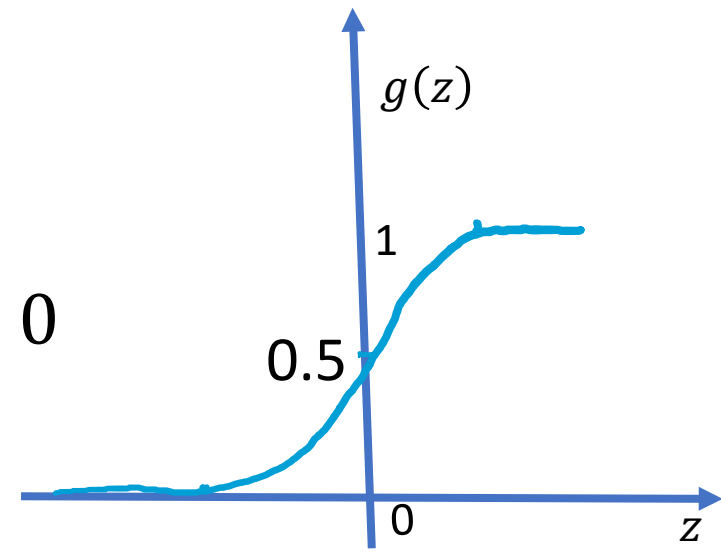
$\theta$ : parameter of model

## Logistic regression model:

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

estimated probability that  $y = 1$  or 0

given input  $x$ ,  $h_{\theta}(x) = P(y|x, \theta)$



Sigmoid/logistic function

Realized by a learning model (logistic regression model)

$$h_{\theta}(x) = g(x^T \theta) \quad \text{where} \quad g(z) = \frac{1}{1 + e^{-z}}$$

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

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

- Probability of predicting  $y$ , given  $x$ , parameterized by  $\theta$

Predict  $y = 1$  if  $h_{\theta}(x) \geq 0.5$

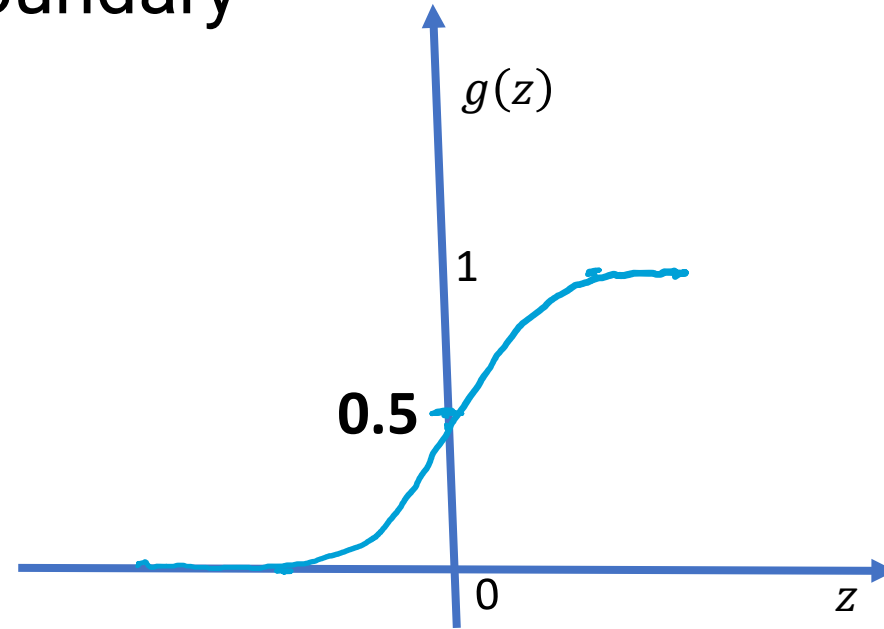
Predict  $y = 0$  if  $h_{\theta}(x) < 0.5$

- For example,  $h_{\theta}(x) = 0.85$ , tell patient that 85% of chance of tumor being cancerous

# Decision Boundary

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

(Let  $z = x^T \theta$ )

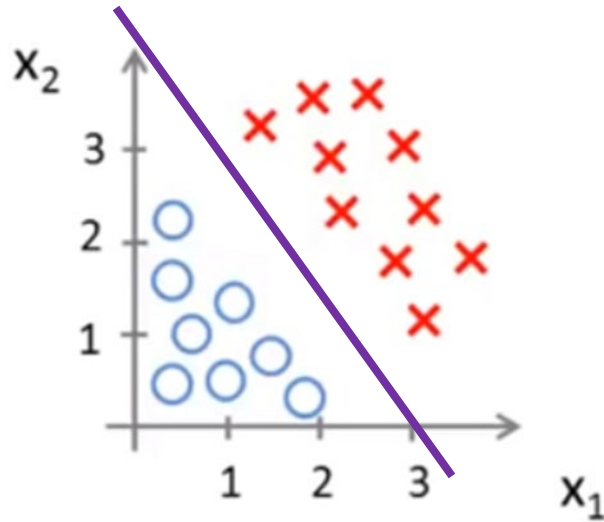


Predict  $y = 1$  if  $h_{\theta}(x) \geq 0.5$  or  $z \geq 0$

Predict  $y = 0$  if  $h_{\theta}(x) < 0.5$  or  $z < 0$

Decision boundary:  $x^T \theta = 0$

## Decision Boundary



$$h_{\theta}(x) = g(-3 + x_1 + x_2)$$

Predict  $y = 1$  if  $h_{\theta}(x) \geq 0.5$  or  $(-3 + x_1 + x_2) \geq 0$

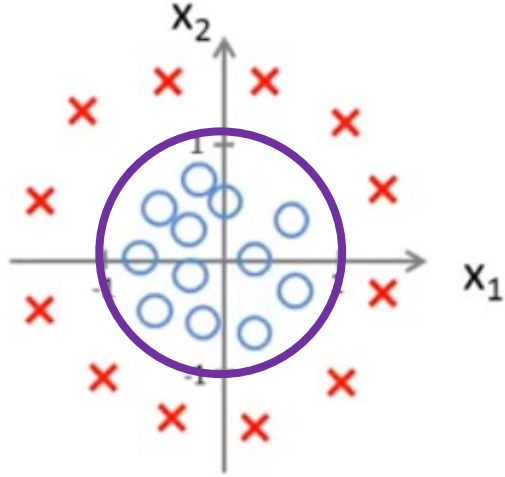
Predict  $y = 0$  if  $h_{\theta}(x) < 0.5$  or  $(-3 + x_1 + x_2) < 0$

**Decision boundary:**

$$x_1 + x_2 = 3$$



## Nonlinear decision boundaries



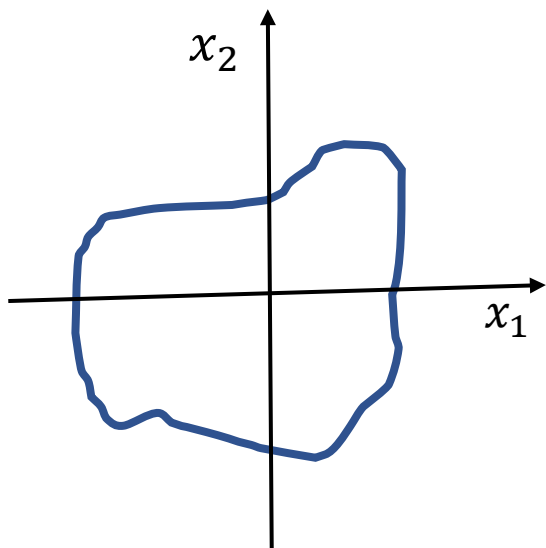
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = [-1, 0, 0, 1, 1]$$

Decision boundary:

$$x_1^2 + x_2^2 = 1$$

## More complicated nonlinear decision boundary



$$h_{\theta}(x) = g \left( \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots \right)$$

**Decision boundary:**

$$\begin{aligned} &\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 \\ &+ \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \\ &\dots = 0 \end{aligned}$$

Given training set for supervised learning:

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

where

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0,1\}$$

In a logistic regression model,

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

How to choose parameter  $\theta$ ?

Cost function consideration:

If we use the cost formulation as in linear regression,

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

Recall:

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

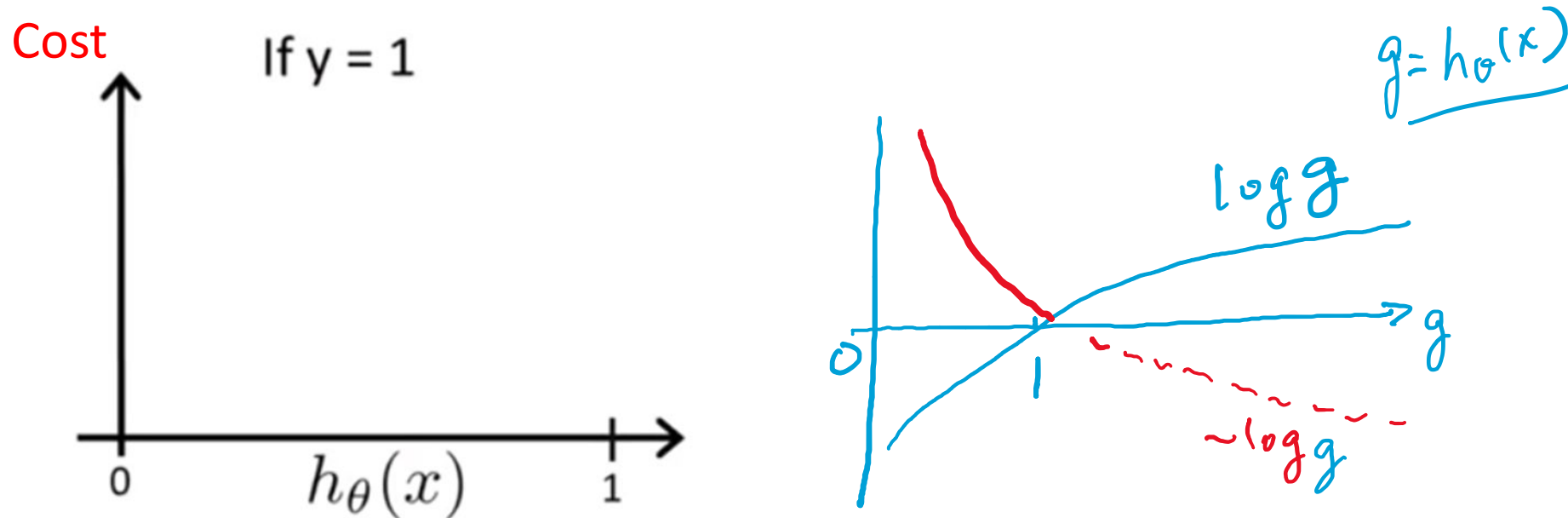
$g$  is a sigmoid/logistic function

Look at the term:

$$\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \implies \text{Nonconvex, thus local minima}$$

Propose “logistic regression cost function” as follows:

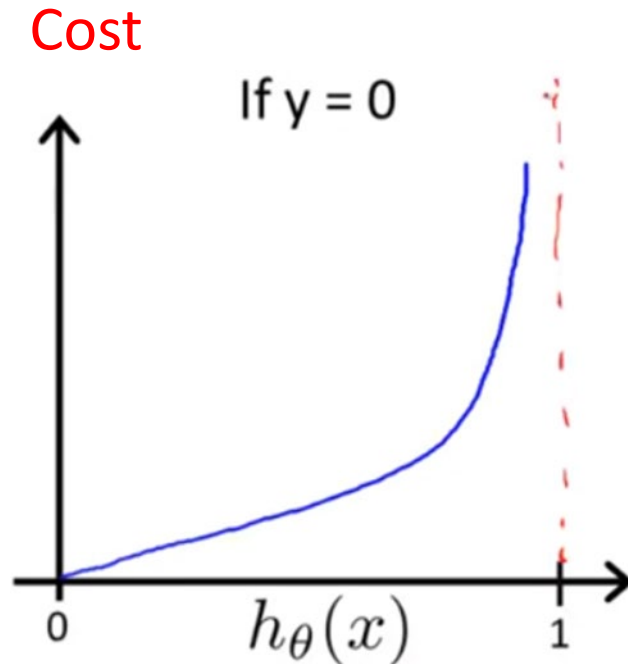
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



$\text{Cost} = 0$  if  $y = 1, h_{\theta}(x) = 1$

$\text{Cost} \rightarrow \infty$  if  $y = 1, \text{ but } h_{\theta}(x) \rightarrow 0$

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



$\text{Cost} = 0$  if  $y = 0, h_{\theta}(x) = 0$

$\text{Cost} \rightarrow \infty$  if  $y = 0$ , but  $h_{\theta}(x) \rightarrow 1$

# Logistic regression cost function

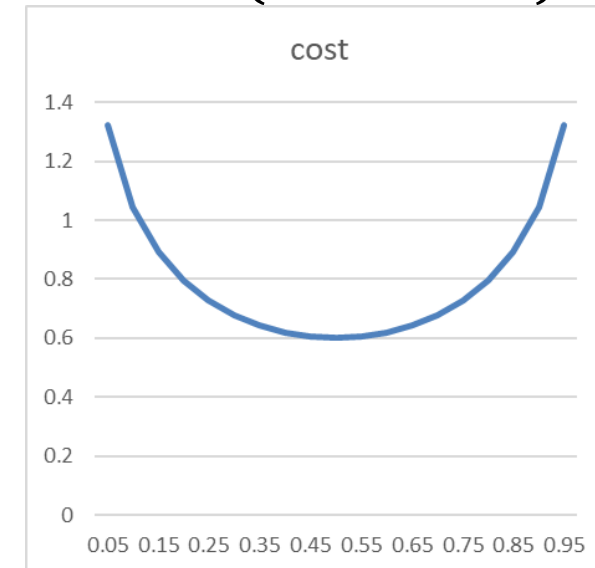
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

Note:  $y = 0$  or  $1$

**Simplification:**

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



## Logistic regression cost function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \end{aligned}$$

To fit parameters  $\theta$ :

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



To get parameters

$$\theta = (\theta_0, \theta_1, \dots, \theta_n)$$

To make a prediction given new  $x$ :

Compute output:

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



Apply gradient descent to reduce the cost measure

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

Update parameters  $\theta_0 \theta_1 \dots \theta_n$  using gradient descent until convergence

Simultaneously update  $\theta_j, j = 0, 1, \dots, n$ , according to

$$\theta_j := \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n)$$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$