



# Introduction to Data Science

## (Lecture 12)

**Dr. Mohammad Pourhomayoun**

Assistant Professor

Computer Science Department

California State University, Los Angeles



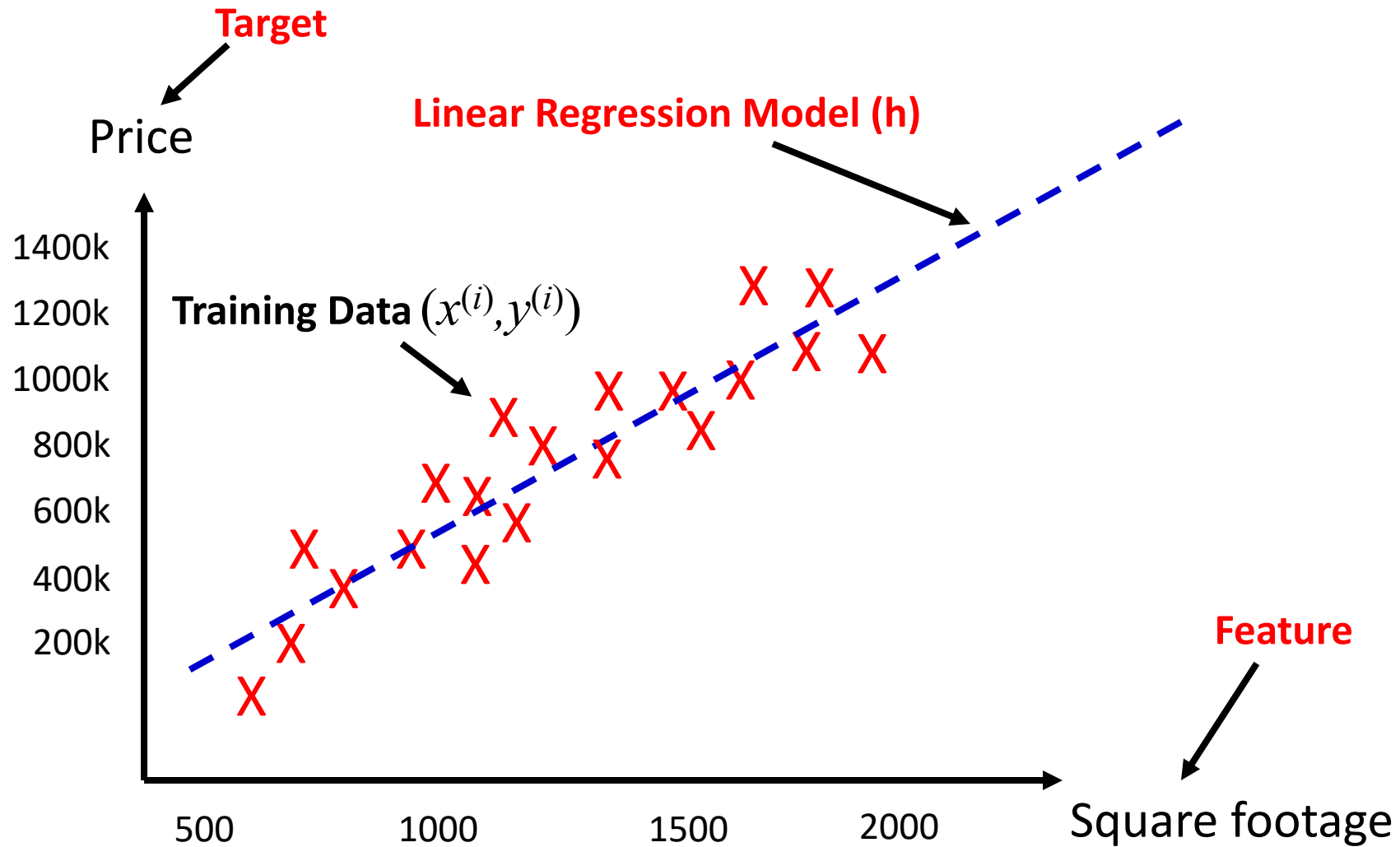


# **Logistic Regression Classifier**

# Review

- **Classification:** Predict a discrete valued output for each item.
  - Labels are discrete (categorical)
  - Labels can be binary (e.g., rainy/sunny, spam/non-spam,) or non-binary (e.g., rainy/sunny/cloudy, Setosa/Versicolor/Virginica)
- **Regression:** Predict a continuous valued output for each item.
  - Labels are continuous (numeric), e.g., stock price, housing price
  - Can define 'closeness' when comparing prediction with true values

# Regression Example: Housing Price

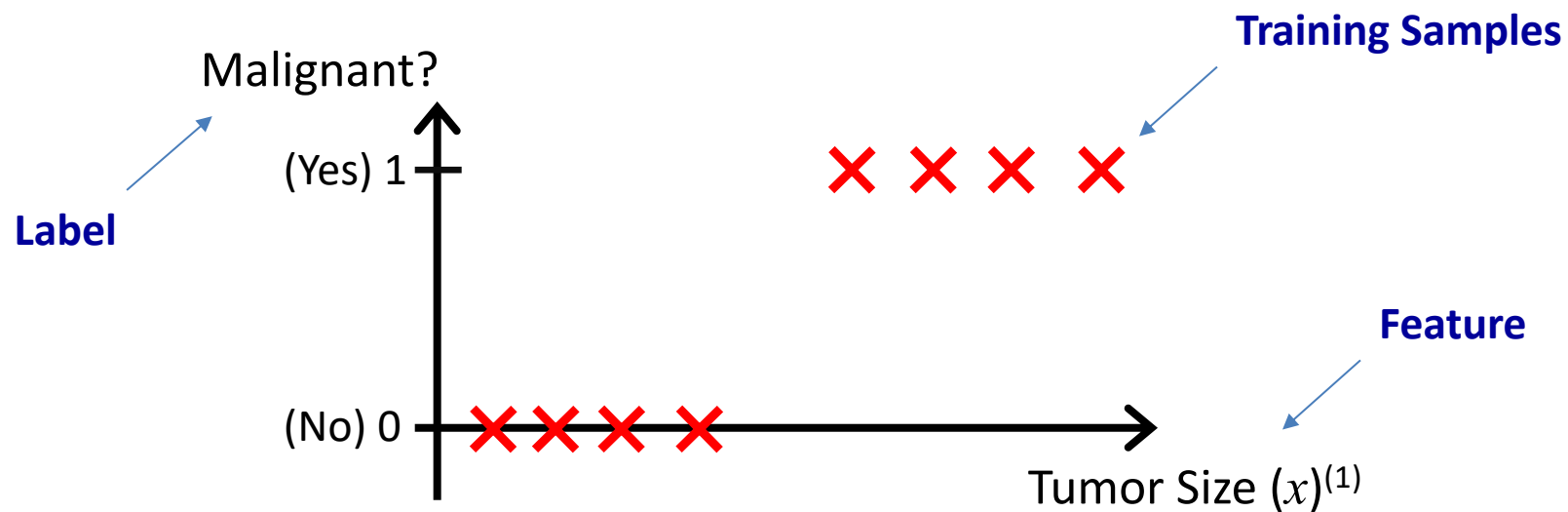


# Deriving a Classifier from a Regression Model

1. We learned how to build a **linear regression** model to predict **continuous-valued** outputs.
2. Well, in theory we can convert Categorical Labels into Numerical Labels:
  - Sunny  $\rightarrow$  0 , Rainy  $\rightarrow$  1    thus: Sunny/Rainy  $\rightarrow$  0/1
  - Not-Cancer  $\rightarrow$  0 , Cancer  $\rightarrow$  1    thus: Not-Cancer / Cancer  $\rightarrow$  0/1
3. Thus, Let's take advantage of the **linear regression** model to develop a **Classifier!**  
(We will see that discretizing the output of a linear regression is not always a good idea! So, we need to define new hypothesis model and a new cost function)
4. Since it is a classifier built based on linear regression algorithm, we call it:  
**"Logistic Regression"!**

# Logistic Regression Classifier with One Feature

**Example:** Predicting if a Tumor is Malignant or Not based on tumor size (so, we need a classifier):

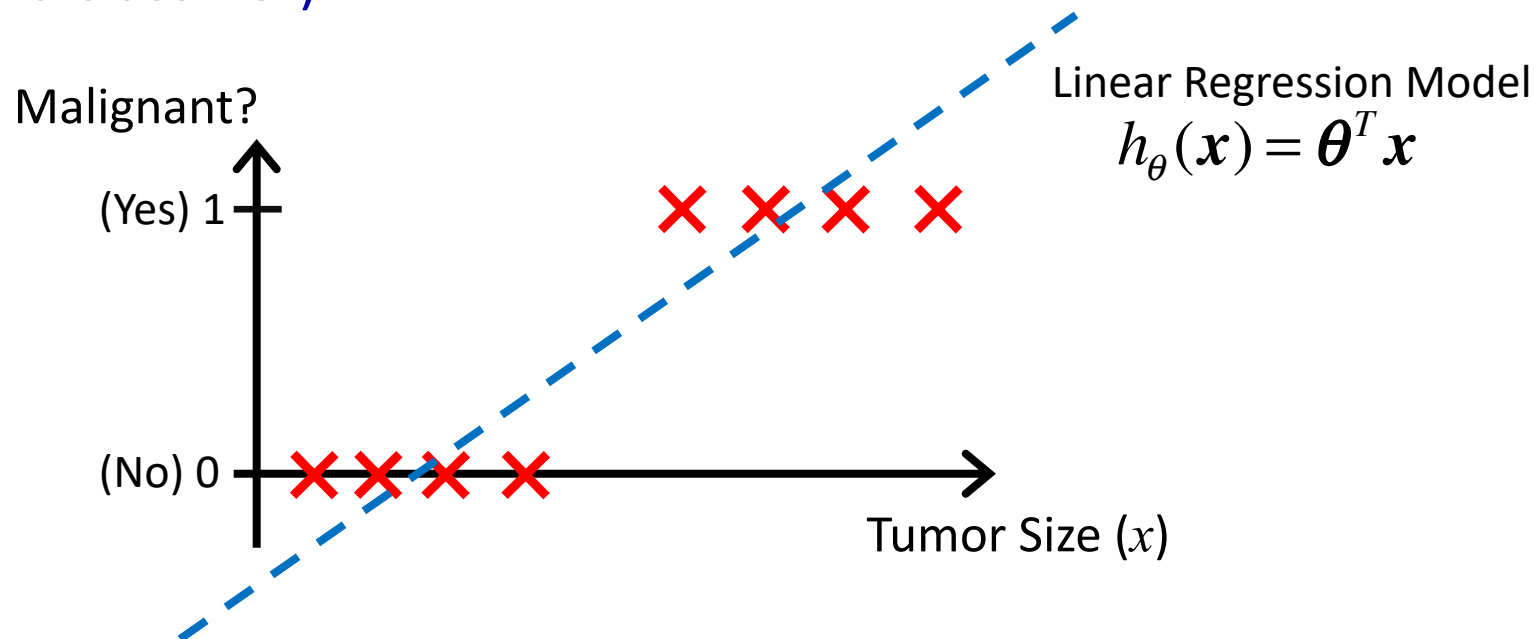


**This is a discrete-valued data. Nonetheless, let's apply our Linear Regression model on it!**

\* Example from Andrew Ng, Machine Learning, Stanford University.

# Logistic Regression Classifier with One Feature

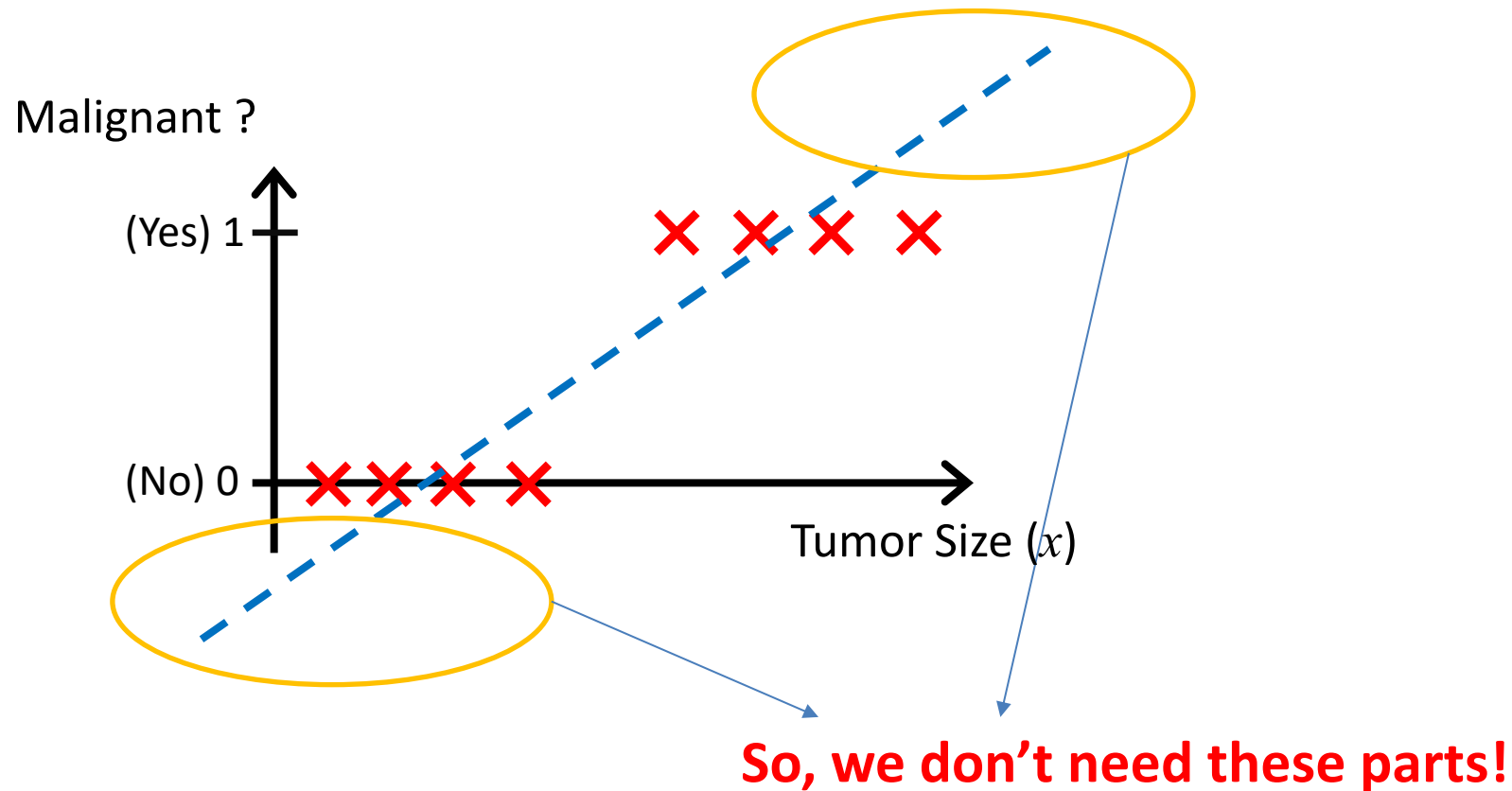
**Example:** Predicting if a Tumor is Malignant or Not based on tumor size (so, we need a classifier):



**This is a discrete-valued data. Nonetheless, let's apply our Linear Regression model on it!**

# Logistic Regression Classifier with One Feature

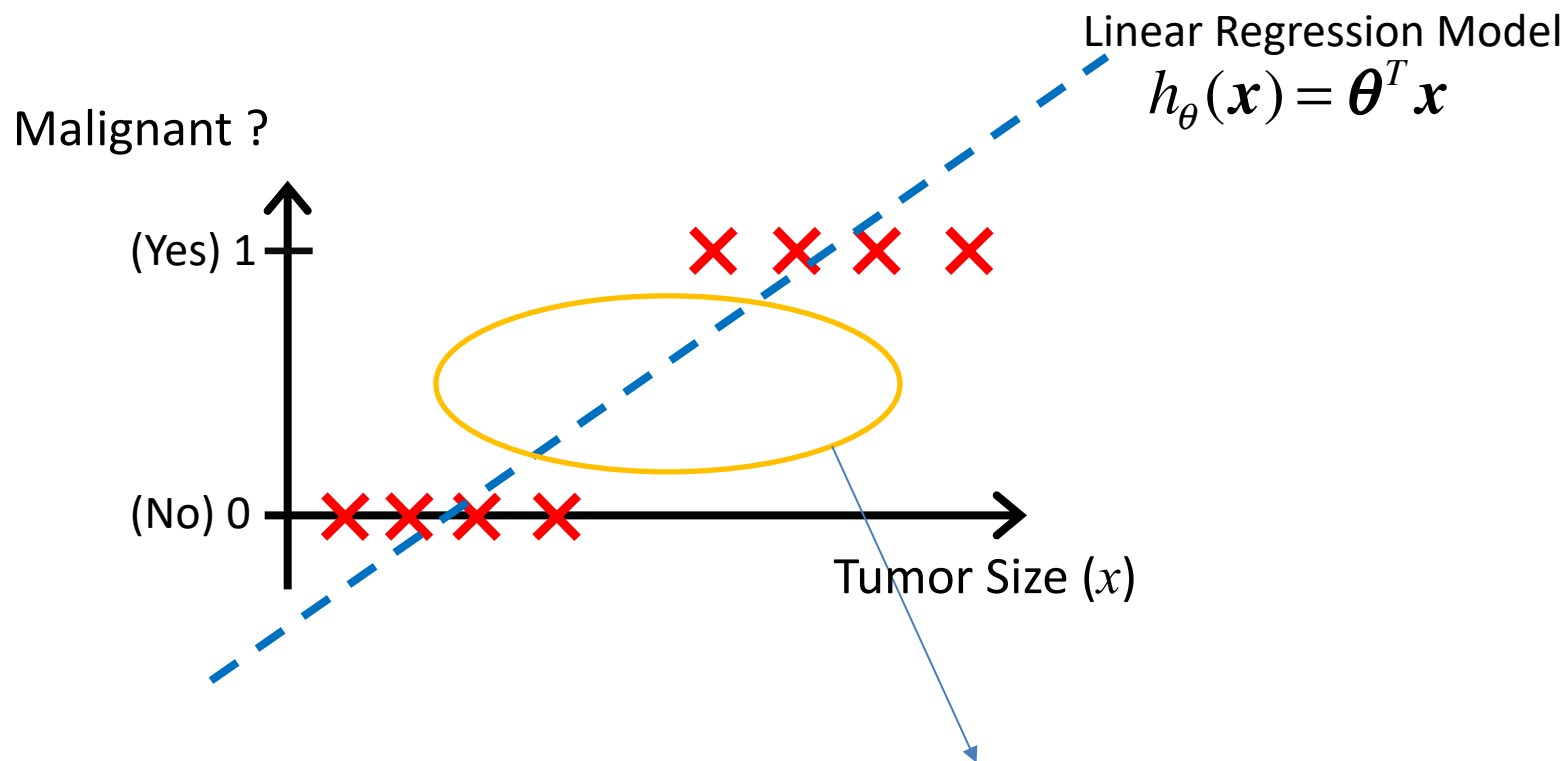
- we know that for our classifier, the output should be either 1 or 0!





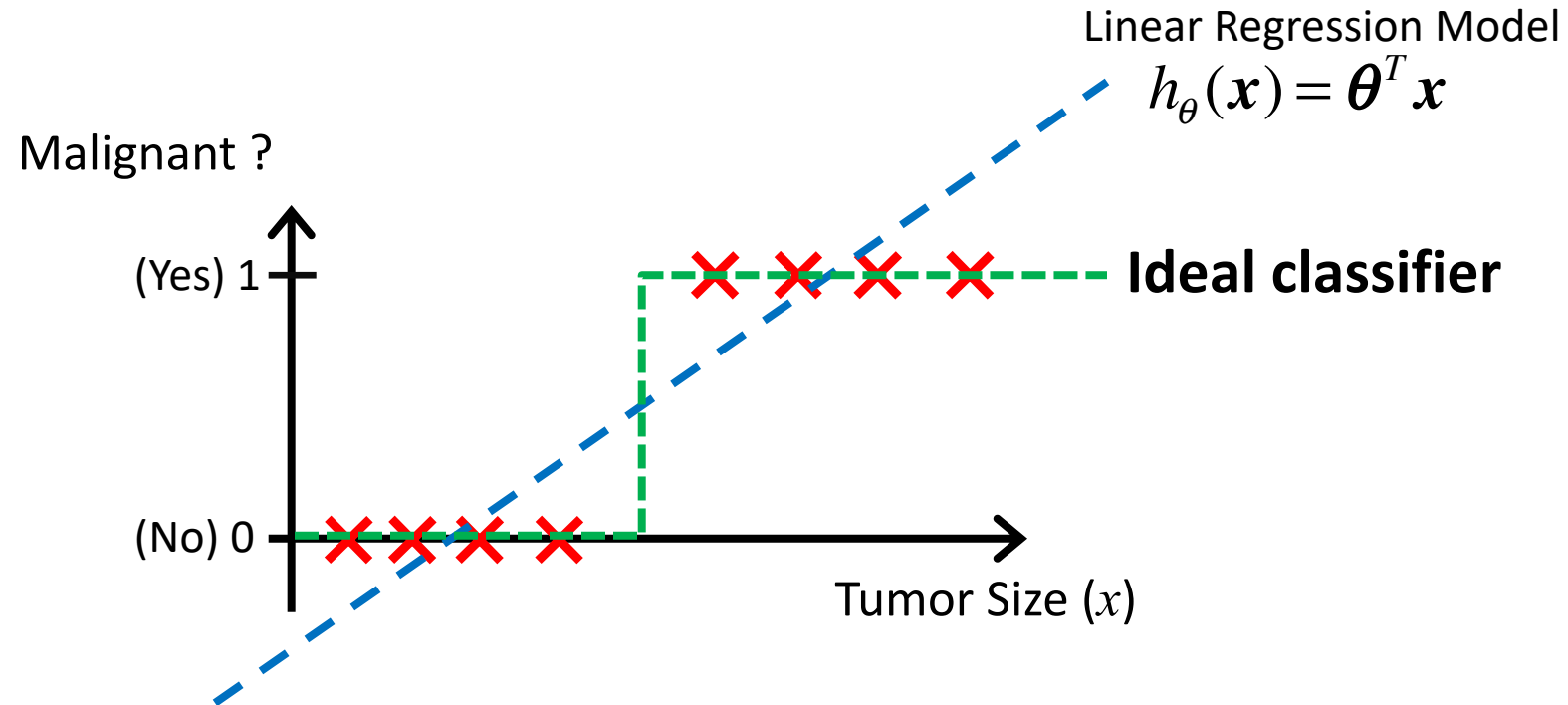
# Logistic Regression Classifier with One Feature

- we know that for our classifier, the output should be either 1 or 0!



So, We also need a sharp transition here

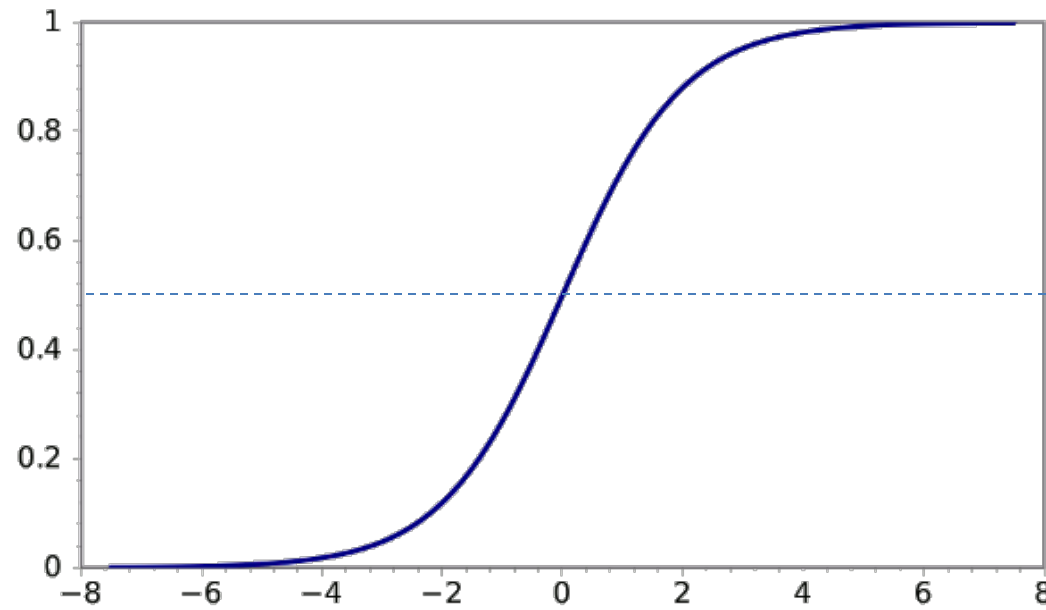
# Logistic Regression Classifier with One Feature



# Sigmoid Function

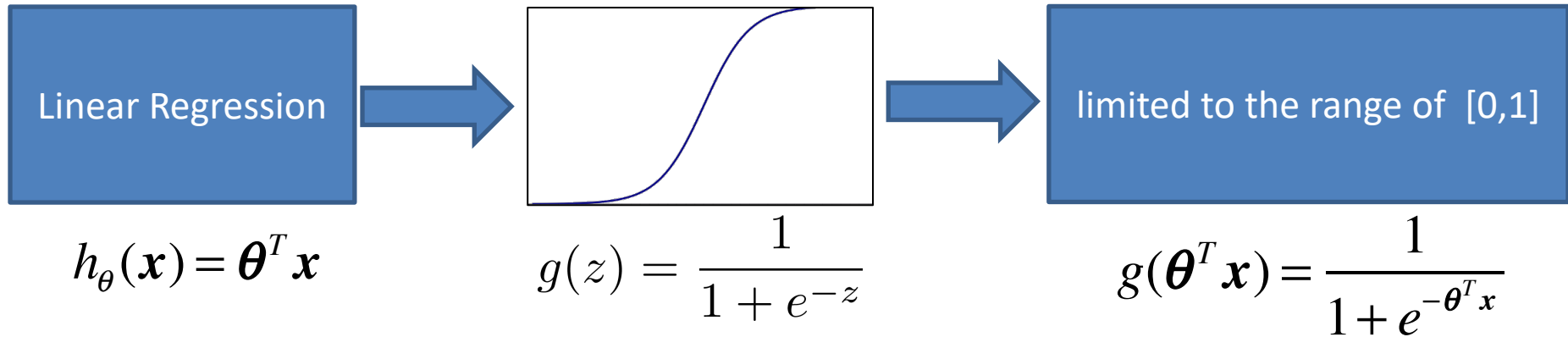
- Sigmoid Function (Logistic Function):

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



# Logistic Regression Model

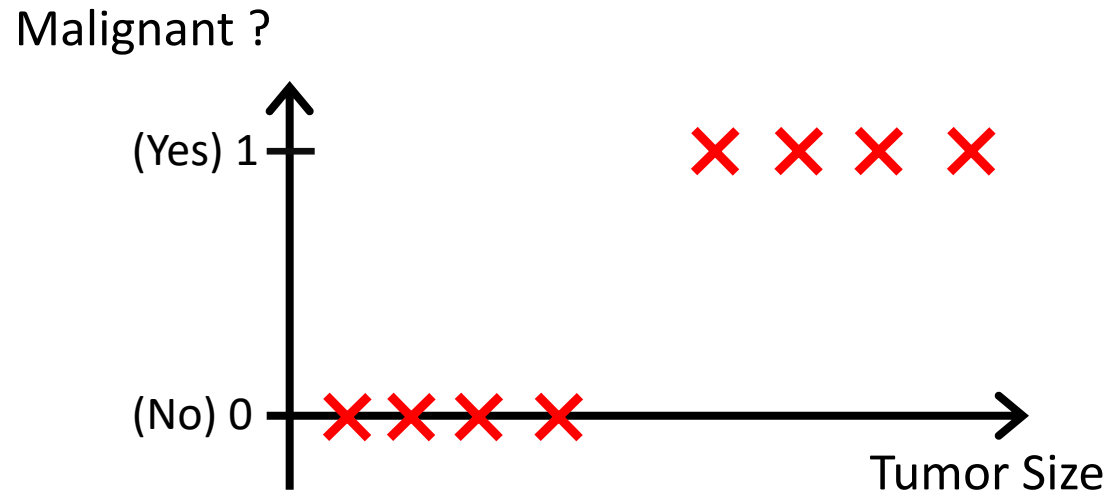
- Using Sigmoid Function (Logistic Function):



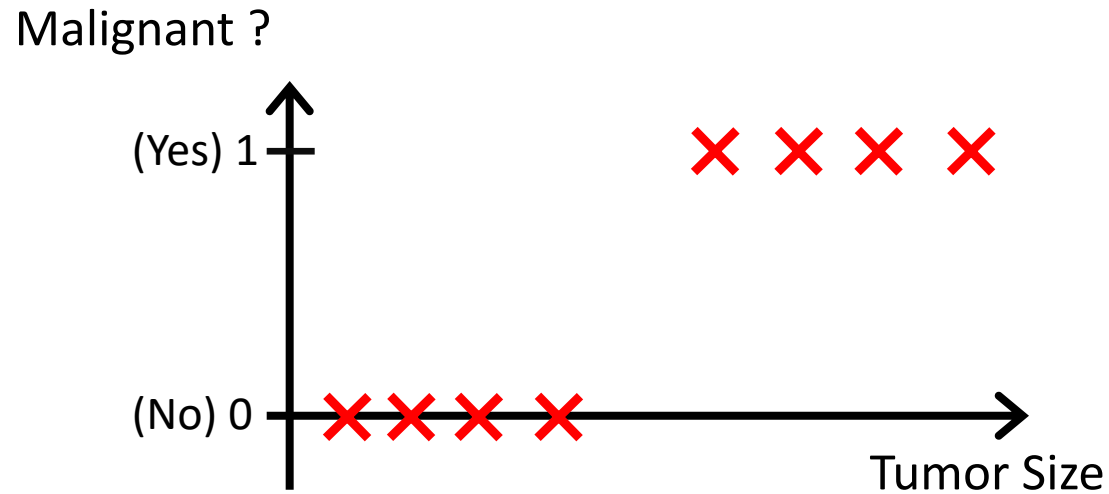
**New approach for output prediction:**

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

**So, Now the NEW  $h_{\theta}(x)$  is limited to the range of [0,1].**

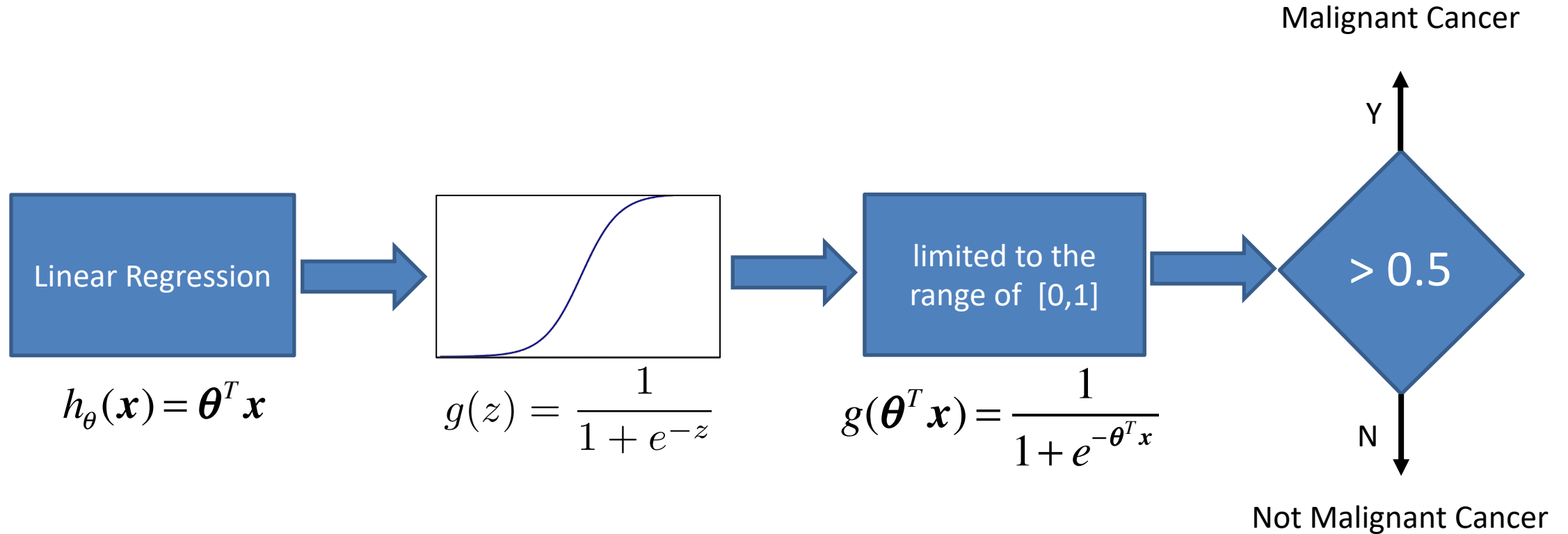


- New approach: 
$$h_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$
- After applying the Sigmoid function,  $h_{\theta}(\mathbf{x})$  will be limited in the range of  $[0,1]$ . But, still can take any value between 0 and 1!
- So, it is like representing the **Probability** of happening each label!
  - E.g. 30% chance of rain, 5% chance of malignant cancer, ...



- New approach: 
$$h_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$
- After applying the Sigmoid function,  $h_{\theta}(\mathbf{x})$  will be limited to the range of  $[0,1]$ . But, still can take any value in this range!
- **Now, to generate class output (binary output), we can compare the results with a threshold (e.g. 0.5) to discretize the output:**
$$\begin{cases} \text{predict "y = 1" if } h_{\theta}(\mathbf{x}) \geq 0.5 \\ \text{predict "y = 0" if } h_{\theta}(\mathbf{x}) \leq 0.5 \end{cases}$$

# After limiting the range and discretizing



# How to Select the Parameters?

Training set (m training samples):

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

Feature Vector:  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$   $x_0 = 1, y \in \{0, 1\}$

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

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

**How to choose parameters  $\theta$ ?**



# Gradient Descent for Logistic Regression

$$h_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

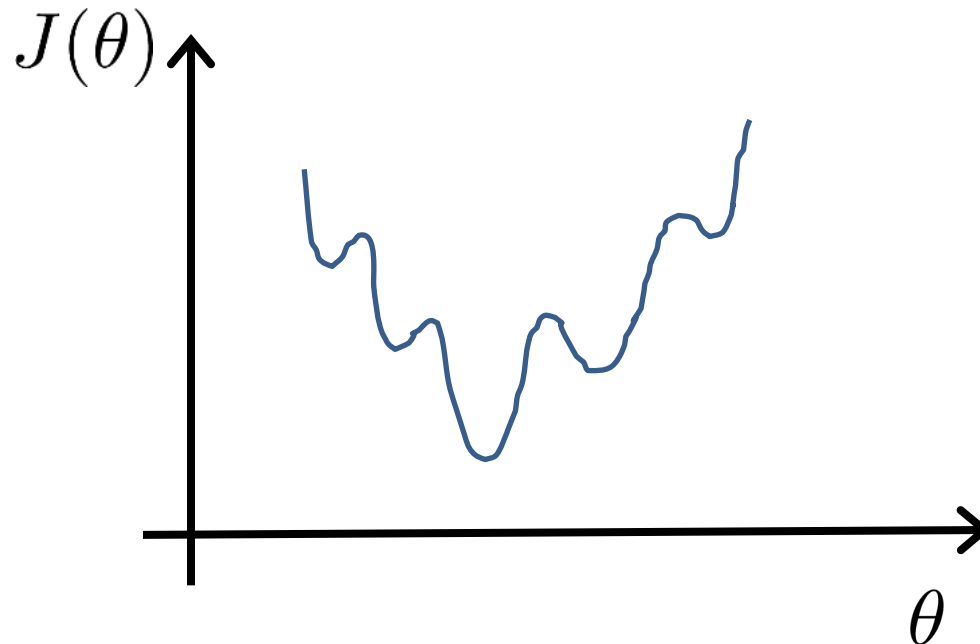
- Previous definition of Cost function for linear regression:

$$J(\boldsymbol{\theta}) = J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

# Logistic Regression Cost Function

- Previous definition of Cost function for linear regression:

$$J(\boldsymbol{\theta}) = J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$



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

**This Cost Function is not Convex for Logistic Regression (It has many local minimums)!**

- Since the Cost Function with this definition **is not convex** for logistic regression, there is **no guarantee** for gradient descent to find the **global minimum** for error. Thus, we take advantage of a “**log**” function to define an alternative **convex** cost function (The mathematical details why *log* is a good option is beyond the scope of the class! But, an intuitive idea is available in next page!).

- **New Cost function:**

$$\begin{aligned} J(\boldsymbol{\theta}) &= J(\theta_0, \theta_1, \dots, \theta_n) \\ &= -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] \end{aligned}$$

Note:  $y = 0$  or  $1$  always

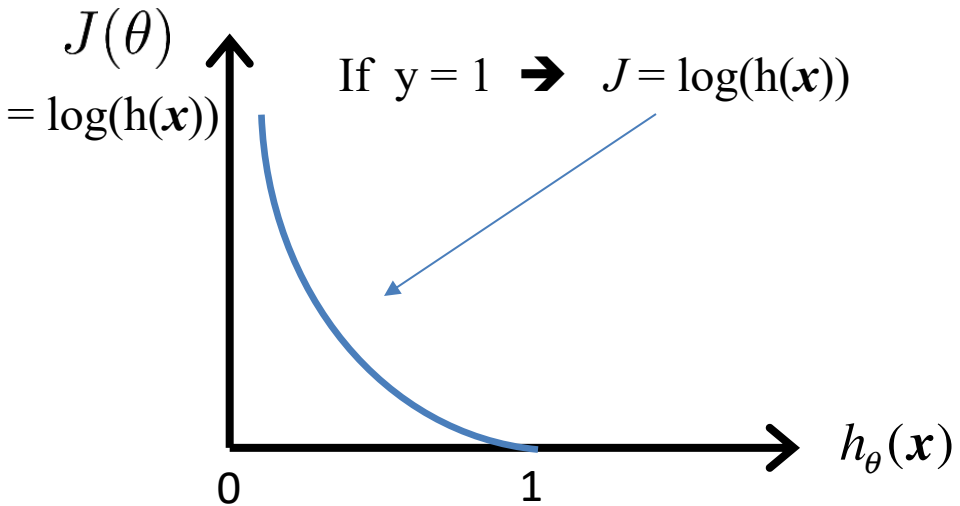
# \*Optional\*

- Why does this Cost function work?

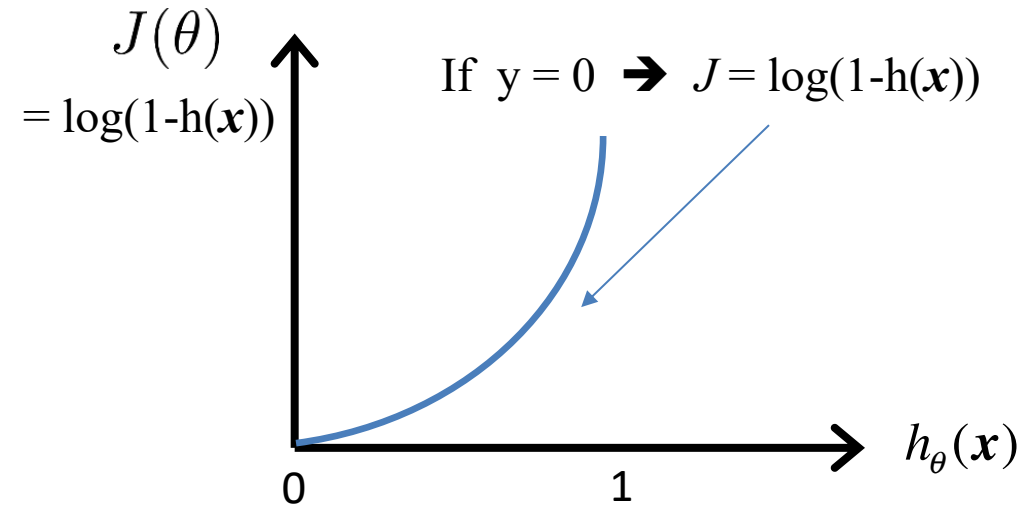
$$J(\boldsymbol{\theta}) = J(\theta_0, \theta_1, \dots, \theta_n)$$

Note:  $y = 0$  or  $1$  always

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



- when  $y=1$ , and  $h(\mathbf{x})$  is close to 1, the cost will be very small, and when  $y=1$ , and  $h(\mathbf{x})$  is close to 0, the cost will be too high.



- when  $y=0$ , and  $h(\mathbf{x})$  is close to 0, the cost will be very small, and when  $y=0$ , and  $h(\mathbf{x})$  is close to 1, the cost will be too high.

# Logistic Regression Cost Function

Cost Function:

$$\begin{aligned} J(\boldsymbol{\theta}) &= J(\theta_0, \theta_1, \dots, \theta_n) \\ &= -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)})) \right] \end{aligned}$$

To find the best parameters  $\theta$ :

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

# Gradient Descent for Multiple Variables

- Output (Probability):  $h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$
- Want  $\min_{\theta} J(\theta)$  to find  $\theta$
- Cost function:  $J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))]$

- **Gradient descent:**

Repeat {

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

}

(simultaneously update for every  $j = 0, \dots, n$ )

# Gradient Descent for Multiple Variables

- Output (probability):  $h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$
- Cost function:  $J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)})) \right]$

- **Gradient descent:**

$$\begin{aligned} \text{Repeat } \{ \quad & \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \\ & = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \\ & \} \quad \text{(simultaneously update for every } j = 0, \dots, n) \end{aligned}$$

- **Algorithm looks identical to linear regression (except for the definition of  $h(\mathbf{x})$ )!!!**

# Gradient Descent for Multiple Variables

- Output (probability):  $h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$
- Cost function:  $J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))]$
- **Gradient descent:**

$$\begin{aligned} \text{Repeat } \{ \quad & \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \\ & = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \\ & \} \quad \text{(simultaneously update for every } j = 0, \dots, n) \end{aligned}$$

- **Algorithm looks identical to linear regression (except for the definition of  $h(x)$ )!!!**





*Thank You!*

**Questions?**