



Machine Learning



1. Supervised – Input and output of each example specified
 1. Regression – Continuous output
 2. Classification – Discrete output
2. Unsupervised – Algorithm finds structures among dataset - clusters

Supervised
Machine Learning



Classification



Am I a cat or a dog?

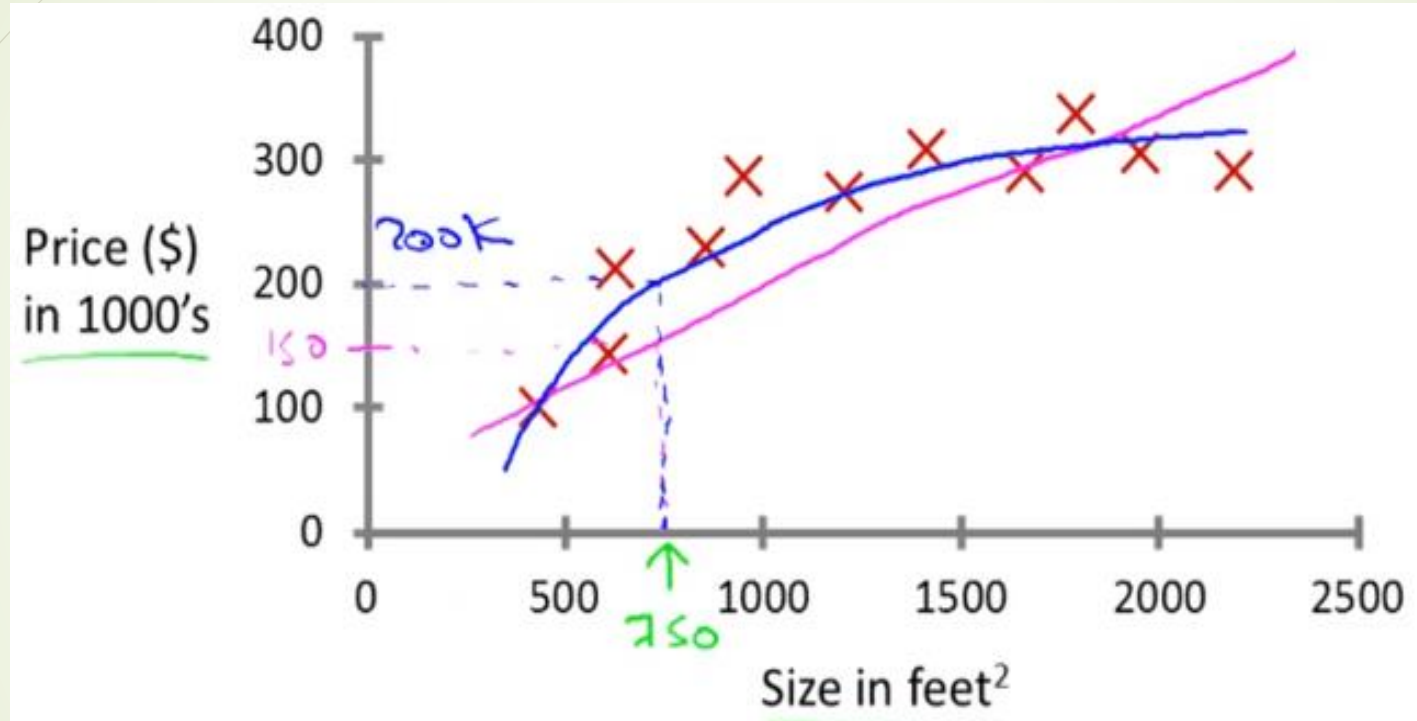
- Dataset consists of 100 pictures of 4 different animals
- Algorithm learns the features of the animals from the training set
- It can then classify the a random picture as one of the animals
- Called **multiclass classification**
- **Discrete output**

Supervised
Machine Learning



Regression

What's the price of the house?

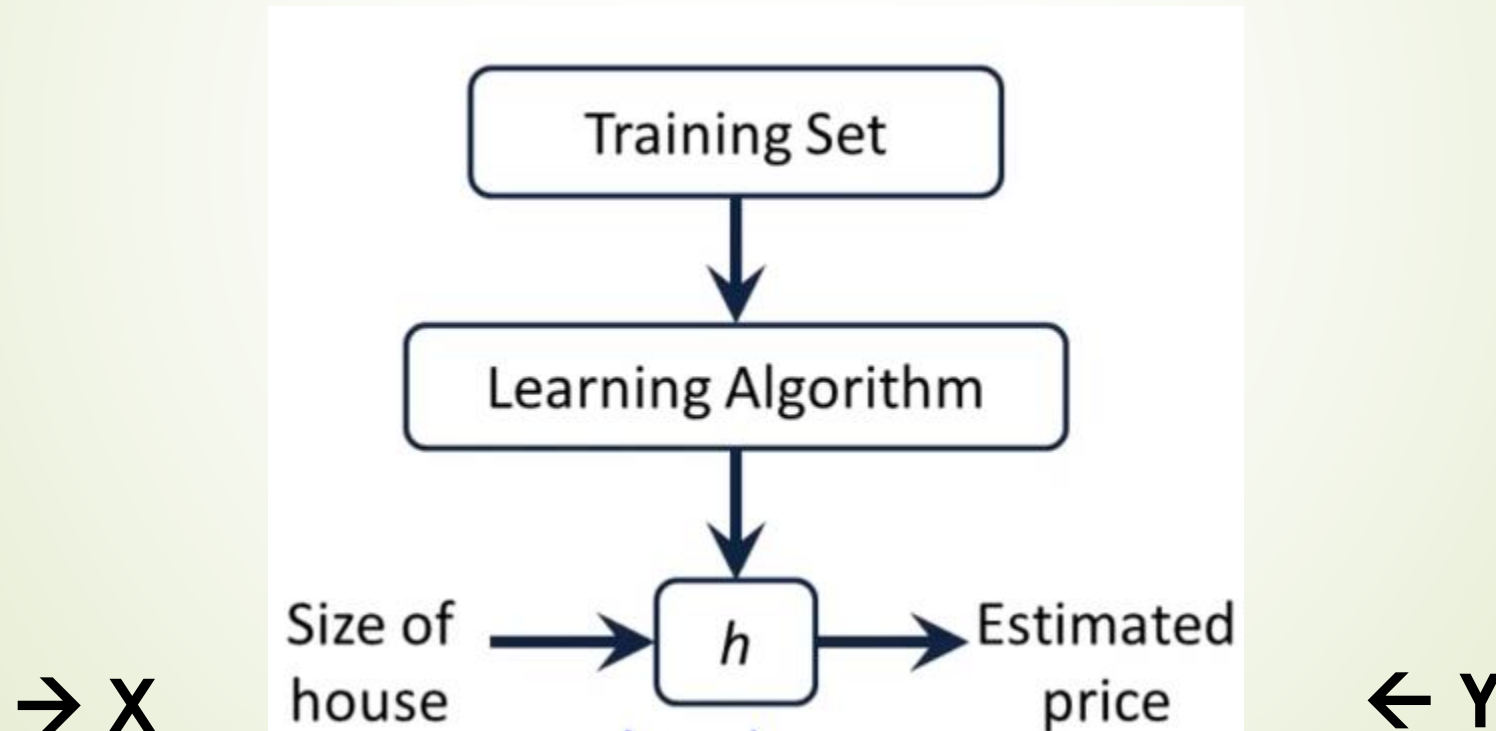


- **Learning algorithm** maps the size of the house to the price
- Predicts the price for the given input house size

Linear Regression – Univariate

- $m \rightarrow$ training examples
- $x \rightarrow$ input variables / features
- $y \rightarrow$ output features

$h \rightarrow$ Hypothesis function



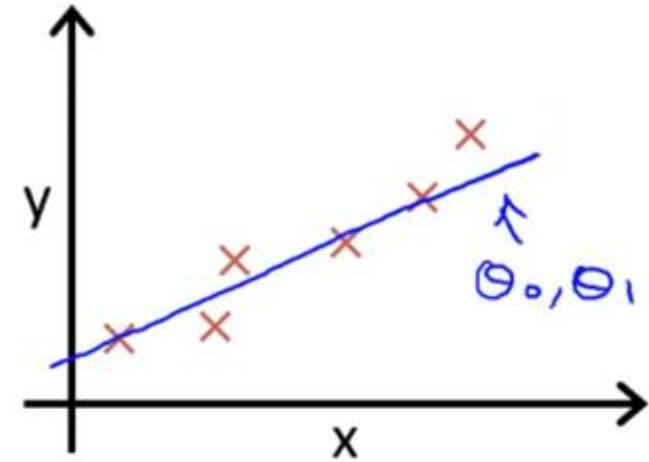
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

θ_0

BIAS – Offset all predictions made

θ_1

SLOPE of the line depicted by hypothesis function



Idea: Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples (x, y)

COST FUNCTION

➤ Error – $[\mathbf{h}(\mathbf{x}) - \mathbf{y}]$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

- Aim – Minimize the cost J
- Mean Squared Error cost function

Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Parameters:

$$\theta_0, \theta_1$$

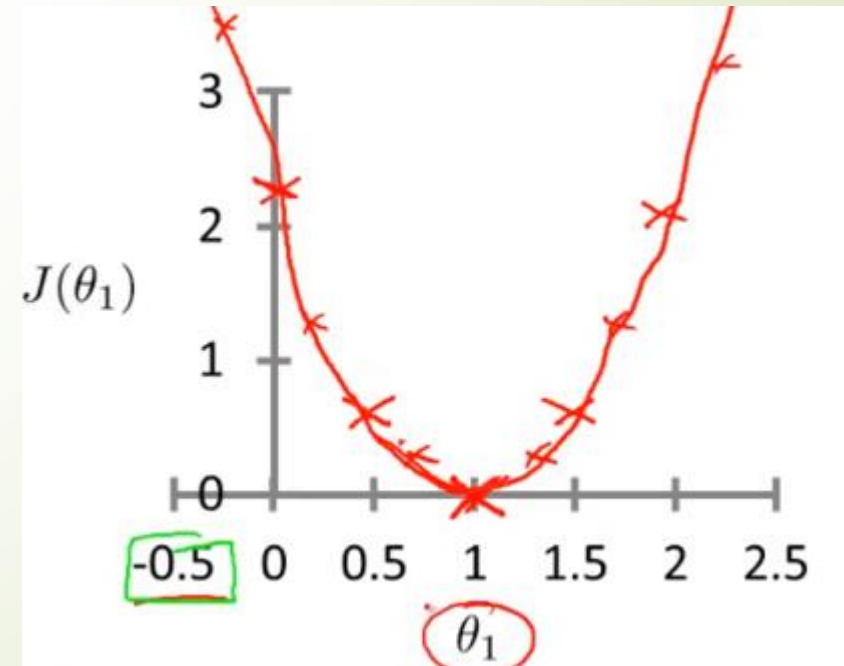
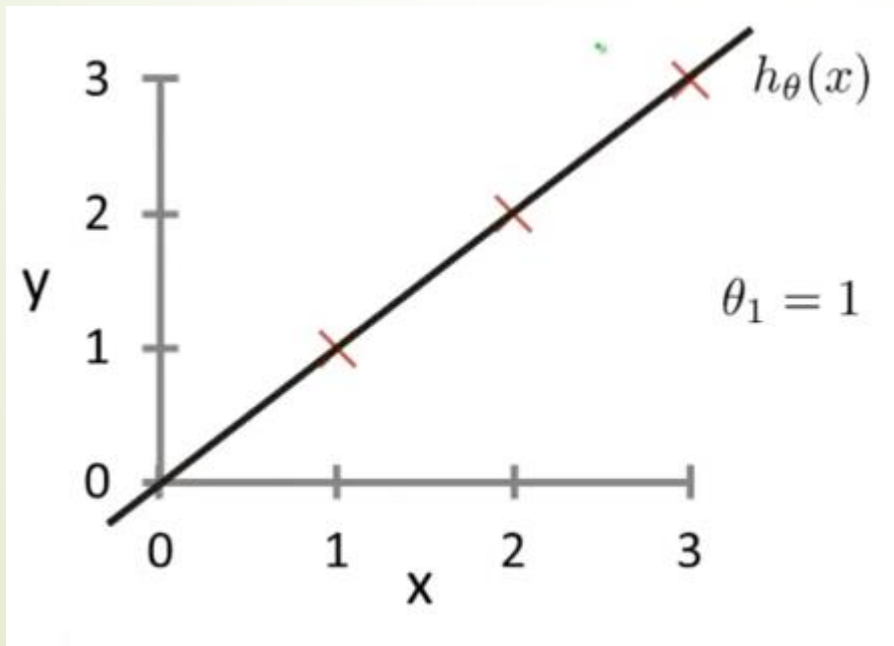
Cost Function:

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

Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

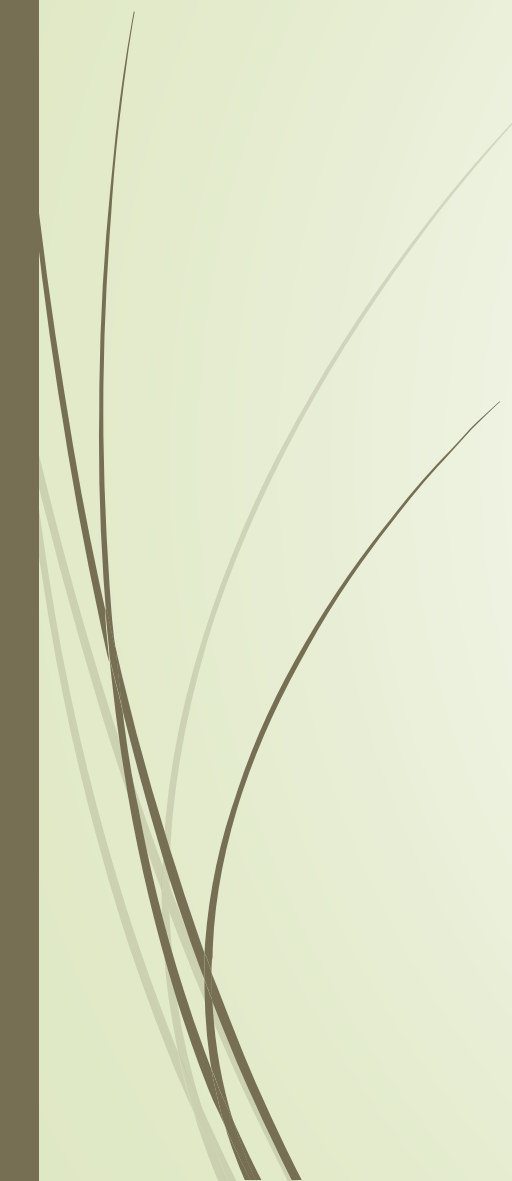
Simplified Model

- $h(x) = \theta_1 x \dots (\theta_0 = 0)$
- When $\theta_1 = 0$
- When $\theta_1 = 0.5$
- When $\theta_1 = 1$
- Find the value of $J(\theta)$ for each value of θ_1





Examples

- ▶ Lemonade business – Revenue is dependent on temperature
 - ▶ Amazon – recommends items to users based on user behavior
 - ▶ Predict number of shoppers passing in front of billboard – predict maximum to bid for advertisement
- 

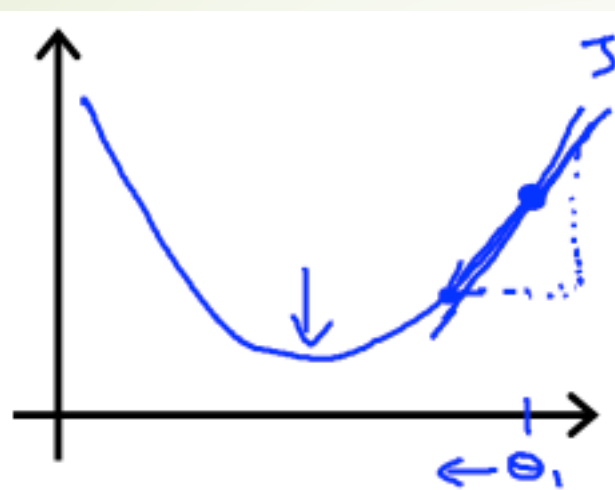
Gradient Descent

- algorithm to minimize cost function

- Start with some parameter values (θ_0, θ_1)
- Change values to reduce $J(\theta_0, \theta_1)$ until minimum

repeat until convergence {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (for $j = 0$ and $j = 1$)
}

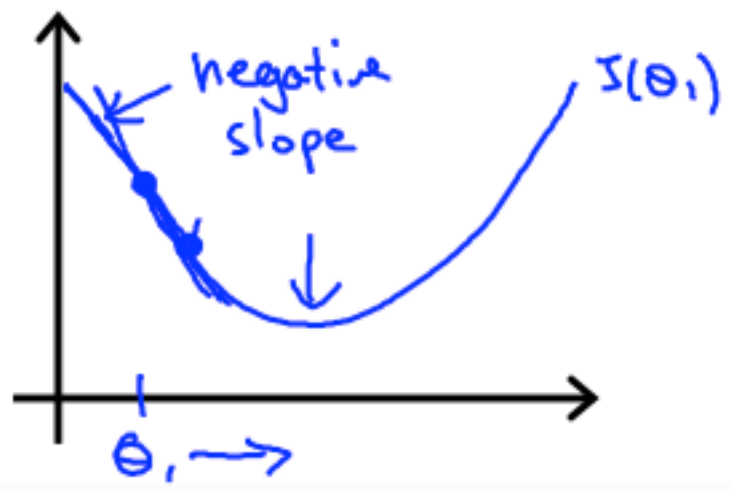
- α - learning rate (controls how big a step is taken downhill)
- Simultaneous Update – Compute the RHS for both and assign to new (θ_0, θ_1)
- Converges even with a **fixed rate** α - the derivative approaches zero as we approach the bottom of the function – gradient descent will take smaller steps



$$\theta_1 := \theta_1 - \alpha \left(\frac{\partial}{\partial \theta_1} J(\theta_1) \right)$$

≥ 0

$$\theta_1 := \theta_1 - \underline{\alpha} \cdot (\text{positive number})$$



$$\theta_1 := \theta_1 - \alpha \left(\frac{\partial}{\partial \theta_1} J(\theta_1) \right)$$

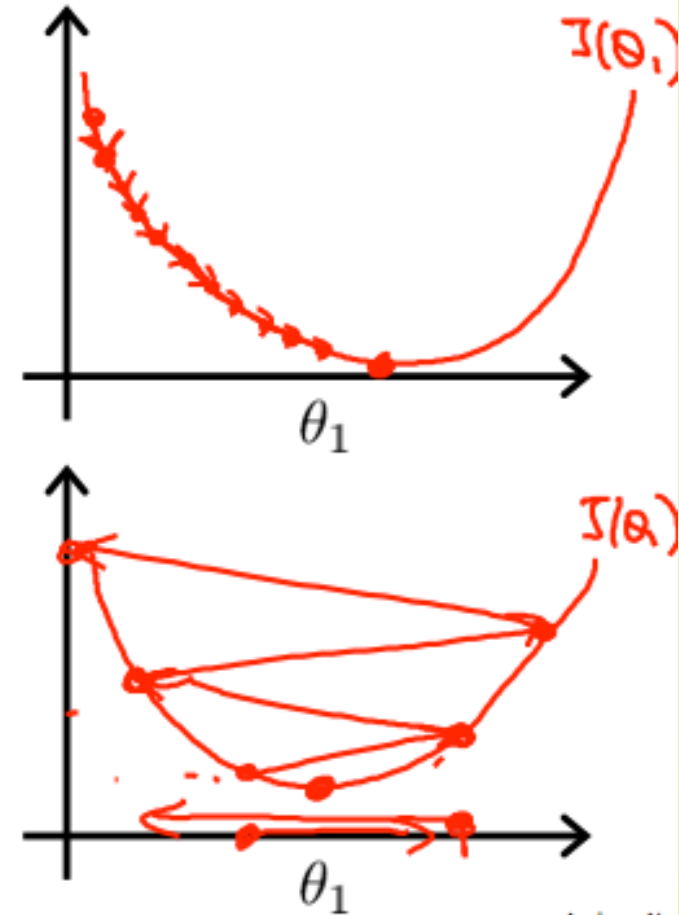
≤ 0

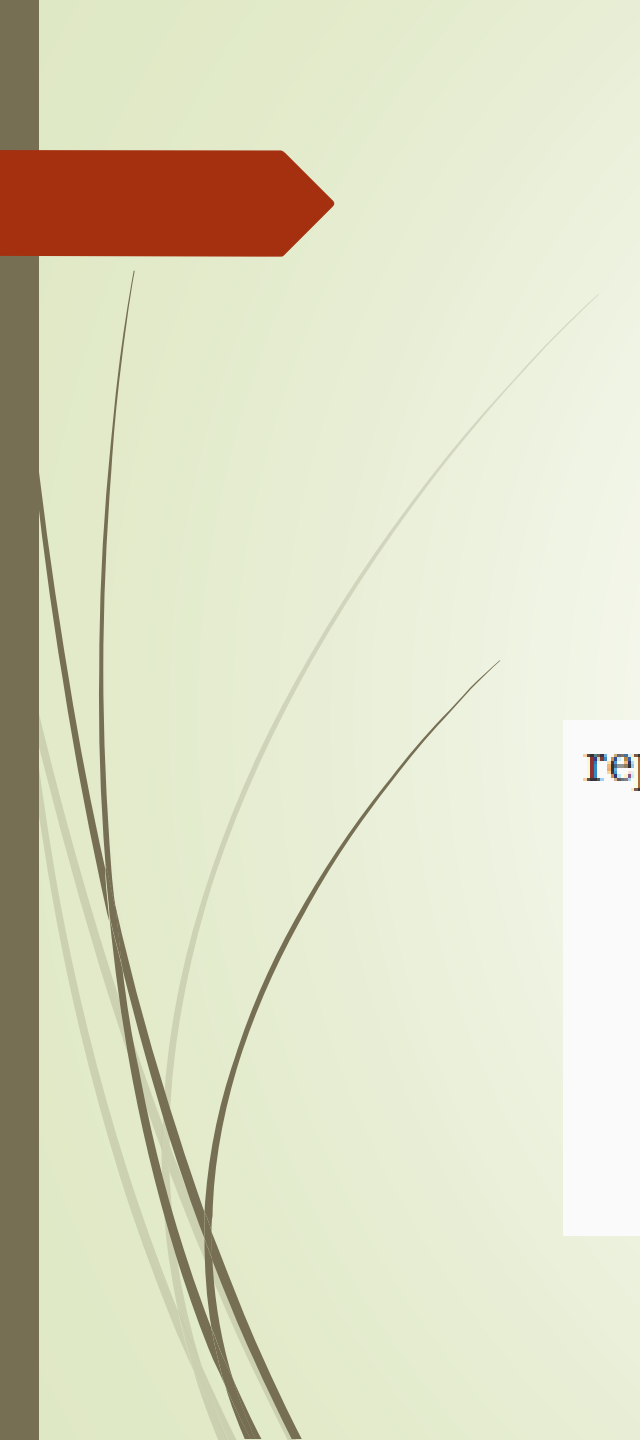
↑ ↑
negative number

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.




$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\ &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) x_j\end{aligned}$$

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_\theta(x_i) - y_i) x_i)$$

}

Multivariate Linear Regression

- ▶ $x_0, x_1, x_2, \dots, x_n \rightarrow$ input
- ▶ $Y \rightarrow$ output
- ▶ E.g. size, number of rooms, age of home, number of floors \rightarrow price

Previously:
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

- ▶ Convenience – introduce $x_0 = 1$, such that **'x' and 'θ'** have **'n+1'** elements each

$$h_{\theta}(x) = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

Gradient Descent - Multivariate

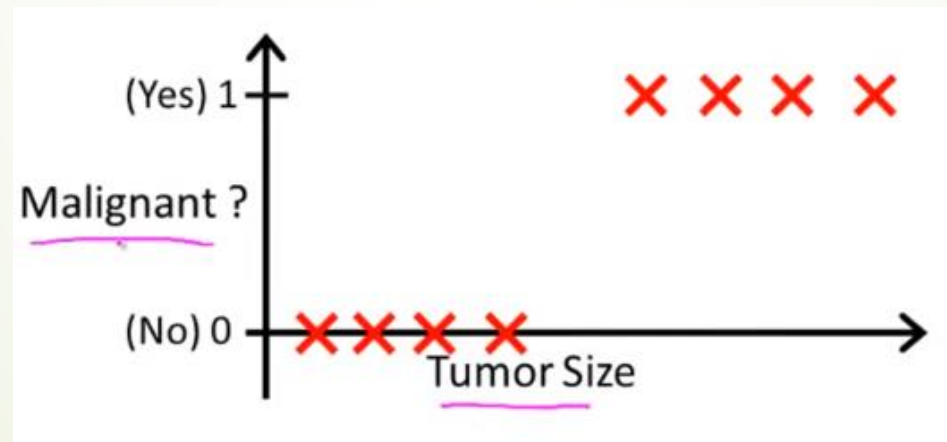
repeat until convergence: {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \text{for } j := 0 \dots n$$

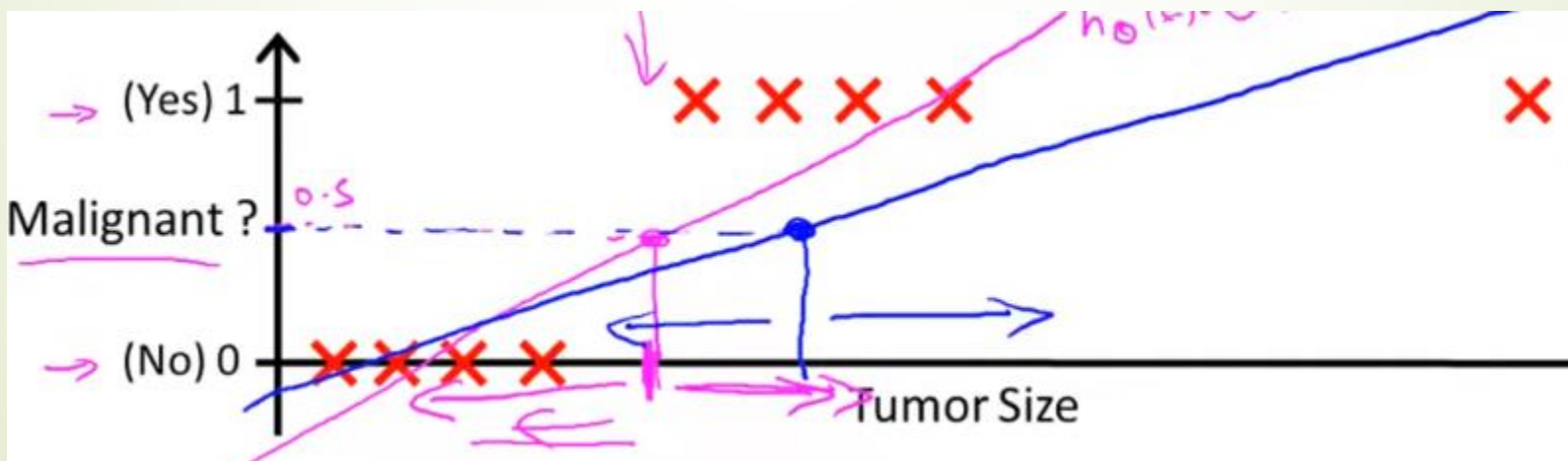
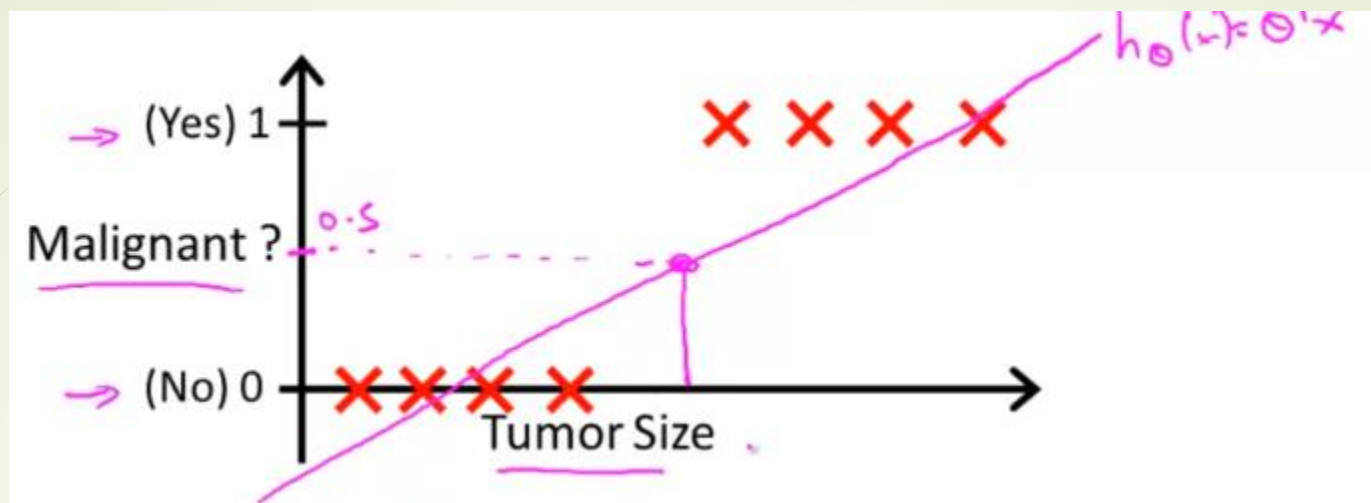
}

Logistic Regression

- Method for classifying data into **discrete** outcomes
- **Classification**, not regression
- Purpose of logistic regression – because linear regression is not useful for classification problems.
- $y = 0$ or $y = 1$, but using Linear Regression, $h(x)$ can be > 1 or < 0 too.
- Using log regression, $0 \leq h(x) < 1$



- Threshold the hypothesis output at 0.5 – if $h(x) > 0.5$, $y = 1$; if $h(x) < 0.5$, $y = 0$

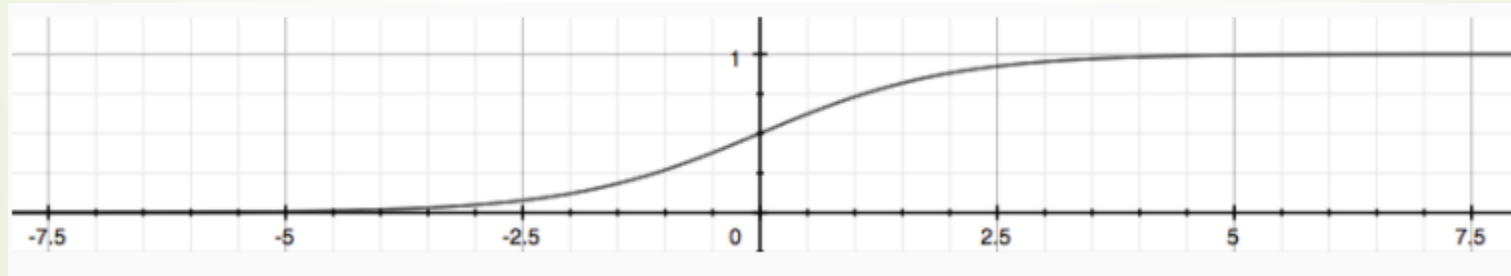


- $h(x) \rightarrow \theta^T x$... linear regression
- $h(x) \rightarrow g(\theta^T x)$... **logistic regression**

- **$g(z) = 1 / (1 + e^{-z})$**

- This is the sigmoid or the logistic function

- **Therefore, $h(x) \rightarrow 1 / (1 + e^{-\theta^T x})$**



- **This function maps any real number to interval (0,1)**
- $h(x)$ – probability that output is 1 – example,

Decision Boundary

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

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

Sigmoid function –

When **input (z)** is greater or equal to 0, output of function i.e. **h(x)** is greater than 0.5

If input to 'g' is $\theta^T x \rightarrow$

$$\theta^T x \geq 0 \Rightarrow y = 1$$

$$\theta^T x < 0 \Rightarrow y = 0$$

Decision Boundary – Line separating area where $y = 0$ and $y = 1$. Created by our hypothesis function

Cost Function – Logistic regression

- It is the penalty the learning algorithm has to pay if it outputs $\mathbf{h}(\mathbf{x})$ and actual label is \mathbf{y} .
- If we use the same cost function as linear regression, the cost function will be a non-convex function (gradient descent won't work as it will have many local optimums) as $h(x)$ is non linear.

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

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

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

- If $y = 1$, cost $\rightarrow \infty$ if prediction is $h(x) = 0$ (**high penalty**) and minimum cost if $h(x) = 1$.
- If $y = 0$, cost $\rightarrow \infty$ if prediction is $h(x) = 1$ (**high penalty**) and minimum cost if $h(x) = 0$.

