

# COMP/EECE 7/8740 Neural Networks

Topics:

## **Linear regression**

- Correlation and linear correlation
- Linear regression with single and multiple variables

Md Zahangir Alom  
Department of Computer Science  
University of Memphis, TN

# Recall: Covariance and interpreting covariance

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{n - 1}$$

$\text{cov}(X, Y) > 0 \rightarrow$  X and Y are positively correlated

$\text{cov}(X, Y) < 0 \rightarrow$  X and Y are inversely correlated

$\text{cov}(X, Y) = 0 \rightarrow$  X and Y are independent

# Correlation coefficient

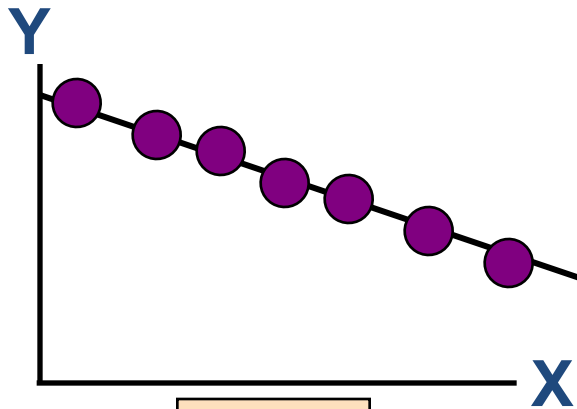
- Pearson's Correlation Coefficient is standardized covariance (unitless):

$$r = \frac{\text{covariance}(x, y)}{\sqrt{\text{var } x} \sqrt{\text{var } y}}$$

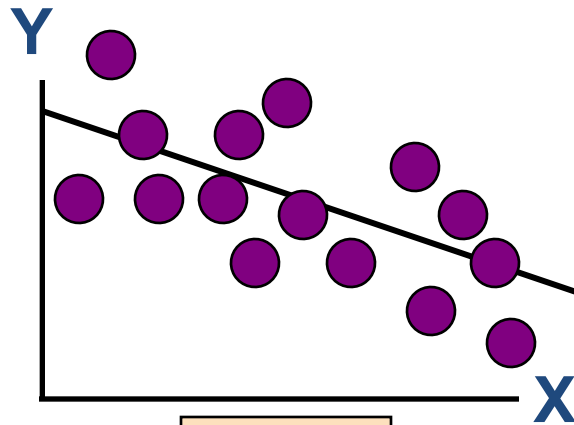
# Correlation

- Measures the relative strength of the *linear* relationship between two variables
  - Ranges between  $-1$  and  $1$
  - The closer to  $-1$ , the stronger the negative linear relationship
  - The closer to  $1$ , the stronger the positive linear relationship
  - The closer to  $0$ , the weaker any positive linear relationship

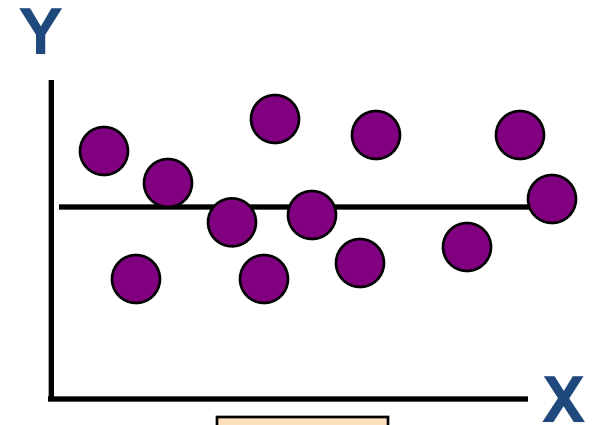
# Scatter Plots of Data with Various Correlation Coefficients



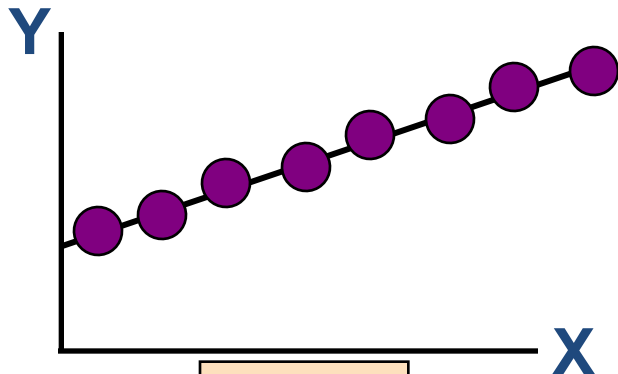
$$r = -1$$



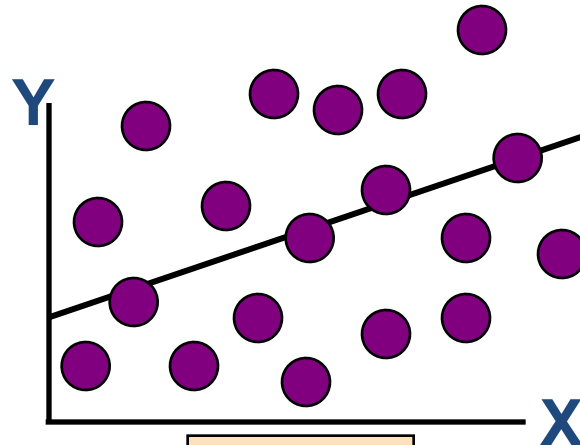
$$r = -.6$$



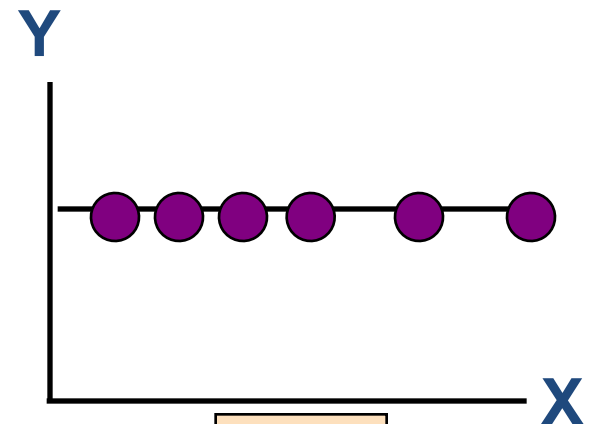
$$r = 0$$



$$r = +1$$



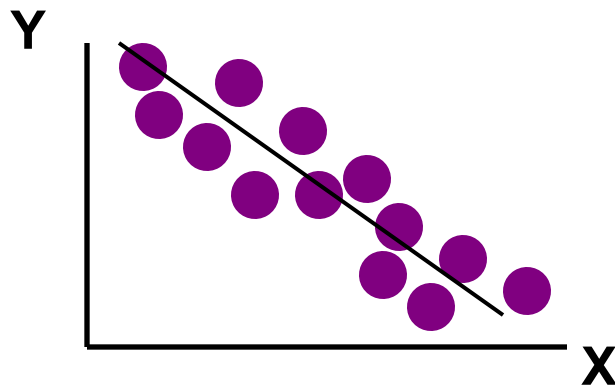
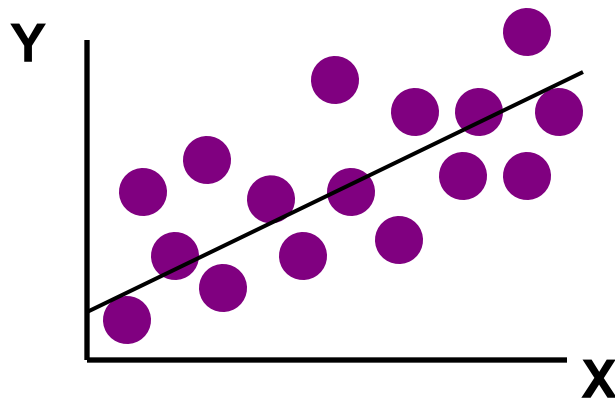
$$r = +.3$$



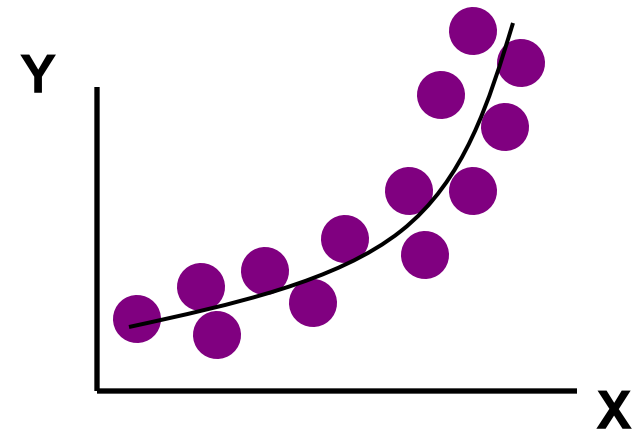
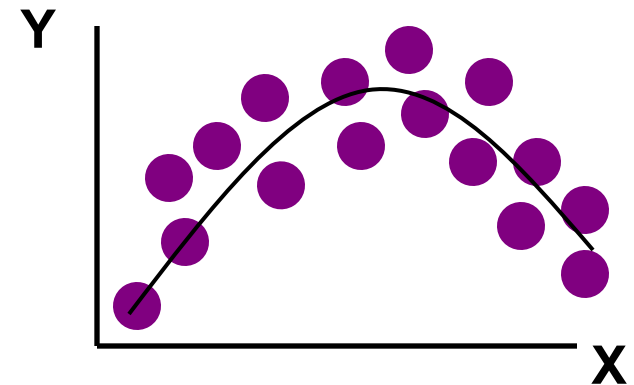
$$r = 0$$

# Linear Correlation

**Linear relationships**

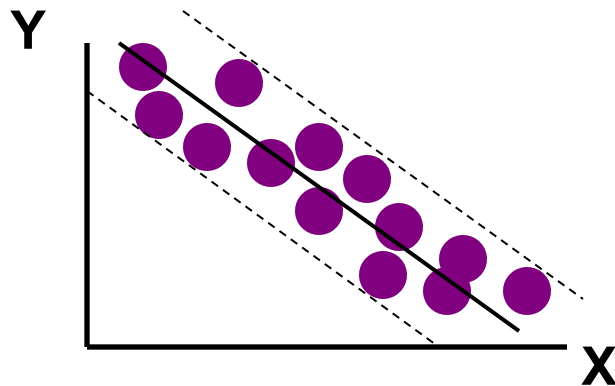
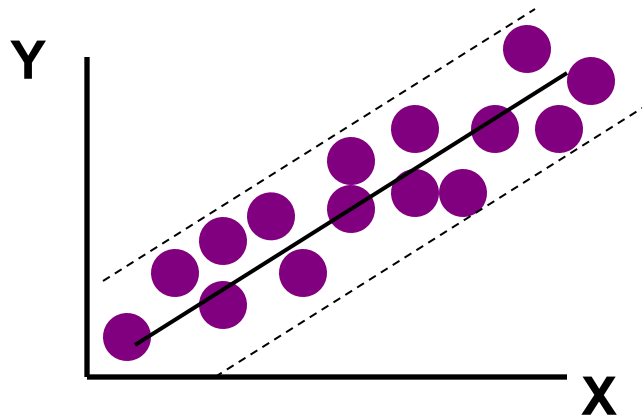


**Curvilinear relationships**

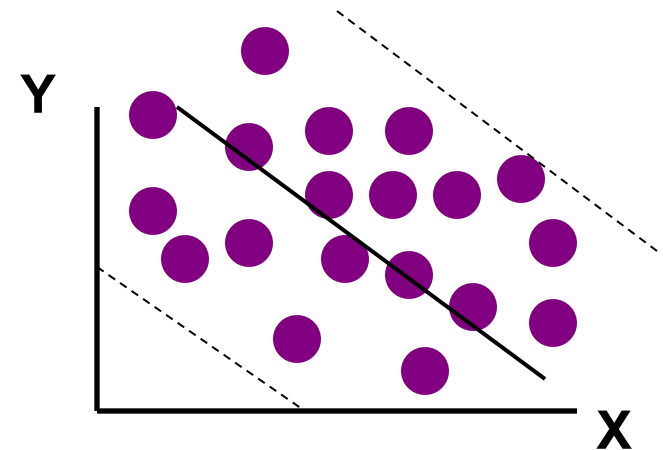
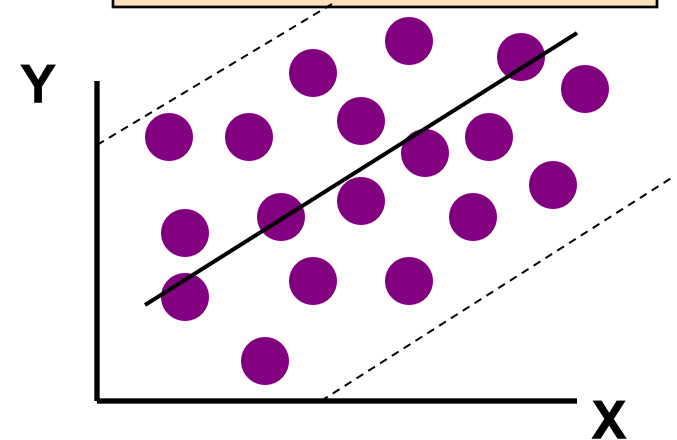


# Linear Correlation

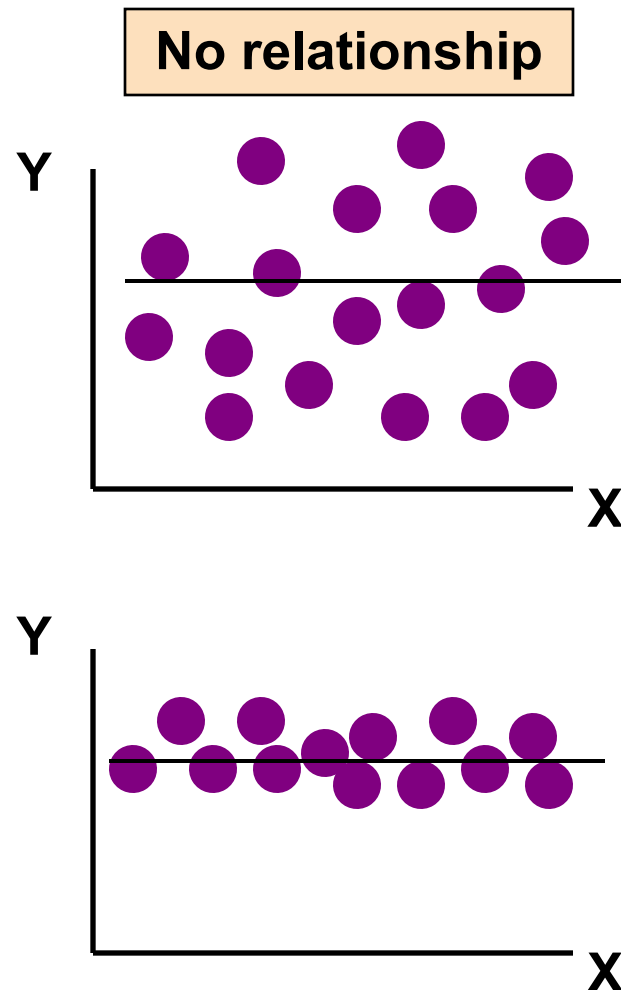
**Strong relationships**



**Weak relationships**



# Linear Correlation





# Linear regression

## What is Linear Regression?

Linear regression is a **basic and commonly used** type of predictive analysis.

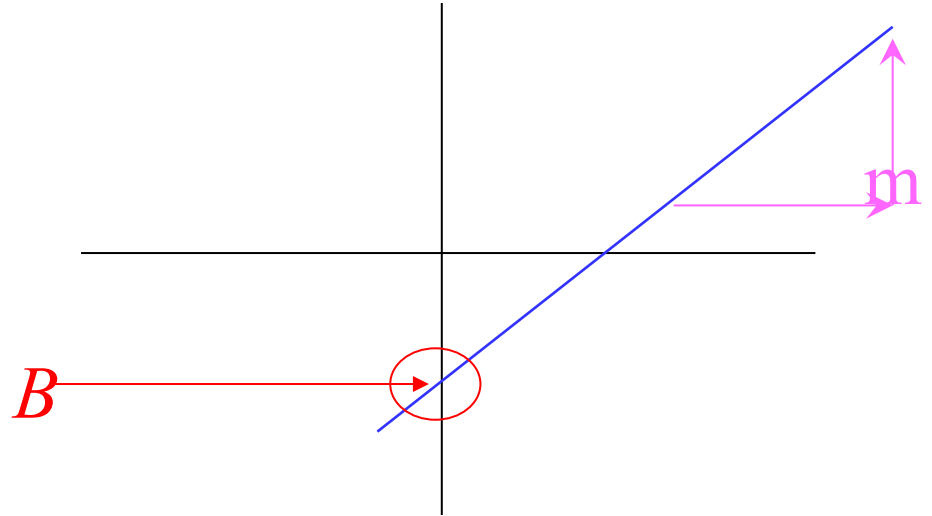
LR examine two things:

- Does a set of predictor variables do a good job in **predicting an outcome (dependent) variable**?
- Which **variables in particular are significant predictors** of the outcome variable?
- LR use to explain the **relationship between one dependent variable and one or more independent variables**.
- Different between linear correlation and regression:
  - In correlation, the two variables **are treated as equals**.
  - In regression, one variable is considered independent (=predictor) variable ( $X$ ) and the other the dependent (=outcome) variable  $Y$ .

# What is “Linear”?

- Remember this:

$$Y=mX+B$$



## What's Slope ( $m$ )?

A slope of 2 means that every 1-unit change in  $X$  yields a 2-unit change in  $Y$ .

**Prediction** : If you know something about  $X$ , this knowledge helps you predict something about  $Y$ . (sound like conditional probabilities?)

# Regression equation...

- Expected value of  $y$  at a given level of  $x$ =

$$E(y_i / x_i) = \alpha + \beta x_i$$

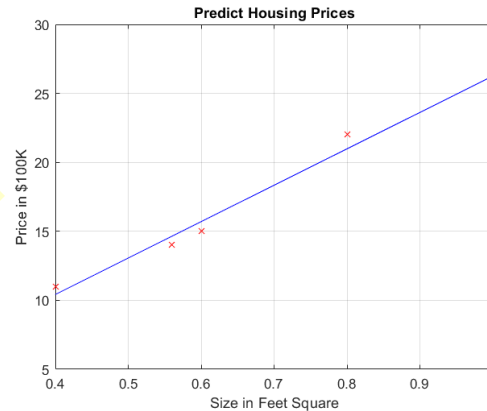
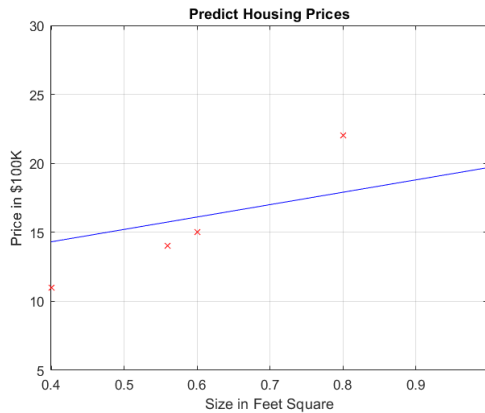
- Predicted value for an individual...

$$y_i = \underbrace{\hat{\alpha} + \beta * x_i}_{\text{Fixed – exactly on the line}} + \boxed{\text{random error}_i}$$

Fixed –  
exactly on  
the line

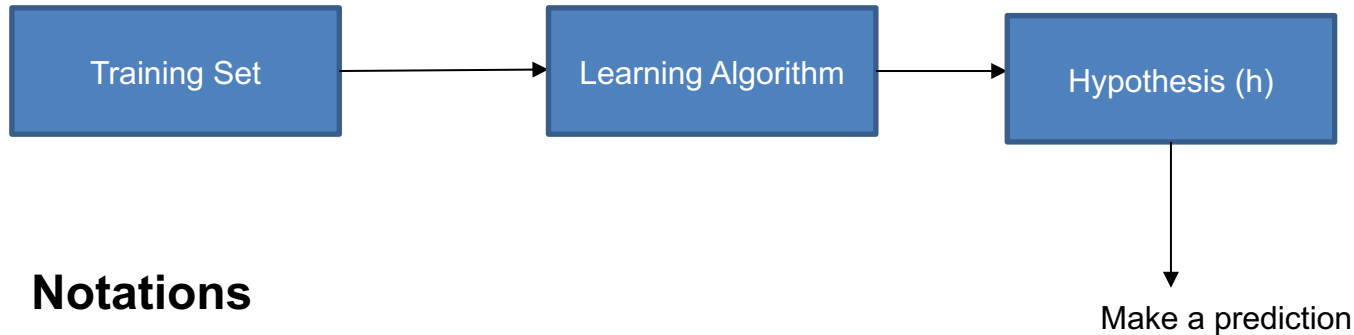
Follows a normal  
distribution

# Linear regression with one variable



X	Y
0.4	11
0.56	14
0.6	15
0.8	22
1	26

# Processing pipeline



## Notations

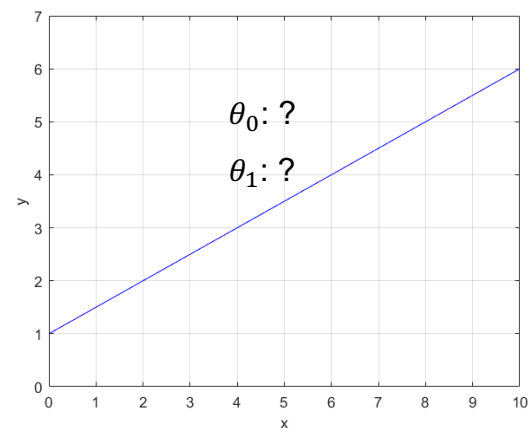
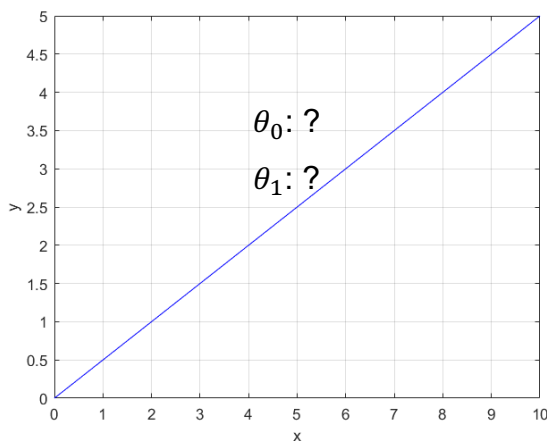
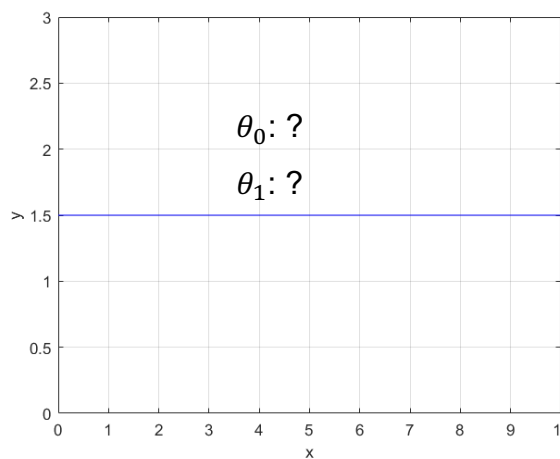
- ❖  $M$  – Number of Training Examples
- ❖  $x$ 's – Input Variables / Features
- ❖  $y$ 's – Output variable / “target” variable
- ❖  $h$  – Hypothesis “Proposal”

Each training example –  $(x,y)$

‘ $h$ ’ is a function that maps ‘ $x$ ’ to ‘ $y$ ’

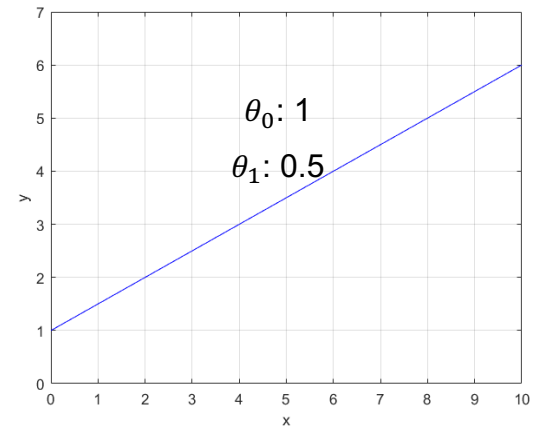
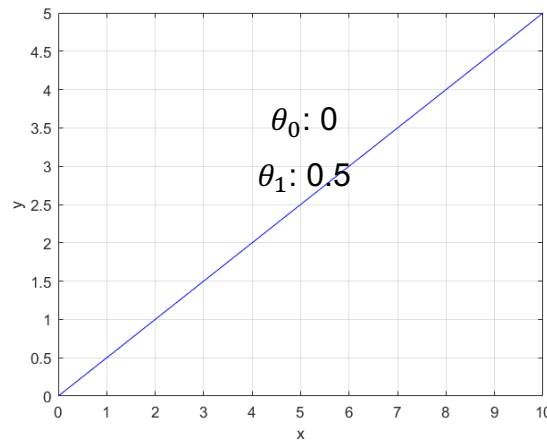
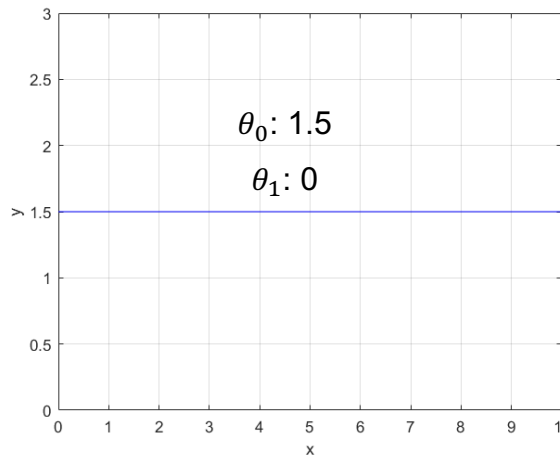
# Linear Regression with one variable

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



# Linear Regression with one variable

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



- We have to choose  $\theta_0$  and  $\theta_1$  such that  $h_{\theta}(x) \cong y$
- Minimize  $h_{\theta}(x) - y$  or  $y - h_{\theta}(x)$
- $(h_{\theta}(x) - y)^2$  is a good function to minimize – Loss Function

# Linear regression with one variable

- We have to optimize the values across the entire dataset: i.e., **every sample needs to contribute to loss function**
- Cost function: J

$$J = \text{minimize}_{\theta_0, \theta_1} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$$



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



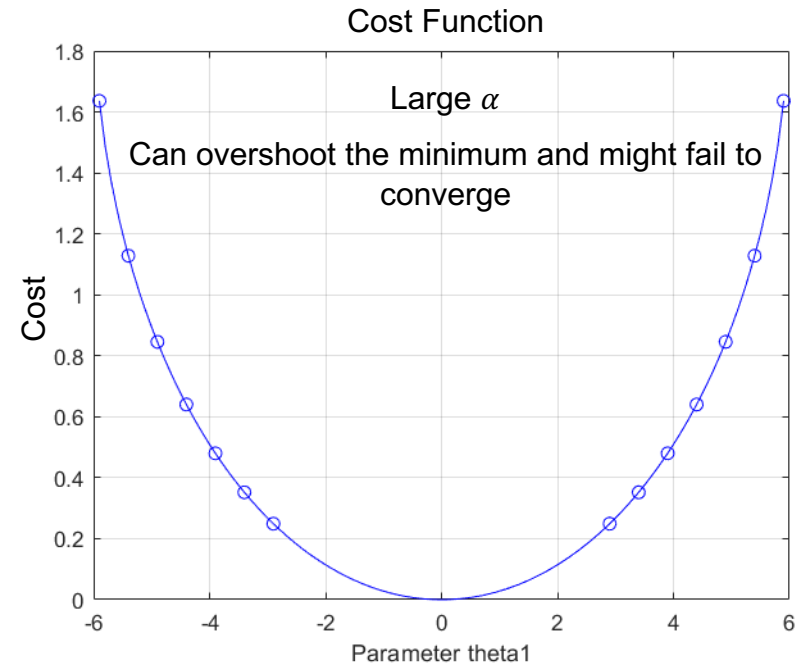
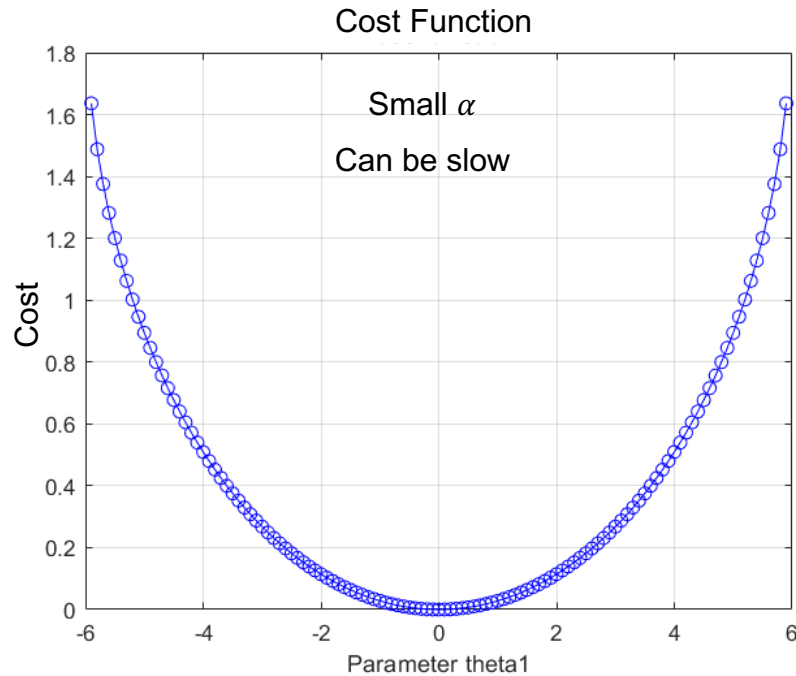
# Gradient Descent

- Minimize cost function:
  - Start with a random
  - Keep updating  $\theta_0$  &  $\theta_1$  to reduce  $J$
  - Find the global minimum
  - Find the values of  $\theta_0, \theta_1$  at which  $J$  is minimum
  - Differentiation is the way to go!

$$\gg \theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

- $\alpha$  – Learning rate – Big steps (Large value), Baby Steps (Small value)

# Learning rate example



# Learning rate example



# Update the parameters : $\theta_0, \theta_1$

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$$

$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2 \quad \text{and} \quad \theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$$

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

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

# Partial differentiation

$$\begin{aligned}\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i)^2 \\ &= \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^m \theta_0^2 + \theta_1^2 x^{i2} + 2\theta_0 \theta_1 x^i \\ &\quad + y^{i2} - 2(\theta_0 + \theta_1 x^i) y^i\end{aligned}$$

$$= \frac{1}{2m} \sum_{i=1}^m 2\theta_0 + 0 + 2\theta_1 x^i + 0 - 2y^i$$

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

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_1} \frac{1}{2m} \sum_{i=1}^m \theta_0^2 + \theta_1^2 x^{i2} + 2\theta_0 \theta_1 x^i + y^{i2} - 2(\theta_0 + \theta_1 x^i) y^i$$

$$= \frac{1}{2m} \sum_{i=1}^m 0 + 2\theta_1 x^{i2} + 2\theta_0 x^i + 0 - 2x^i y^i$$

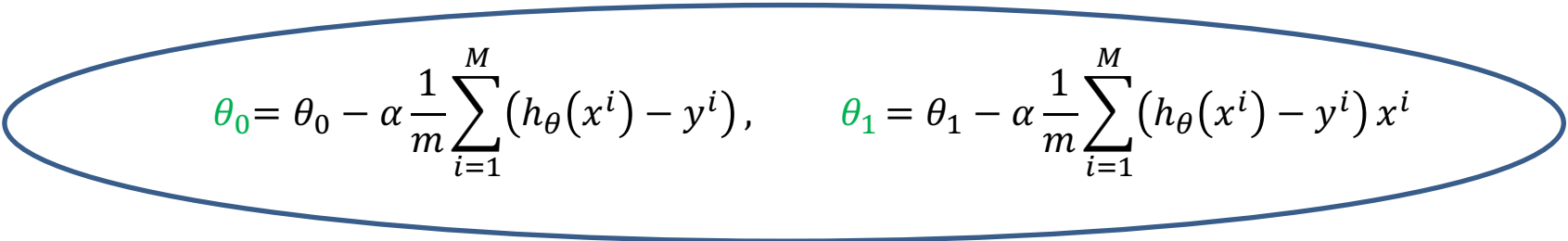
$$= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) x^{(i)}$$

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

# Update the parameters : $\theta_0, \theta_1$

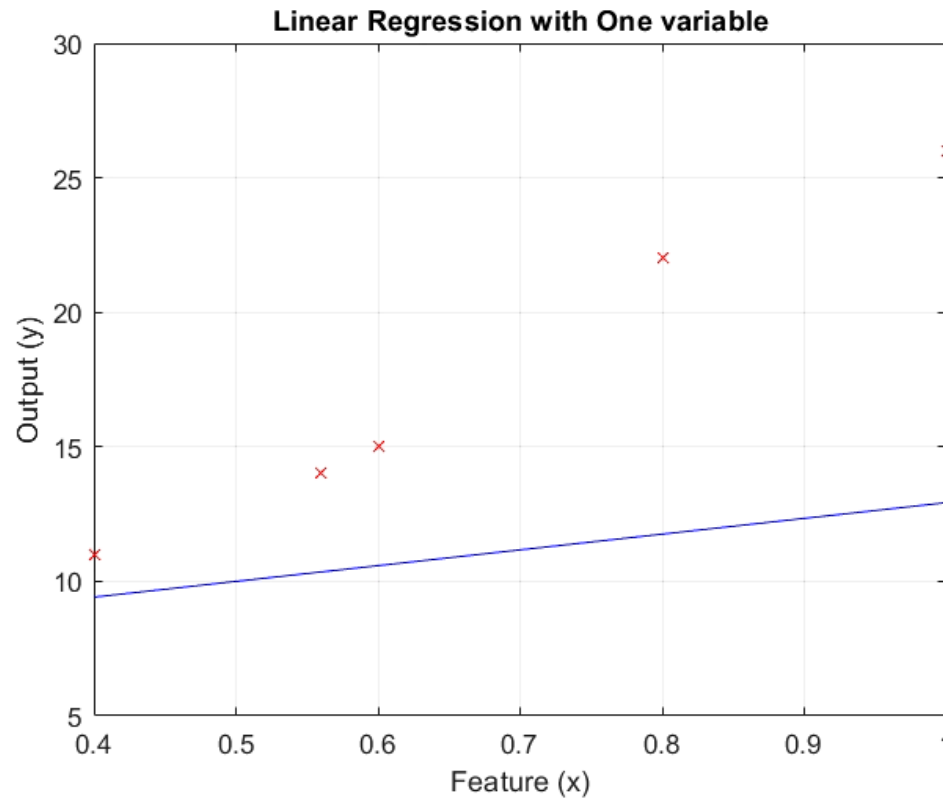
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i),$$

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


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

**Repeat this process for 'N' iterations**

# Linear Regression Demo! (Video)



# Feature engineering?





# Feature engineering

- Size in square feet
- # Bedrooms
- # Restrooms
- School district
- School quality
- Crime rate
- Highway accessibility

And so on...

**We can engineer 'n' number of features and so on.**



# Linear regression with multiple variables

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$



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

- It doesn't make sense to type out this equation manually for every variable,
- You'll have to **modify the equation every time a feature is added** especially in problems where you are engineering features
- Moreover, the **gradient descent function would also** vary as you add more variables



# Linear regression with multiple variables

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

Let's resort to vectorization:

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

- However, size of  $\theta$  is  $(1 \times n+1)$  and  $x$  is  $(M \times n)$
- Note that,  **$x$  is a feature matrix** instead of an array

	<h3>FOR Loops</h3> <pre>J = 0, dw1 = 0, dw2 = 0, db = 0 for i = 1 to m:     z<sup>(i)</sup> = w<sup>T</sup>x<sup>(i)</sup> + b ←     a<sup>(i)</sup> = σ(z<sup>(i)</sup>) ←</pre>
	<h3>Vectorization</h3> <pre>z = w<sup>T</sup>X + b = np.dot(w.T, X) + b A = σ(z)</pre>

# Hypothesis

- Let's introduce  $x_0$  whose value is 1 to make sure the dimensions match – [ $x_0 = 1$  for all samples present in the dataset]

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \dots \dots \theta_n x_n$$

- Now, size of  $\theta$  is **(1 x n+1)** and  $x$  is **(M x n+1)**
- $h_{\theta}(x) = \theta x^T$  should be possible now
- $h_{\theta}(x)$  - Size would be **(1 x M)**

**Hypothesis can be computed for multiple variables now!**

# Cost function

- $J = \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$  - from previous slide
- $y$  is **1 x M** and  $h_{\theta}(x^i)$  is **1 x M**

**Do we need to make any changes to this equation?**

- No need of for-loops for cost calculation

# Gradient Descent

From previous slides,

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

**Is there any generic way we can write such that we can update all weights simultaneously?**

# Gradient Descent

- However, now we have until  $\theta_n$  and  $x$  is  $(M \times N+1)$
- Instead of updating parameters one by one, it's good to update them all simultaneously
- $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$
- $\delta = \frac{\partial}{\partial \theta} J(\theta)$

# Vectorized implementation of gradient descent

- Generic formula for simultaneously updating all  $\theta_j$
- $\delta = \frac{1}{m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i) x^i$

## Notes for programming:

- Dimensions of  $(h_{\theta}(x) - y)$  - (1 x M)
- Dimensions of  $x$  - (M x n+1)
- $\delta = (h_{\theta}(x) - y) x \Rightarrow$  Note that, it is matrix multiplication
- Dimensions of  $\delta$  - (1 x n+1)

$$\theta = \theta - \alpha \delta$$

Gradient descent for multiple variables is complete now!



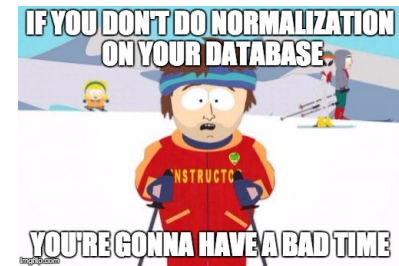
# One extra step to be completed!

- $x_1$  - Size in square feet – **Range of values: (100-10,000)**
- $x_2$  - # Bedrooms – **Range of values: (0-6)**
- $x_3$  - # Restrooms – **Range of values: (0-5)**
- $x_4$  - School quality – **Range of values: (1-100)**
- Indifferent range for every feature – Make it forever to converge or even might make it biased to certain features
- It's important to normalize the data

# Data normalization

- There are multiple ways to normalize the data
- Let's adopt to the one where we can zero center the values

$$\hat{x} = \frac{x - \mu}{\sigma}$$



$\mu$  – mean for each feature,  $\sigma$  – standard deviation for each feature

- No need of for-loops for this as well.

**Note that, please do this normalization step before adding  $x_0$**

# Steps for linear regression with multiple variables

- Read the dataset : (Input features and Output values)
- Normalize features
  - $\hat{x} = \frac{x - \mu}{\sigma}$
  - Add  $x_0 = 1$  to all samples (Normalized  $\hat{x}$ )
  - new Dimensions of  $x$  –  $M \times (n+1)$
- Implement cost function after computing hypothesis
  - $h_{\theta}(x) = \theta x^T$
  - $J = \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$
- Implement gradient descent function for multiple variables
  - $\theta = \theta - \alpha \delta$
  - $\delta = \frac{1}{m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i) x^i$

*\*No need any for-loops other than one for iteration*

# Assignment 1 for practice at home

- Compute cost Function
  - Determine the number of samples
  - Compute  $h_{\theta}(x) = \theta_0 + \theta_1 x$
  - Compute  $J = \frac{1}{2m} \sum_{i=1}^M (h_{\theta}(x^i) - y^i)^2$
- Gradient Descent Function
  - Determine the number of samples
  - Loop for iterations
  - Compute  $h_{\theta}(x) = \theta_0 + \theta_1 x$
  - Compute cost using your own “compute cost” function
  - Calculate updated  $\theta_0$  and  $\theta_1$
  - Repeat for all iterations!

# Summary

- Linear correlation
- Linear regression with and single and multiple variables
- What's next ?
  - Logistic Regression (LR)
  - Variants...