

# Linear Regression

---

Dr Mehrdad Ghaziasgar

## Content - Part 1

- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

# Simple Linear Regression

## Model Representation

---

- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

## Model Representation - Introduction

- Running example: predicting housing prices
- We obtain information about:
  - Houses e.g. size (sq. metres), no. of bedrooms, no. of bathrooms, no. of garages, frontage (metres), no. of storeys, garden size (sq. metres) etc. etc.
  - Price that each house last (recently) sold for
- Use a learning algorithm to build a model to predict housing prices

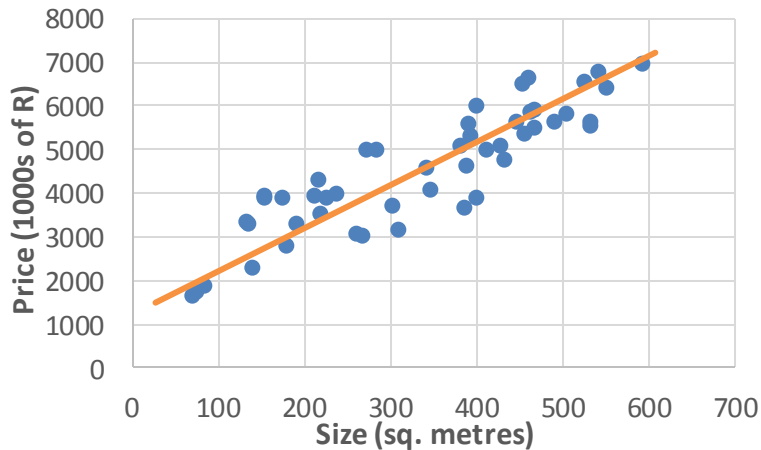
## Model Representation - Introduction

Size (sq. m) ( $x$ )	Price (1000s of R) ( $y$ )
460	6639
70	1681
155	3969
429	5095
...	...

Training set of house sizes and prices

- We obtain the data above

## Model Representation - Introduction



### Supervised learning:

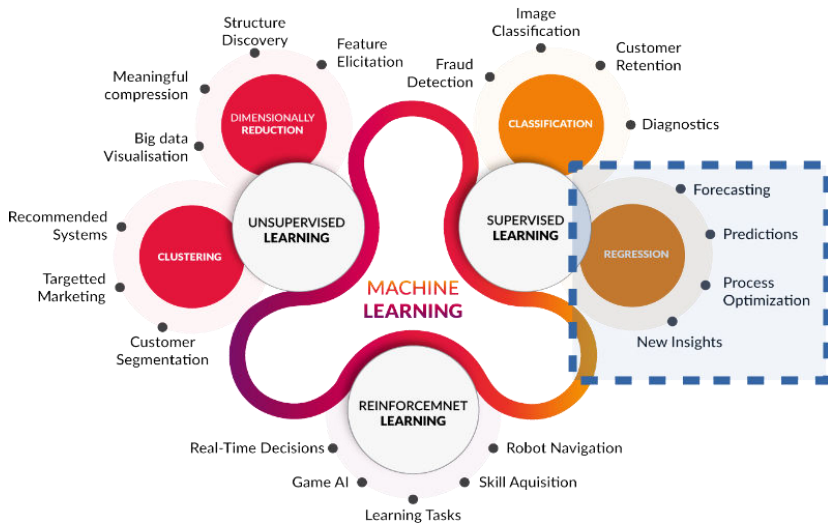
- Given the “right” answer for each data point

### Regression:

- Predict a real-valued (continuous) output e.g. price

(Side note: the other branch of supervised learning is classification i.e. predict discrete output e.g. category)

# Model Representation - Introduction



## Model Representation - Notation

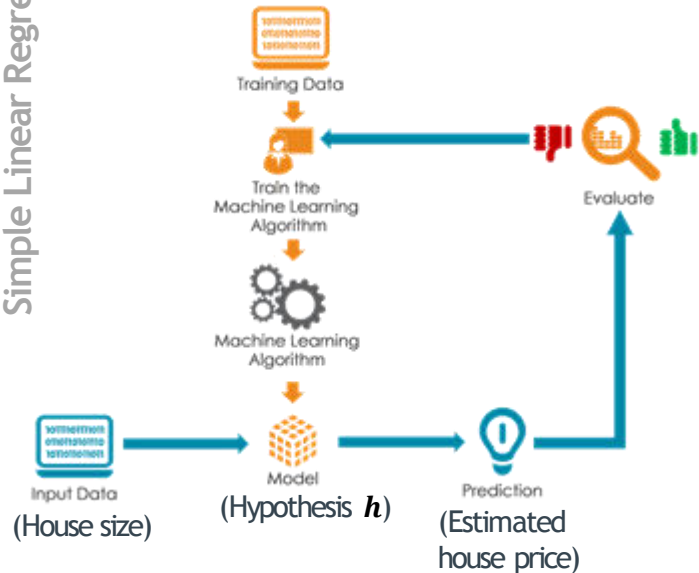
Size (sq. m) ( $x$ )	Price (1000s of R) ( $y$ )
460	6639
70	1681
155	3969
429	5095
...	...

Training set of house sizes and prices

- $x$ : Input variable; “features” used to make predictions e.g. 460 sq. metres
- $y$ : Output variable; “target” output e.g. R 6639000
- $m$ : Number of examples in the training set
- $(x, y)$ : A specific sample in the training set
- $(x^{(i)}, y^{(i)})$ : the  $i$ th training example e.g.  $x^{(1)} = 460$ ,  $x^{(2)} = 70$ ... ;  $y^{(1)} = 6639$ ,  $y^{(2)} = 1681$

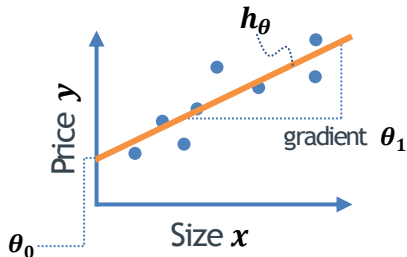


# Model Representation - Project Design



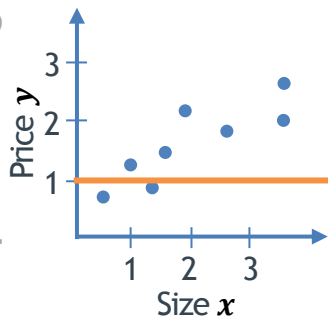
- $h$ : Mapping of  $x$  values (sizes) onto  $y$  values (prices)

What does  $h$  look like?

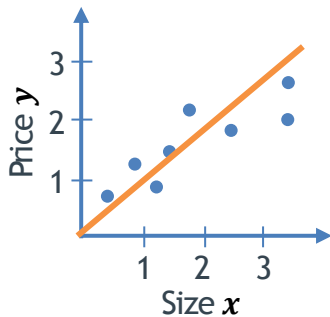


$$h_\theta(x) = \theta_0 + \theta_1 x$$

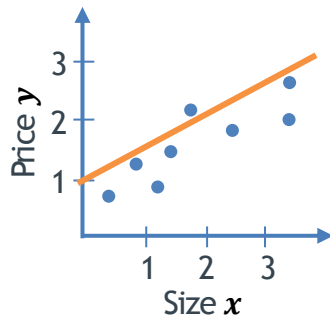
- Equation of a straight line with gradient  $\theta_1$  and y-intercept  $\theta_0$

Model Representation - Deciphering  $h_{\theta}(x)$ 

$$\theta_0 = 1; \theta_1 = 0$$

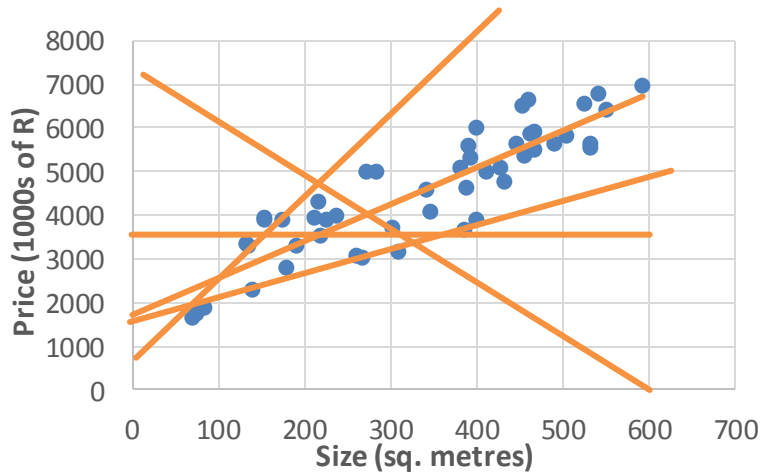


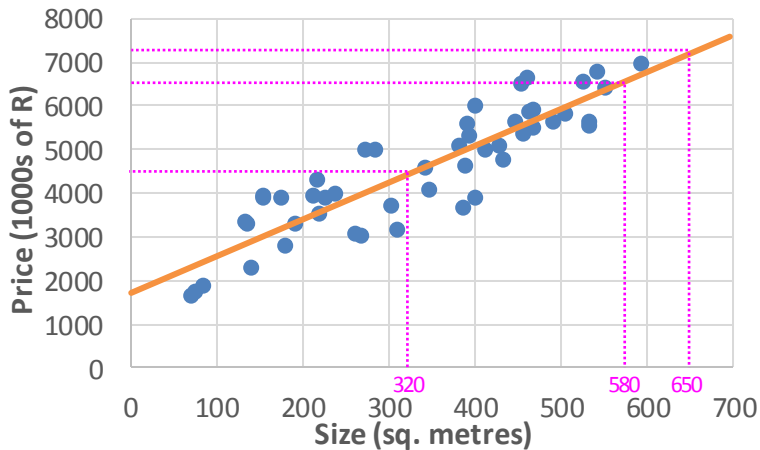
$$\theta_0 = 0; \theta_1 = 1$$



$$\theta_0 = 1; \theta_1 = 0.5$$

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

Model Representation - Deciphering  $h_{\theta}(x)$ 

Model Representation - Using  $h_{\theta}(x)$  to Make Predictions

Given the best fit  $\theta_0 = 1740$ ;  
 $\theta_1 = 8.4$ , predict:

- $h_{\theta}(320)$   
 $= 1740 + 8.4(320)$   
 $= 4428$
- $h_{\theta}(580)$   
 $= 1740 + 8.4(580)$   
 $= 6612$
- $h_{\theta}(650)$   
 $= 1740 + 8.4(650)$   
 $= 7200$

## Model Representation - Summary

- We need to fit a line ( $h_{\theta}$ ) onto the housing data (price versus size)
  - The line is determined by the parameters  $\theta_0$  and  $\theta_1$
  - Once we have an appropriate line, we can make price predictions on any unknown values in future
  - The values of  $\theta_0$  and  $\theta_1$  will determine how well the line fits the training data
  - Very important: the parameter values also determine how accurate future price predictions based on size may be
  - This is a simple / univariate linear regression problem
- 
- Golden question: How do we determine the best  $\theta_0$  and  $\theta_1$  to use??

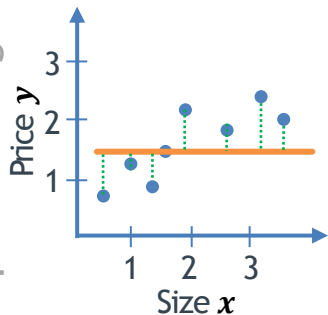
# Simple Linear Regression

## Cost Function

---

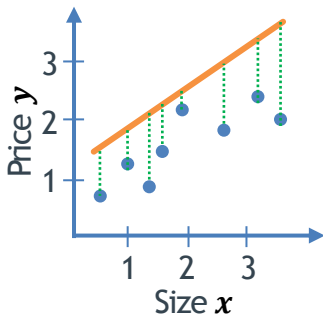
- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

## Cost Function - Formulation



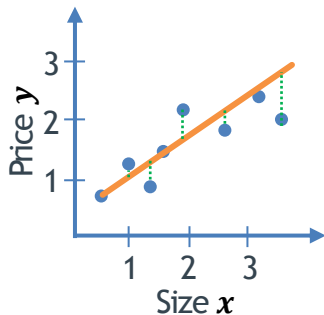
$$\theta_0 = 1.5; \theta_1 = 0$$

$$h_{\theta}(x) = 1.5$$



$$\theta_0 = 1; \theta_1 = 1$$

$$h_{\theta}(x) = 1 + x$$



$$\theta_0 = 0.5; \theta_1 = 0.5$$

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

- Given the points in the training set  $(x, y)$ , choose  $\theta_0$  and  $\theta_1$  such that the line's predicted prices ( $h_{\theta}(x)$ ) are "close" to the each actual price ( $y$ )

## Cost Function - Formulation

- Given the points in the training set  $(x, y)$ , choose  $\theta_0$  and  $\theta_1$  such that the line's predicted prices ( $h_{\theta}(x)$ ) are "close" to the each actual price ( $y$ )
- Compute the distance between  $h_{\theta}(x)$  and  $y$  for all the data points:

$$h_{\theta}(x^{(i)}) - y^{(i)} \text{ for all } (i) \text{ from } 1 \text{ to } m$$

- Add/sum the distances up

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



## Cost Function - Formulation

- Similar to adding/summing the **squares** of the distances up

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

- Similar to adding/summing the squares of the distances up and dividing by the number of points ***m*** to get the total **average** square distance

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

- Similar to multiplying by a half to get **half** the total average square distance

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

The smaller this is, the better the line fits  
 The larger this is, the worse the line fits  
 This is the “cost function”

## Cost Function - Formulation

- Similar to multiplying by a half to get half the total average square distance

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

The smaller this is, the better the line fits  
The larger this is, the worse the line fits  
This is the “cost function”

- Remember that:

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

- So actually:

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

- So the cost function (how well the line fits the data) depends directly on  $\theta_0$  and  $\theta_1$ 
  - In mathematical terms: it is a “function” of  $\theta_0$  and  $\theta_1$

## Cost Function - Formulation

- So the cost function (how well the line fits the data) depends directly on  $\theta_0$  and  $\theta_1$ 
  - In mathematical terms: it is a “function” of  $\theta_0$  and  $\theta_1$
  - Conventionally, the function is denoted as  $J$
  - Also known as the squared error function

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

or

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

## Cost Function - Formulation

- To choose the best  $\theta_0$  and  $\theta_1$  for a given training set, find  $\theta_0$  and  $\theta_1$  for which the cost function  $J$  has the smallest value i.e. minimize  $J$
- Mathematically:

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

or

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

or

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

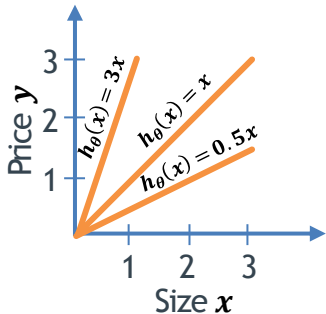
## Cost Function - Deciphering the Cost Function $J$

- To understand what minimizing the cost function means, let's assume:

$$\theta_0 = 0$$

- $h_{\theta}(x)$  is a line that passes through the origin (0,0), with changing gradients so:

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

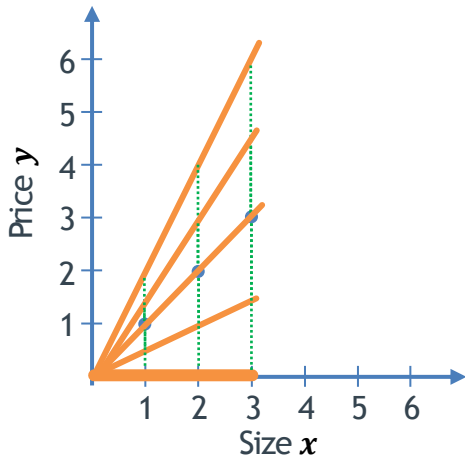


- Therefore the cost function now only depends on  $\theta_1$ :

$$\underset{\theta_1}{\text{minimize}} \quad J(\theta_1)$$

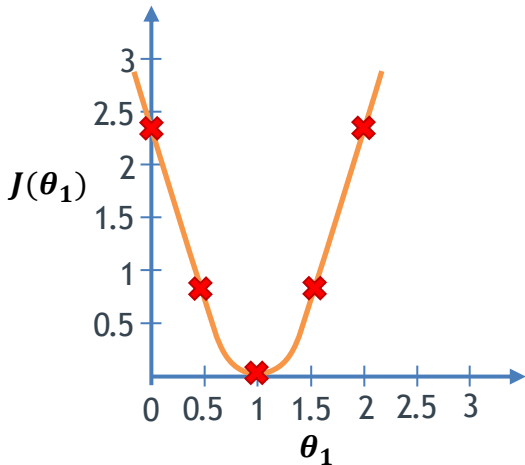
$$\underset{\theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\underset{\theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

Cost Function - Deciphering the Cost Function  $J$ 

$$\theta_1 = 0.5$$

$$h_{\theta}(x) = 0.5x \quad (\theta_0 = 0)$$



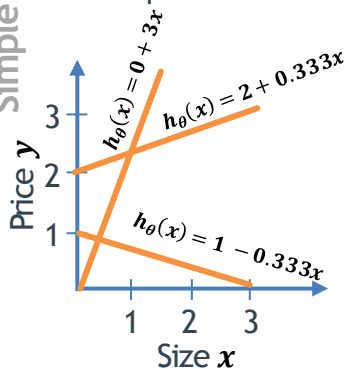
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

## Cost Function - Deciphering the Cost Function $J$ With Two Params

- To understand what minimizing the cost function means, now let's take the original cost function :

$$J(\theta_0, \theta_1)$$

- An arbitrary line with any gradient and y-intercept so:



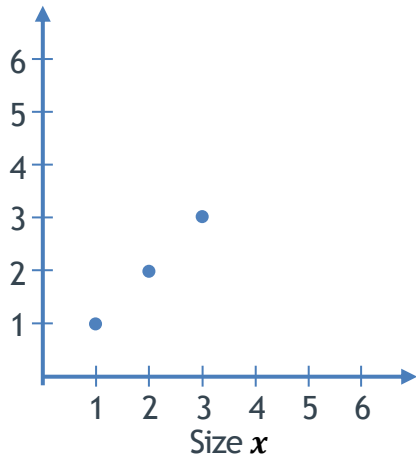
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- The cost function depends on both  $\theta_0$  and  $\theta_1$ :

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

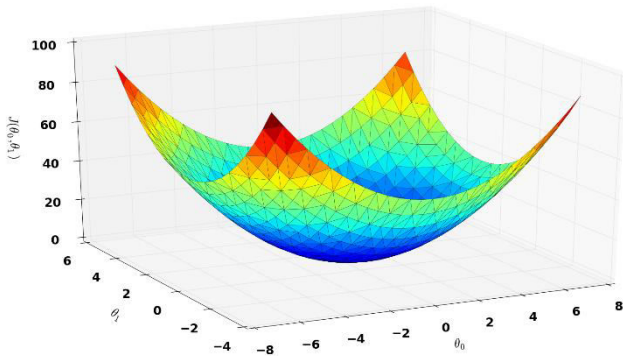
$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

Cost Function - Deciphering the Cost Function  $J$  With Two Params

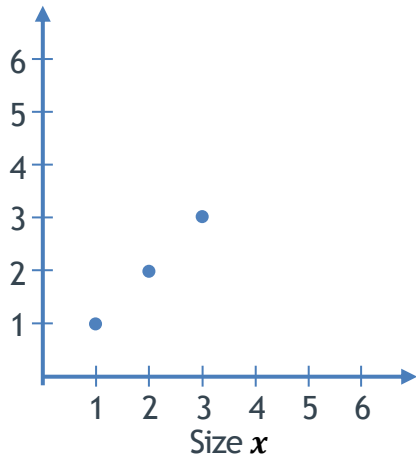
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- The cost function is a 3D bowl-shaped surface:



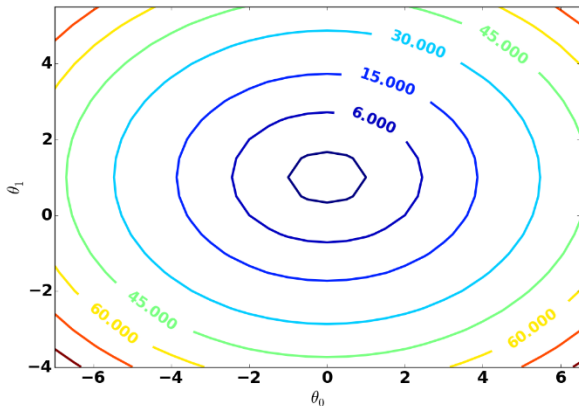
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$



Cost Function - Deciphering the Cost Function  $J$  With Two Params

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

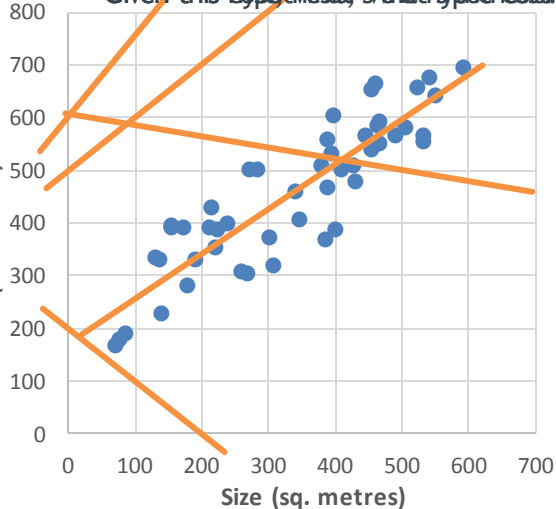
- The cost function as a contour plot:



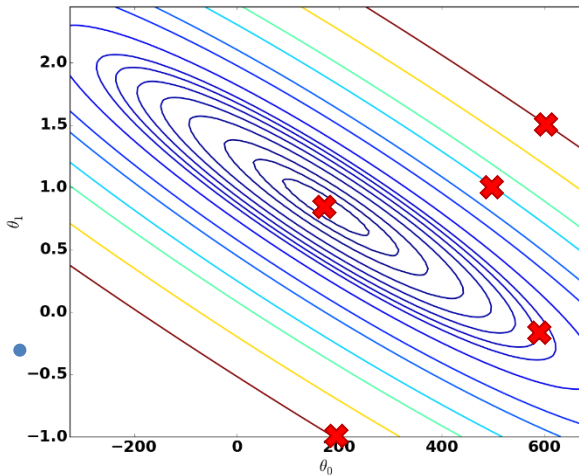
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

# Cost Function - Deciphering the Cost Function $J$ With Two Params

- Given this hypothesis, what hypothesis?



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



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

# Gradient Descent

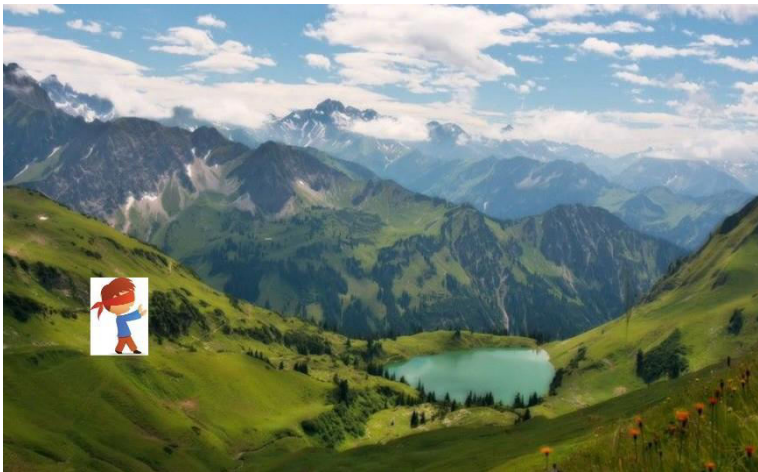
## Formulation

---

- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

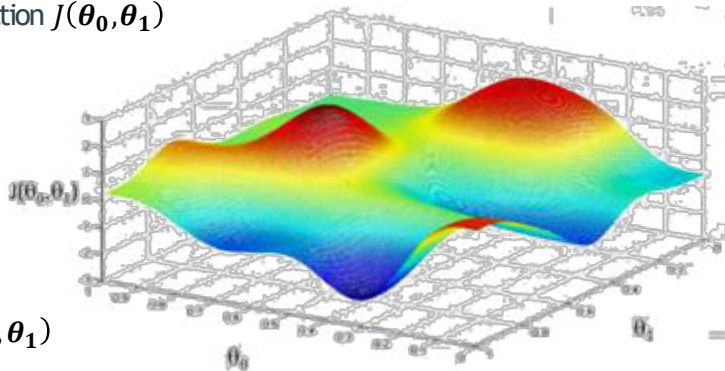
## Formulation - Introduction

- How would you go about reaching the minimum (water) below?



## Formulation - Minimizing the Cost Function

- Given some arbitrary cost function  $J(\theta_0, \theta_1)$



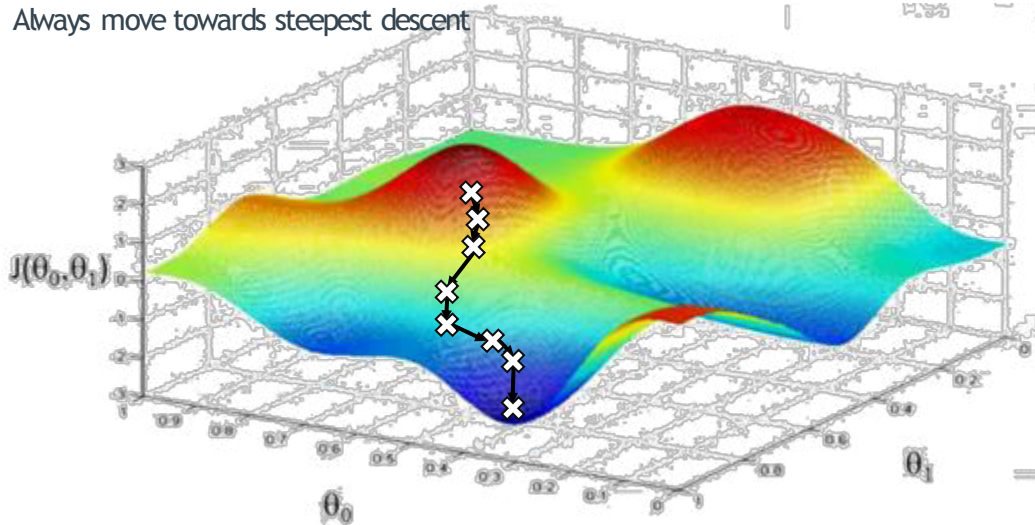
- We want to achieve:

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

- Strategy:
  - Initialize  $\theta_0, \theta_1$  to some random values
  - Continuously make updates to  $\theta_0$  and  $\theta_1$  in the direction of “descent”
  - Until the minimum is reached

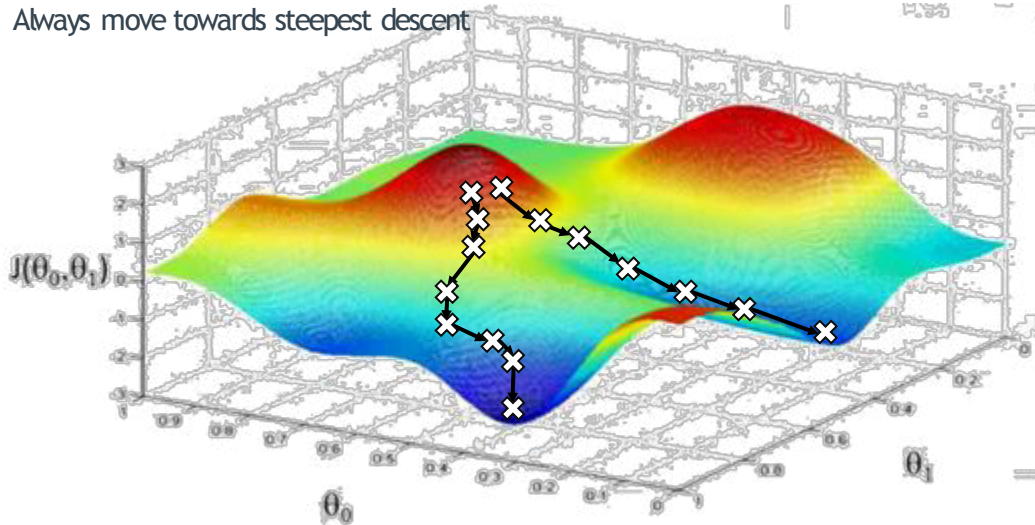
## Formulation - Minimizing the Cost Function

- Starting at  $\theta_0 = 0.7$  and  $\theta_1 = 0.4$
- Always move towards steepest descent



## Formulation - Minimizing the Cost Function

- Starting at  $\theta_0 = 0.67$  and  $\theta_1 = 0.43$
- Always move towards steepest descent



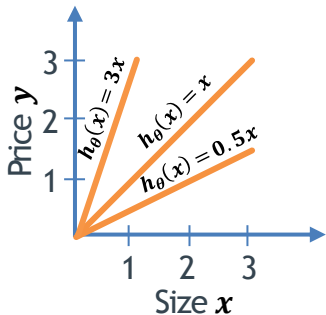
## Formulation - Gradient Descent With One Parameter

- To show you how to formulate gradient descent, let's again assume:

$$\theta_0 = 0$$

- $h_{\theta}(x)$  is a line that passes through the origin (0,0), with changing gradients so:

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



- Therefore the cost function now only depends on  $\theta_1$ :

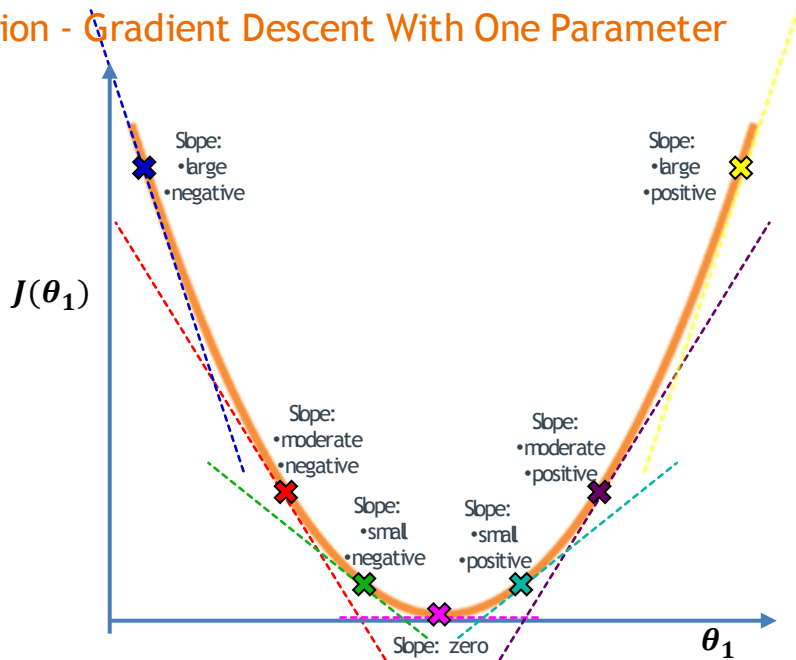
$$\underset{\theta_1}{\text{minimize}} \quad J(\theta_1)$$

$$\underset{\theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

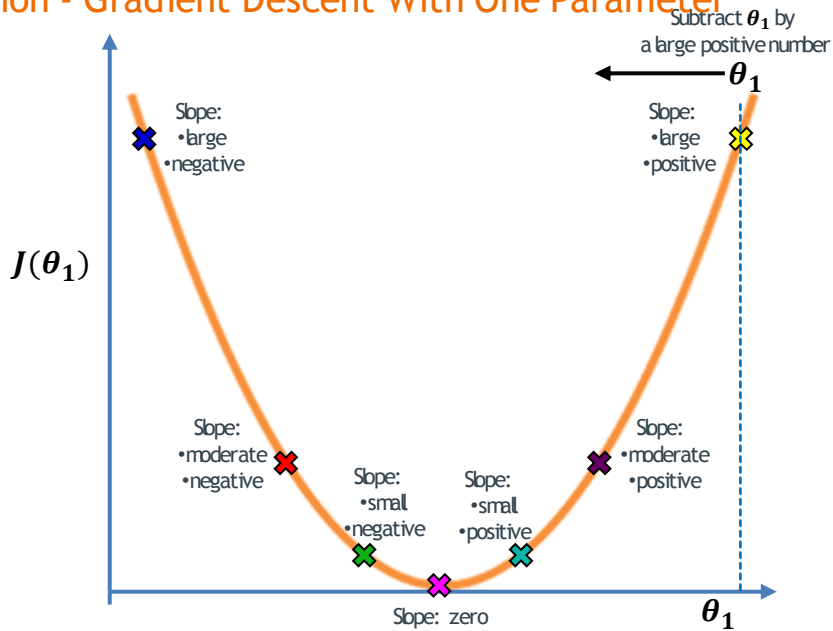
$$\underset{\theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$



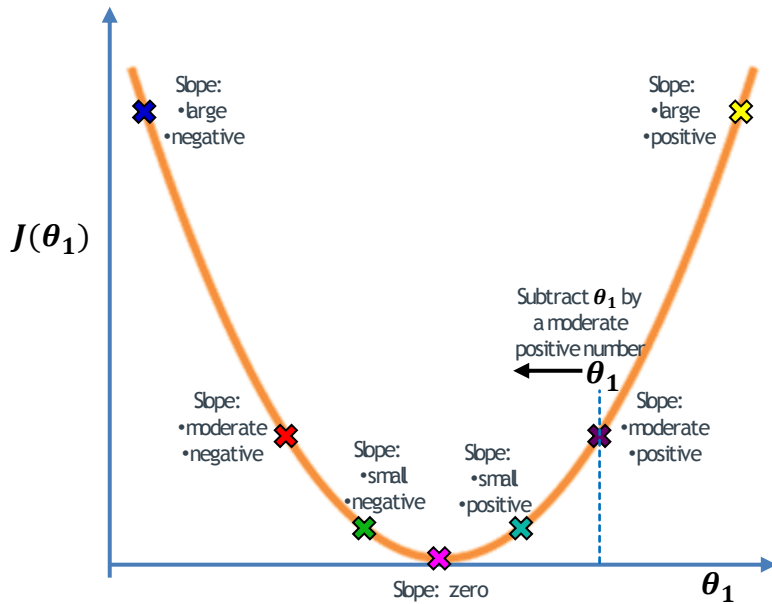
## Formulation - Gradient Descent With One Parameter



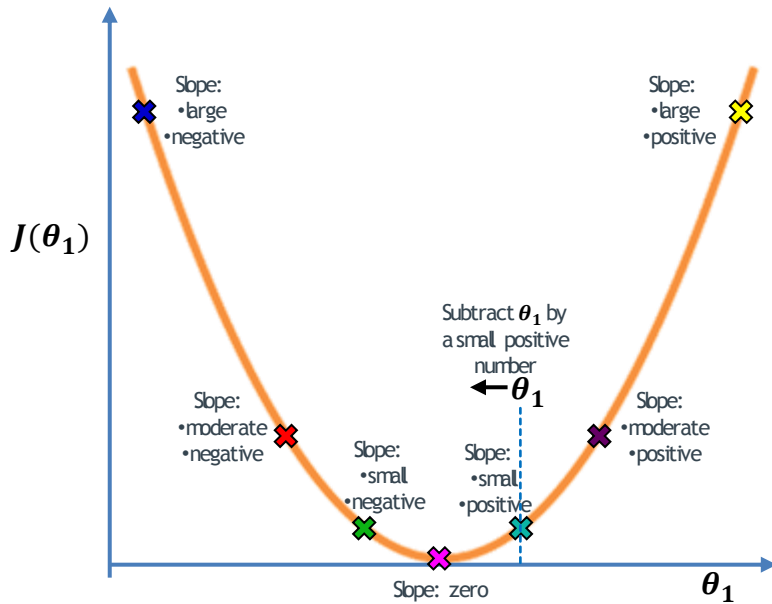
## Formulation - Gradient Descent With One Parameter



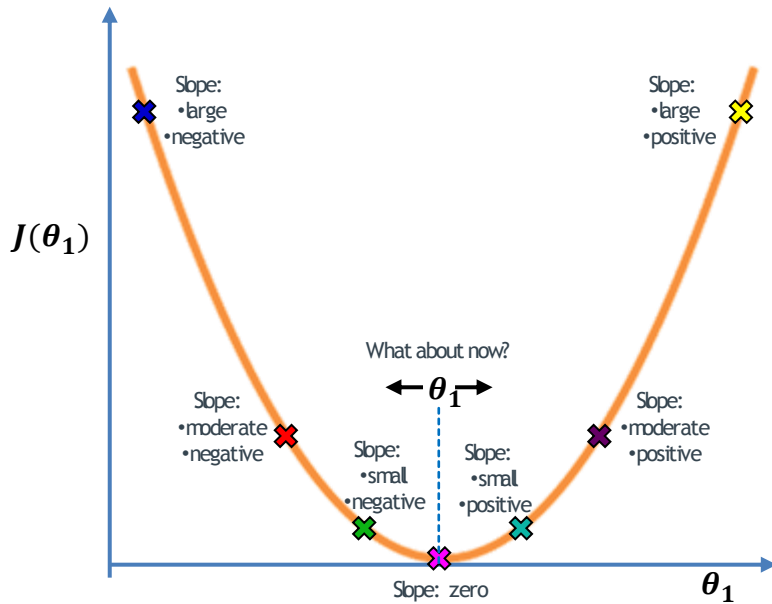
## Formulation - Gradient Descent With One Parameter



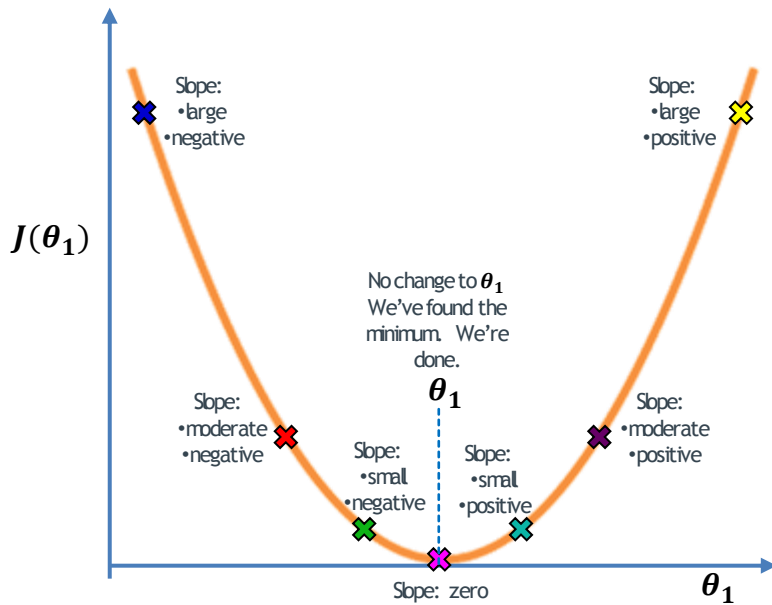
## Formulation - Gradient Descent With One Parameter



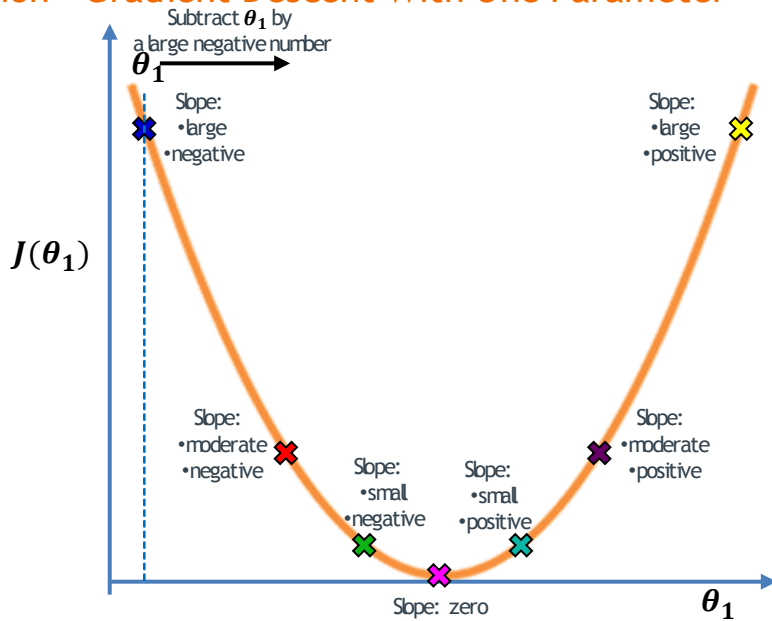
## Formulation - Gradient Descent With One Parameter



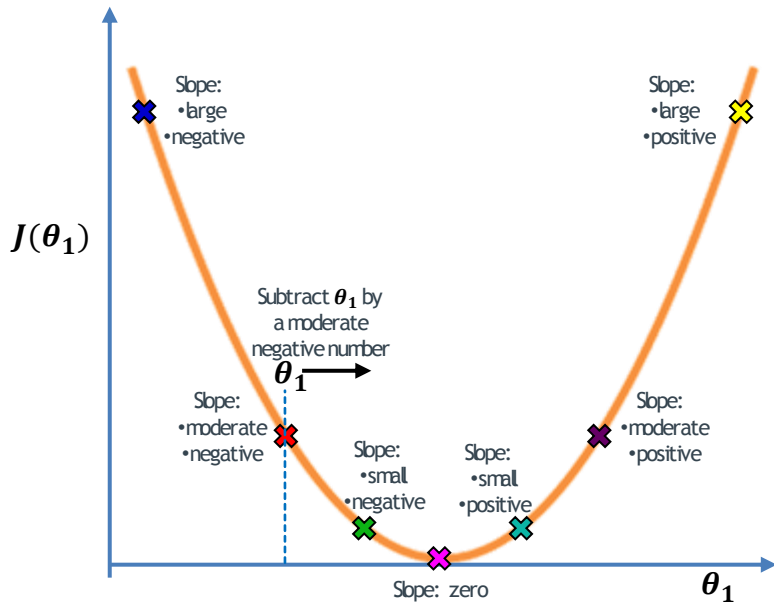
## Formulation - Gradient Descent With One Parameter



## Formulation - Gradient Descent With One Parameter

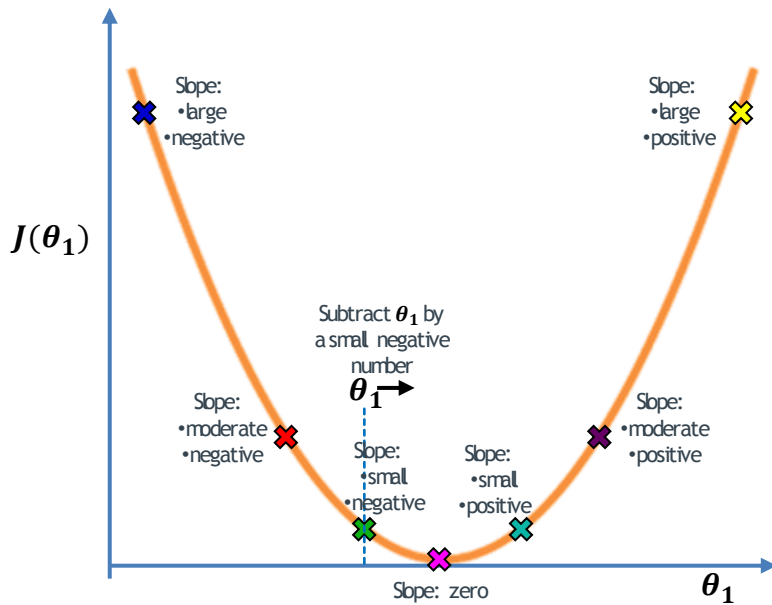


## Formulation - Gradient Descent With One Parameter





## Formulation - Gradient Descent With One Parameter



# Gradient Descent Algorithm

---

- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

## Algorithm - Gradient Descent With One Parameter

- For a cost function  $J(\theta_1)$  :

Repeat until convergence:

Update  $\theta_1$ :

$$\theta_1 \leftarrow \theta_1 - \frac{d}{d\theta_1} J(\theta_1)$$

Learning rate

Slope

## Algorithm - Deciphering Learning Rate

- With no learning rate:

$$\theta_1 \leftarrow \theta_1 - \frac{d}{d\theta_1} J(\theta_1)$$

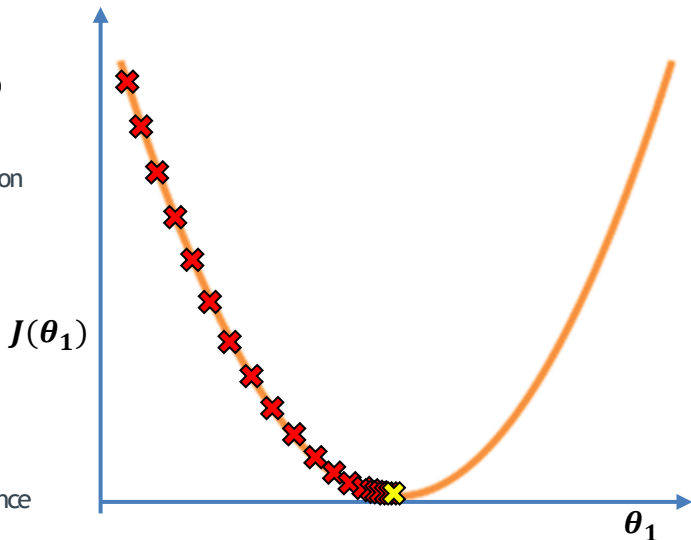
- No learning rate: no controls on convergence

- Could be too slow
- Could be too fast

- Learning rate is very cost-function specific

- In this case:

- Too slow
- 20 steps to convergence



## Algorithm - Deciphering Learning Rate

- With appropriate learning rate e.g.

$\alpha = 1.5$  :

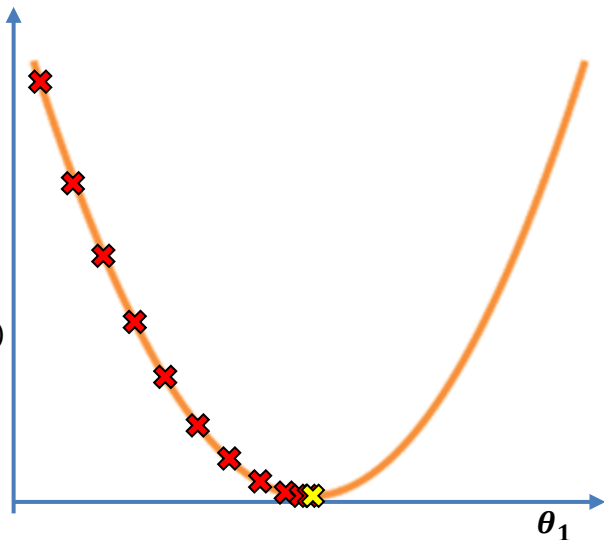
- fast(er) convergence

$$\theta_1 \leftarrow \theta_1 - 1.5 \frac{d}{d\theta_1} J(\theta_1)$$

- 12 (vs 20) steps to convergence (for this specific cost function)

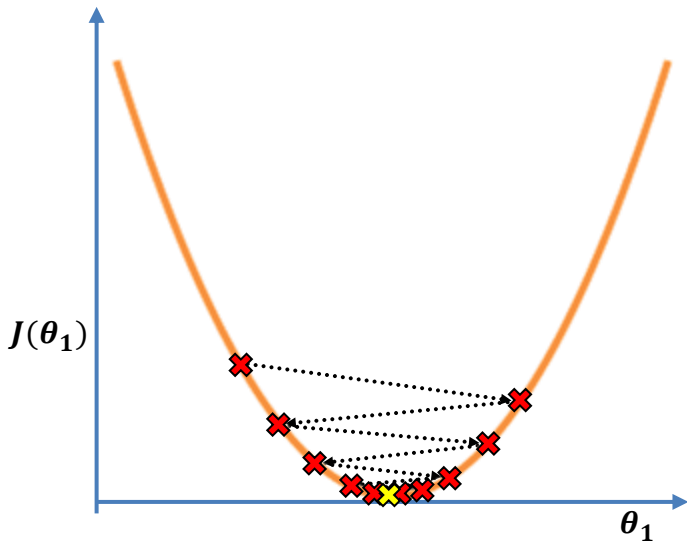
$J(\theta_1)$

- Note: learning rate is cost-function specific
  - Higher learning rate doesn't necessarily mean faster convergence



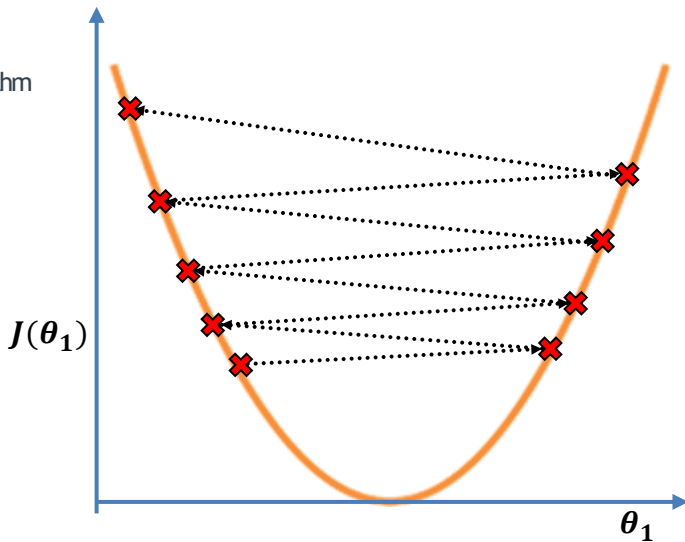
## Algorithm - Deciphering Learning Rate

- With learning rate too large:
- **First possibility:** convergence becomes very slow



## Algorithm - Deciphering Learning Rate

- With learning rate too large:
- **Second possibility:** the algorithm diverges (doesn't converge)



## Algorithm - Gradient Descent With Two Parameters

- For a cost function  $J(\theta_0, \theta_1)$ :

Repeat until convergence:

Update  $\theta_0, \theta_1$  simultaneously:

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

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

Slope of  $J$  in the  $\theta_1$  direction

Repeat until convergence:

$$\text{tmptheta0} = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{tmptheta1} = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 = \text{tmptheta0}$$

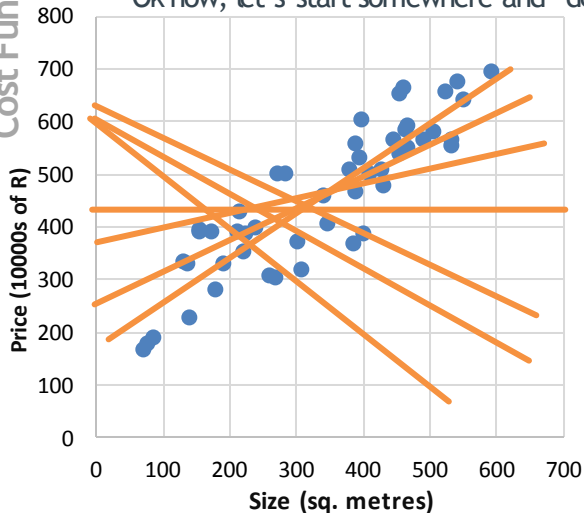
$$\theta_1 = \text{tmptheta1}$$

Slope of  $J$  in the  $\theta_0$  direction

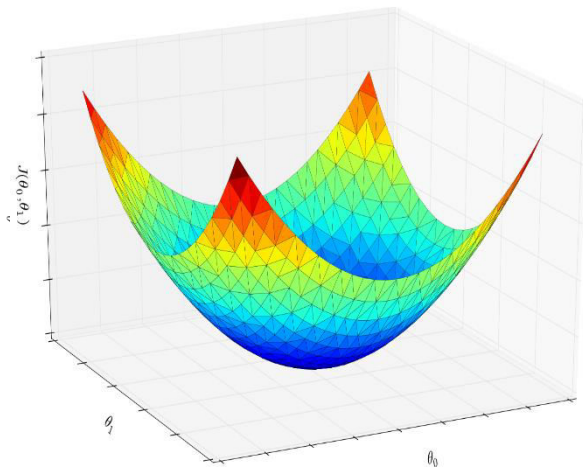


## Algorithm - Illustration

- Ok now, let's start somewhere and "descend" the cost function down to the minimum



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



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

## Algorithm - Gradient Descent With Generic Cost Function

- For a cost function  $J(\theta_0, \dots, \theta_n)$  with parameters  $\theta_0, \dots, \theta_n$ :

Repeat until convergence:

Update all  $\theta_j$  simultaneously:

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

Repeat until convergence:

$$\text{tmptheta0} = \text{theta0} - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \dots, \theta_n)$$

...

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

...

$$\text{tmpthetan} = \text{thetan} - \alpha \frac{\partial}{\partial \theta_n} J(\theta_0, \dots, \theta_n)$$

$$\text{theta0} = \text{tmptheta0}$$

...

$$\text{thetaj} = \text{tmpthetaj}$$

...

$$\text{thetan} = \text{tmpthetan}$$

# Gradient Descent

## Gradient Descent Applied to Linear Regression

- Simple (Univariate) Linear Regression
  - Model Representation
  - Cost Function
- Gradient Descent
  - Formulation
  - Algorithm
  - Application to Linear Regression

## Gradient Descent Applied to Linear Regression

- For a cost function  $J(\theta_0, \theta_1)$ :

Repeat until convergence:

Update  $\theta_0, \theta_1$  simultaneously:

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

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

?

?

- Linear Regression:

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

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

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = ?$$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = ?$$

## Gradient Descent Applied to Linear Regression

- For any parameter  $\theta_j$ :

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

- For  $\theta_0$ :

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

- For  $\theta_1$ :

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

Repeat until convergence:

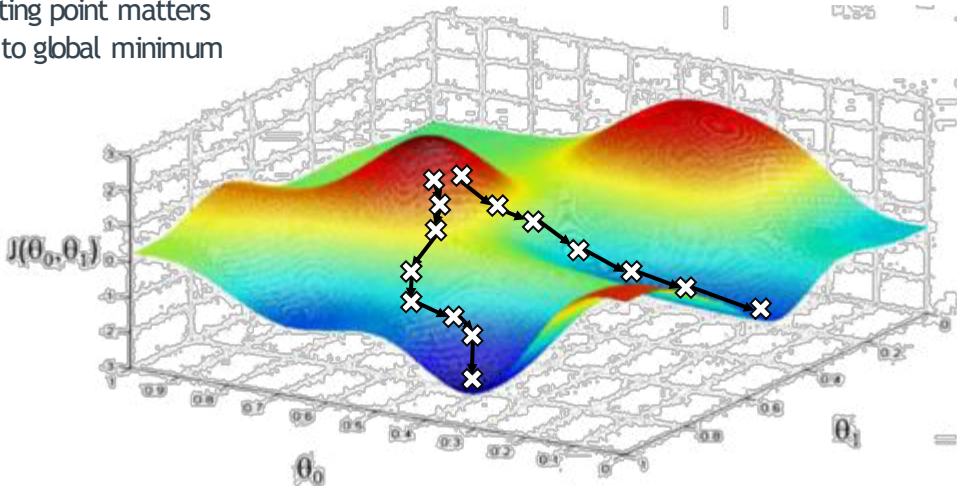
Update  $\theta_0, \theta_1$  simultaneously:

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

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

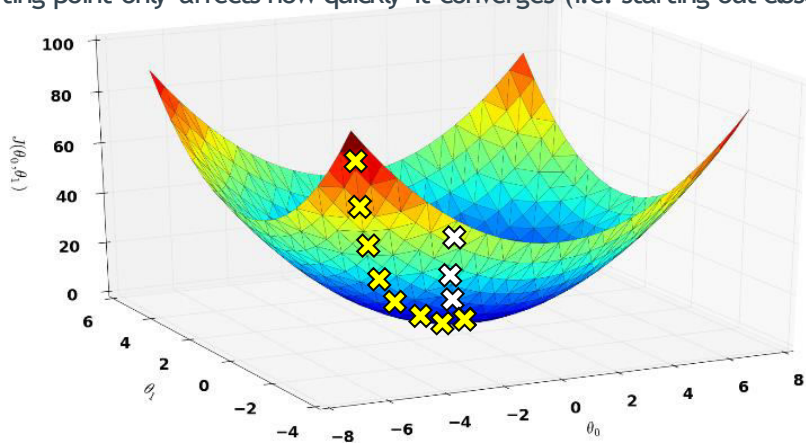
## Gradient Descent Applied to Linear Regression

- Gradient descent works for minimizing any generic cost function
- Cost function may have many different local minima
  - Starting point matters
- May not get to global minimum



## Gradient Descent Applied to Linear Regression

- The cost function of linear regression is “convex” i.e. it has only one minimum - global minimum
- With the right learning rate  $\alpha$ , it always converges
- Starting point only affects how quickly it converges (i.e. starting out close VS far)



## Note on Gradient Descent

- This kind of gradient descent is actually called “Batch Gradient Descent”
  - One update uses all the training samples
- Other types of gradient descent (not covered in this course - feel free to look them up):
  - Mini-Batch Gradient Descent
  - Stochastic Gradient Descent



## Content - Part 2

- Linear Regression With Multiple Variables
  - Model Representation
  - Gradient Descent With Multiple Variables
- Practical Tips On Implementing Linear Regression
  - Diagnosing the Learning Rate
  - Feature Scaling

# Linear Regression With Multiple Variables

## Model Representation

---

- Linear Regression With Multiple Variables
  - Model Representation
  - Gradient Descent With Multiple Variables
- Practical Tips On Implementing Linear Regression
  - Diagnosing the Learning Rate
  - Feature Scaling

## Previously (With Simple Linear Regression)

$x$	$y$
Size (sq. m)	Price (1000s of R)
460	6639
70	1681
155	3969
429	5095
...	...

Training set of house sizes and prices

- $x$ : Input variable; “features” used to make predictions e.g. 460 sq. metres
- $y$ : Output variable; “target” output e.g. R 6639000
- $m$ : Number of examples in the training set
- $(x, y)$ : A specific samples in the training set
- $(x^{(i)}, y^{(i)})$ : the  $i$ th training example e.g.  $x^{(1)}=460$ ,  $x^{(2)}=70$ ... ;  $y^{(1)}=6639$ ,  $y^{(2)}=1681$

## Model Representation - Notation With Multiple Variables

$x$	$x_2$	$x_3$	$x_4$	$y$
Size (sq. m)	No. of Rooms	Age (years)	No. of Garages	Price (1000s of R) ( $y$ )
460	4	12	2	6639
70	1	5	0	1681
155	3	8	2	3969
429	6	10	3	5095
...	...	...	...	...

Training set of house features and prices

- $y$ : Output variable; “target” output e.g. R6639000
- $m$ : Number of examples in the training set
- $(x^{(i)}, y^{(i)})$ : the  $i$ th training example
- $y^{(i)}$ : the  $i$ th target e.g.  $y^{(1)} = 6639$ ,  $y^{(2)} = 1681$  etc.

## Model Representation - Notation With Multiple Variables

$x_1$	$x_2$	$x_3$	$x_4$	$y$
Size (sq. m)	No. of Rooms	Age (years)	No. of Garages	Price (1000s of R) ( $y$ )
460	4	12	2	6639
70	1	5	0	1681
155	3	8	2	3969
429	6	10	3	5095
...	...	...	...	...

Training set of house features and prices

- $n$ : Number of features; in the above case  $n = 4$
- $x_j$ : the  $j$ th input variable (column); e.g.  $x_1$  is the size,  $x_2$  is the no. of rooms etc.
- $x^{(i)}$ : is now a set of  $n$  values i.e a vector
- $x_j^{(i)}$ : is the  $j$ th feature of the  $i$ th example i.e. a number

## Model Representation - Notation With Multiple Variables

$x_1$	$x_2$	$x_3$	$x_4$	$y$
Size (sq. m)	No. of Rooms	Age (years)	No. of Garages	Price (1000s of R) ( $y$ )
460	4	12	2	6639
70	1	5	0	1681
155	3	8	2	3969
429	6	10	3	5095
...	...	...	...	...

Training set of house features and prices

$$X = \begin{bmatrix} 460 & 4 & 12 & 2 \\ 70 & 1 & 5 & 0 \\ 155 & 3 & 8 & 2 \\ 429 & 6 & 10 & 3 \\ \dots & \dots & \dots & \dots \end{bmatrix}$$


$$y = \begin{bmatrix} 6639 \\ 1681 \\ 3969 \\ 5095 \\ \dots \end{bmatrix}$$

## Model Representation - Hypothesis

- The previous hypothesis (with simple linear regression) with only one variable  $x$ :

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

- With more variables:  $x_1, \dots, x_4$ :


$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

- E.g.

$$h_{\theta}(x) = 174 + 0.84x_1 + 1.8x_2 - 1.5x_3 + 3x_4$$


## Model Representation - Notation With Multiple Variables

- In the generic case with  $n$  features:

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

- $n$ -dimensional hypothesis
  - We can't plot/visualize it for  $n > 2$
- For consistency in the notation we add a 0<sup>th</sup> feature  $x_0 = 1$  so that:

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

  
= 1



## Model Representation - Vectorized Representation

- In the generic case with  $n$  features:

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

Everything in this column = 1

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & \cdots & x_n^{(3)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

- Given the two matrix-representations above,  $h_{\theta}(x)$  can be computed using the matrix operation:

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

## Model Representation - Vectorized Representation

- In the specific case with 4 features (in the example):

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$X = \begin{bmatrix} 1 & 460 & 4 & 12 & 2 \\ 1 & 70 & 1 & 5 & 0 \\ 1 & 155 & 3 & 8 & 2 \\ 1 & 429 & 6 & 10 & 3 \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

- Given the two matrix-representations above,  $h_{\theta}(x)$  can be computed using the matrix operation:

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

# Linear Regression With Multiple Variables

## Gradient Descent With Multiple Variables

- Linear Regression With Multiple Variables
  - Model Representation
  - Gradient Descent With Multiple Variables
- Practical Tips On Implementing Linear Regression
  - Diagnosing the Learning Rate
  - Feature Scaling

## Gradient Descent With Multiple Variables

- In the generic case with  $n$  features/variables:

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

$$X = \begin{bmatrix} \overset{=1}{x_0^{(1)}} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \dots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

- Cost function:

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

- Generic gradient descent update equation (as seen before):

Repeat until convergence:

Update all  $\theta_j$  simultaneously:

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

## Gradient Descent With Multiple Variables

- For any parameter  $\theta_j$ :

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

- Previously with only one variable  $x$ :  $\theta_0, \theta_1$ :

Repeat until convergence:

Update  $\theta_0, \theta_1$  simultaneously:

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

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

- With variables  $n > 1$ :  $\theta_0, \theta_1, \dots, \theta_n$ :

Repeat until convergence:

Update all  $\theta$  simultaneously:

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 \leftarrow \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x_1^{(i)}$$

...

$$\theta_n \leftarrow \theta_n - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x_n^{(i)}$$

## Gradient Descent With Multiple Variables - Vectorized Representation

- Given:

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & \dots & x_n^{(3)} \\ \dots & \dots & \dots & \dots & \dots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \dots \\ y^{(m)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \dots \\ \theta_n \end{bmatrix}$$

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

- The updates to all theta can be made simultaneously using:

$$\theta \leftarrow \theta - \frac{\alpha}{m} [X^T (X\theta - y)]$$

## Gradient Descent With Multiple Variables - Vectorized Representation

- As by the way: given:

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & \dots & x_n^{(3)} \\ \dots & \dots & \dots & \dots & \dots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \dots \\ y^{(m)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \dots \\ \theta_n \end{bmatrix}$$

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

- The cost  $J(\theta)$  can be computed as:

$$J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$$

## Practical Tips

### Diagnosing the Learning Rate

---

- Linear Regression With Multiple Variables
  - Model Representation
  - Gradient Descent With Multiple Variables
- Practical Tips On Implementing Linear Regression
  - Diagnosing the Learning Rate
  - Feature Scaling

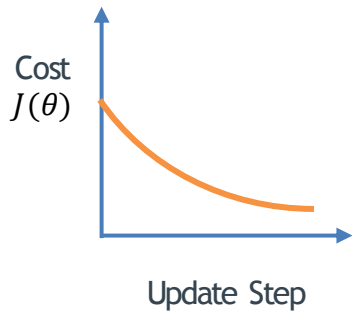


## Diagnosing the Learning Rate

- Learning rate is specific to each problem
- Too small: slow convergence
- Too large:
  - Slow convergence
  - Divergence (no convergence)
- Helpful tip: Proven that if the learning rate  $\alpha$  is small enough, the cost  $J$  will reduce in every iteration of gradient descent
- Possible strategy: keep track of the cost at each update
  - Plot a graph of the cost at each update
  - Use the plot to determine if  $\alpha$  is:
    - Too small
    - Too large
  - Use a series of plots to pick the best  $\alpha$

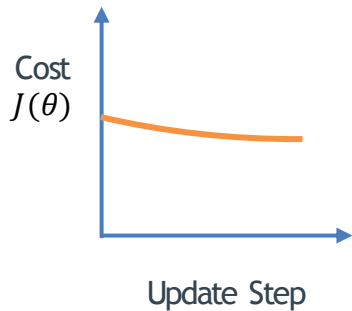
## Diagnosing the Learning Rate

- If  $\alpha$  is just right:



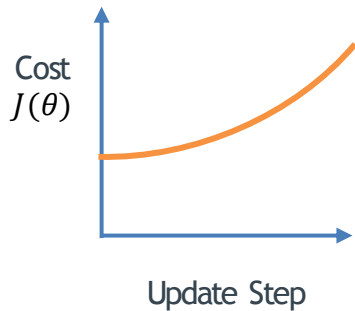
## Diagnosing the Learning Rate

- If  $\alpha$  is too small:



## Diagnosing the Learning Rate

- If  $\alpha$  is too large :



## Diagnosing the Learning Rate - Strategy

- Try a range of values for  $\alpha$  on an increasing scale e.g.  
0.001   0.003   0.01   0.03   0.1   0.3   1
- Plot the curve of  $J(\theta)$  over the data for a number of iterations for every  $\alpha$
- Pick the  $\alpha$  value that converges the fastest

## Practical Tips

### Feature Scaling

---

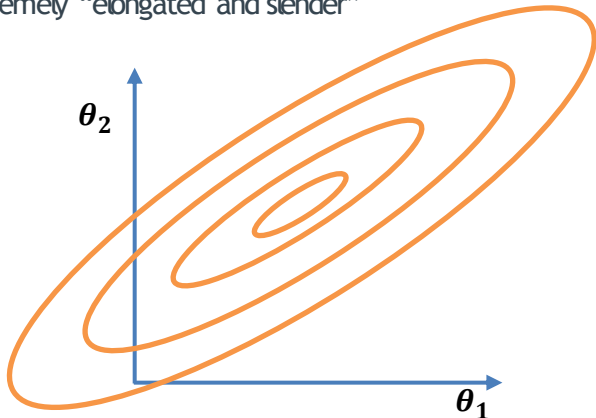
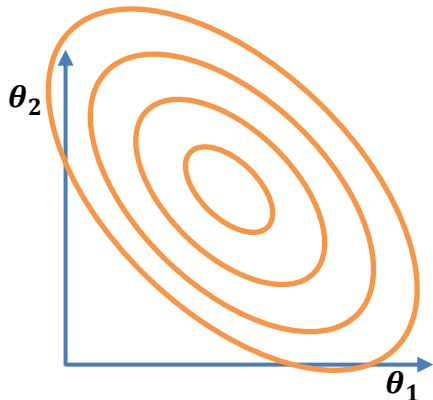
- Linear Regression With Multiple Variables
  - Model Representation
  - Gradient Descent With Multiple Variables
- Practical Tips On Implementing Linear Regression
  - Diagnosing the Learning Rate
  - Feature Scaling

## Feature Scaling - Concept

- Concept: Having features on a similar scale helps gradient descent converge faster
- Opposite: Having features on very different scales slows gradient descent down
- Concept: scale features to similar scales
- Features may be on different scales
- E.g.:
  - House size: 0 - 5000 sq. metres                      i.e. 0 - 5000
  - No of bedrooms: 1 - 12 rooms                      i.e. 0 - 12
  - Age: 0 - 200 years                      i.e. 0 - 200

## Feature Scaling - Intuition

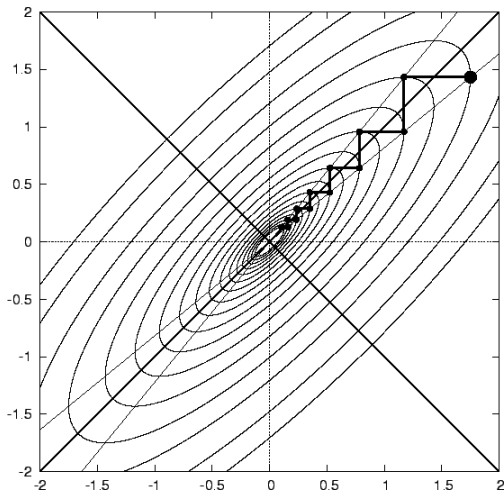
- One example E.g.:
  - No of bedrooms: 0 - 12
  - House size: 0 - 5000
- More extreme example E.g.:
  - Feature 1: 0 - 5
  - Feature 2: 0 - 20000000
- This causes the cost function to be extremely “elongated and slender”





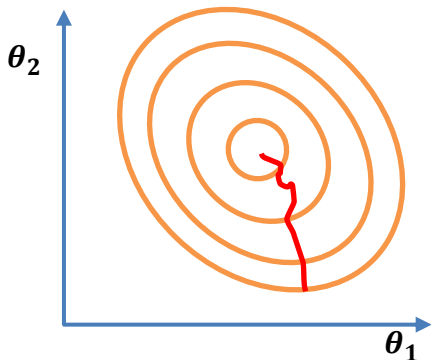
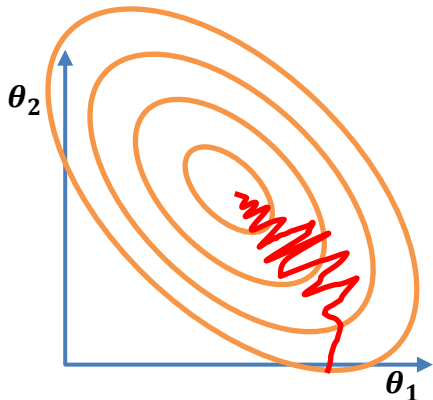
## Feature Scaling - Intuition

- This can cause gradient descent to take a long time to reach the minimum
- May also meander around a lot



## Feature Scaling - Intuition

- One example E.g.:
  - No of bedrooms: 0 - 12
  - House size: 0 - 5000
- Scale the features:
$$x_1 = \frac{\text{Bedrooms}}{12} \qquad x_2 = \frac{\text{Size}}{5000}$$
- With a smaller and more circular cost function, gradient descent can converge faster



## Feature Scaling - Intuition

- Applying the simple division:

$$x_1 = \frac{\text{Bedrooms}}{12}$$

$$x_2 = \frac{\text{Size}}{5000}$$

- Has the effect of re-scaling  $x_1$  and  $x_2$ .
  - $x_1$ : re-scaled to 0 - 1
  - $x_2$ : re-scaled to 0 - 1

## Feature Scaling - Concept

- $x_0 = 1$  in all cases, so no scaling is needed
- For all other features: aim is to re-scale every feature to **approximately**  $-1 \leq x_1 \leq 1$  range
- Approximately means all of the following are **okay**:
  - $0 \leq x_1 \leq 4$
  - $-2 \leq x_2 \leq 1$
  - $-0.5 \leq x_3 \leq 3$
- The following are examples of **not okay**:
  - $-80 \leq x_4 \leq 90$
  - $-0.0001 \leq x_5 \leq 0.0004$

## Feature Scaling - Method

- Applied to every feature except  $x_0$
- Replace  $x_j$  with  $x_j - \mu_j$ 
  - where  $\mu_j$  is the average of all values of  $x_j$
  - This results in a 0 mean for this feature
- Divide by  $S_j$  which can be either of:
  - The range of values i.e.  $\max(x_j) - \min(x_j)$  OR
  - The standard deviation of values of  $x_j$
- Results in the approximate range:
  - $-0.5 \leq x_j \leq 0.5$

$$x_j \leftarrow \frac{x_j - \mu_j}{S_j}$$

$$X = \begin{bmatrix} 1 & 460 & 2 \\ 1 & 70 & 0 \\ 1 & 155 & 2 \\ 1 & 429 & 3 \end{bmatrix}$$

$$X_{scaled} = \begin{bmatrix} 1 & 0.46 & 0.08 \\ 1 & -0.53 & 0.58 \\ 1 & -0.32 & 0.08 \\ 1 & 0.38 & 0.42 \end{bmatrix}$$

## Content - Part 3

- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

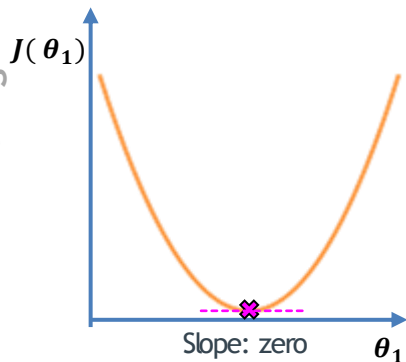
# Linear Regression Using the Normal Equation

---

- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

## Intuition

- A second analytical alternative to solving for parameters  $\theta$ 
  - (First one was gradient descent - algorithmic solution to  $\theta$ )
- For a very simple hypothesis  $h_{\theta}(x) = \theta_1 x_1$  the cost function looks like the below:



- Alternative to getting  $\theta$ : solve for  $\theta$  in:

$$\frac{d}{d \theta_1} J(\theta_1) = 0$$

- The only unknown is  $\theta_1$ . Solvable
- With many more features, concept is the same:
  - Solve for  $\theta$  in:

$$\frac{d}{d \theta_1} J(\theta_0, \theta_1, \dots, \theta_n) = 0$$



## Normal Equation Definition

- Given a feature matrix  $X$ , and the corresponding output matrix  $y$
- The following equation solves for  $\theta$  that best fits the data:

$$\theta = (X^T X)^{-1} X^T y$$

- Where:
  - $X^T$  is the transpose of  $X$
  - $(X^T X)^{-1}$  is the inverse of  $(X^T X)$
  - $X$  includes feature  $x_0$  which is a column of 1s

- E.g.

$$X = \begin{bmatrix} 1 & \text{Size} & \text{\#Rms} & \text{Age} & \text{\#Grgs} \\ 1 & 460 & 4 & 12 & 2 \\ 1 & 70 & 1 & 5 & 0 \\ 1 & 155 & 3 & 8 & 2 \\ 1 & 429 & 6 & 10 & 3 \end{bmatrix}$$

$$y = \begin{bmatrix} \text{Price} \\ 6639 \\ 1681 \\ 3969 \\ 5095 \end{bmatrix}$$

## Normal Equation - Notes

- To get the inverse In SciPy:
- First import `scipy.linalg` as `spla`, then either
  - `spla.inv(X)` OR:
  - `spla.pinv(X)`
  - Rather use `spla.pinv(X)`
- To get matrix transpose in SciPy:
  - Easiest way: `X.T`
- No need for scaling at all
  - It's not descending a function

## Normal Equation VS Gradient Descent

### Normal Equation

#### Pros

- Doesn't need any iterations - finds the solution in one go immediately
- No need to choose  $\alpha$
- Preferable over gradient descent if possible

#### Cons

- Can be impossibly slow to compute  $(X^T X)^{-1}$  when  $n$  gets large ( $n > 10000$ )
  - $O(n)^3$  time
- $(X^T X)^{-1}$  can be non-invertible
  - E.g. when  $n \gg m$

### Gradient Descent

#### Pros

- Works well even for very large  $n$

#### Cons

- Need to choose  $\alpha$
- Could take many iterations to find best  $\theta$  fit

# Complex (Non-Linear) Features

---

- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

## Concept

- Up to now, we've used simple linear features e.g. Size, No of Rooms etc.
- Possible to use linear regression to learn more complex features:
  - Combinations of features
  - Higher-order features

## Combinations Of Features

- Assume we've got 4 features: size, no of rooms, height and width of properties
- Up to now: Linear combination of features

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \text{Size} + \theta_2 x_2 \text{Rooms} + \theta_3 x_3 \text{Height} + \theta_4 \theta_4 \text{Width}$$

- Given insight into the problem, you can create new features using these basic features
  - If you think they are more “telling” / “feature-ful”

## Combinations Of Features

- E.g. Divide Size by the No of Rooms to get the “Size-to-rooms-ratio” (STRR)
  - (Maybe) cramming more rooms in a smaller area (smaller STRR) means lower quality i.e. price
  - Conversely, (maybe) placing fewer rooms in a larger area (larger STRR) means higher quality i.e. price
- We can replace the two features  $x_1$  and  $x_2$  with just a single STRR feature:

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \text{STRR} + \theta_2 \cdot \text{Rooms} + \theta_3 \cdot \text{Height} + \theta_4 \cdot \text{Width}$$

- Where STRR is now feature  $x_1 = (\text{Size} \cdot \text{Rooms})$

## Combinations Of Features

- E.g. 2: Multiply the frontal-width of the house by the height of the house to get the area of the frontal display (frontal-display area - FDA)
  - (Maybe) Having a larger frontal display area leads to a better first impression → higher demand → higher price
  - Conversely (maybe) a smaller frontal display area leads to a let-down first impression → lower demand → lower price
- We can replace the two features  $x_2$  and  $x_3$  with just a single FDA feature:

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \text{STRR} + \theta_2 \cdot \text{FDA} + \theta_3 \cdot \text{Width}$$

- Where FDA is now feature  $x_2 = (\text{Height} \cdot \text{Width})$



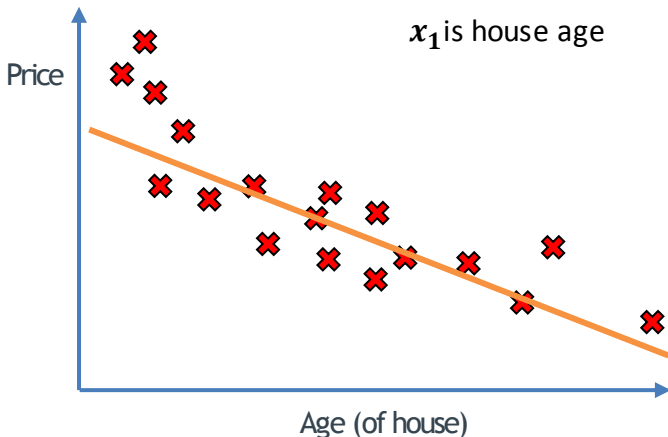


## Higher-Order Features

- Closely related to previous idea

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

$x_1$  is house age



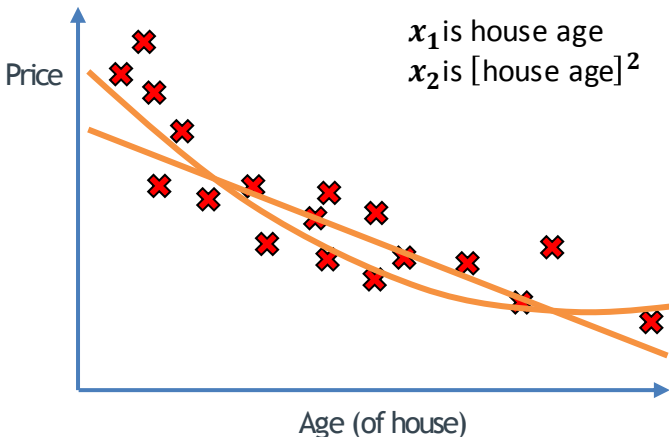
## Higher-Order Features

- Maybe a second-order polynomial (quadratic function) makes more sense i.e. fits this data better

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

$x_1$  is house age

$x_2$  is  $[\text{house age}]^2$



## Higher-Order Features

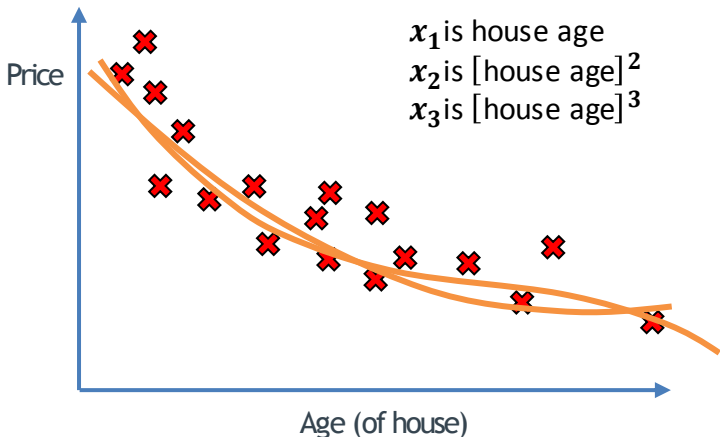
- Or maybe a third-order polynomial makes more sense i.e. fits this data better

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$x_1$  is house age

$x_2$  is [house age]<sup>2</sup>

$x_3$  is [house age]<sup>3</sup>

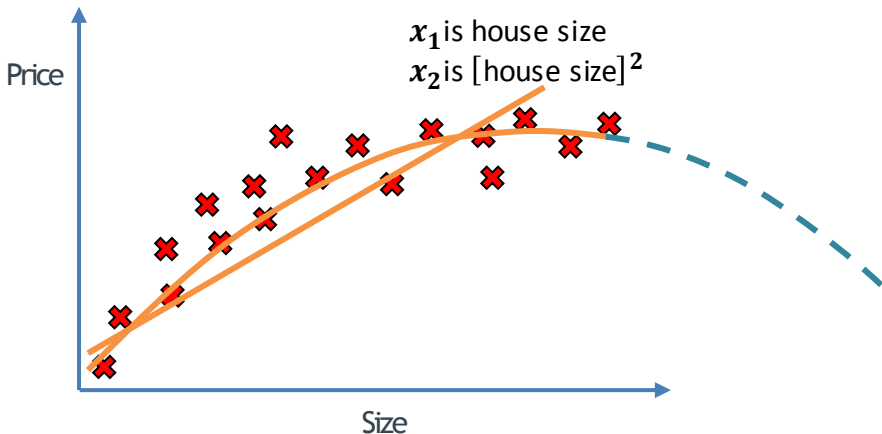


## Higher-Order Features

- Another example: Price vs Size
- Maybe a quadratic function fits better

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

$x_1$  is house size  
 $x_2$  is [house size]<sup>2</sup>



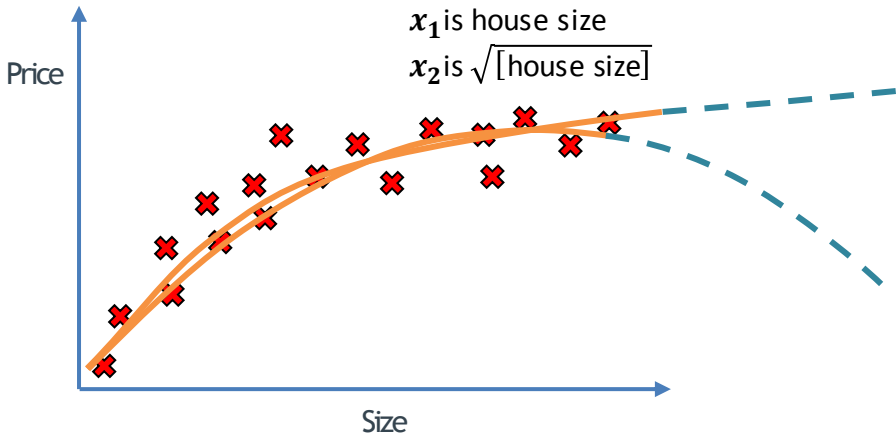
## Higher-Order Features

- Another example: Price vs Size
- Maybe a square root function actually makes more sense

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

$x_1$  is house size

$x_2$  is  $\sqrt{[\text{house size}]}$



## Complex Features - Important Points

- The list of feature possibilities (combinations / higher-order features) are endless
  - Provides flexibility VS Challenging to choose good features
- Extremely important to scale features if you're using gradient descent with complex features:
  - If  $x_1$  has the range  $0 \leq x_1 \leq 1000$
  - Then:
    - $x_1^2$  has the range  $0 \leq x_2 \leq 1,000,000$
    - $x_1^3$  has the range  $0 \leq x_3 \leq 1,000,000,000$
    - $\sqrt{x_1}$  has the range  $0 \leq x_3 \leq 32$
- When trying to predict on a set of test data, produce the same combination / higher-order features
  - Format has to be exactly the same as those used in training

# Model Evaluation

## Testing a Model

---

- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

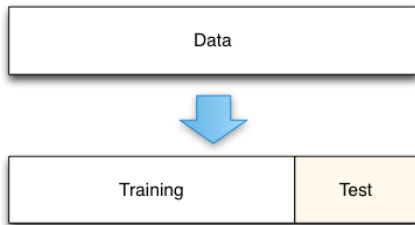
## Testing a Model

- Given a training set: train a predictive model
- Question: how well does this predict?
- One strategy: Test on the training data:
  - Model is tailor-made for the training data
  - Will (likely) give a good result
  - Not a good indicator of accuracy
- Another strategy: Collect more data
  - Very impractical
  - Expensive



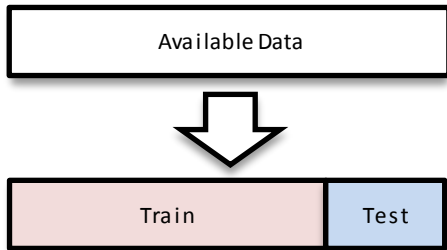
## Testing a Model

- Good strategy: Divide up the data that you have
  - Training portion (between 50% and 90%)
  - Testing portion (between 50% and 10%)
  - VERY important to randomize first
- Do not use testing data in training at all
  - It is “unseen” to the model
- Once model is trained, pass testing data to model to make predictions
  - Then apply a metric to determine performance



## Testing a Model

- Good strategy: Divide up the data that you have
  - Training portion (between 50% and 90%)
  - Testing portion (between 50% and 10%)
  - VERY important to randomize first



$X =$	460	4	12	2	$y =$	6639	Training portion
	70	1	5	0		1681	
	155	3	8	2		3969	
	429	6	10	3		5095	
	313	3	4	1		3412	
	56	2	1	0		1234	
	434	1	9	3		2334	Testing portion
	123	7	3	1		2786	

## Testing a Model - Metrics

- Three metrics:

- Mean Absolute Error:

$$\mathbf{MAE} = \frac{1}{m} \sum_{i=1}^m |y^{(i)} - y_{pred}^{(i)}|$$

$0 - \infty$   
The smaller the better

- Mean Squared Error:

$$\mathbf{MSE} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - y_{pred}^{(i)})^2$$

$0 - \infty$   
The smaller the better

- Chi-squared error  $R^2$ :

[Placeholder - Look it up]

$-\infty - 1$   
 $-\infty$  - Very poor fit  
1 - Perfect fit

# Evaluating a Model

## Comparing Models

---

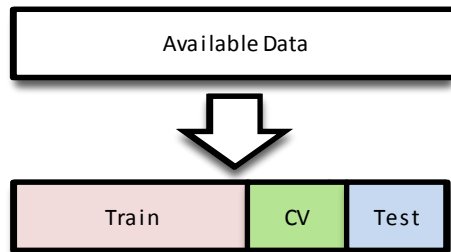
- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

## Comparing Models

- The list of features to try are endless
- Given several models/options to compare e.g.
  - quadratic features VS square root features
  - Using Size + No Of Rooms VS Size + No Of Garages VS No of Rooms + No Of Garages
- How to compare them?
- Possible strategy:
  - Use the same strategy used to test a single model i.e.
    - Divide available data into training and testing sets
    - Train all models on training sets
    - Test all models on testing sets
    - Conclude which is the best
  - Technically not statistically valid; more valid strategy:
    - After picking the “winner” we need to test it on a final piece of data to quote its performance

## Comparing Models

- Valid strategy: Divide up the data that you have
  - Training portion (between 50% and 80%)
  - Cross-validation portion (between 25% and 10%)
  - Testing portion (between 25% and 10%)
  - VERY important to randomize first



$X =$	460	4	12	2	$y =$	6639	Training portion
	70	1	5	0		1681	
	155	3	8	2		3969	
	429	6	10	3		5095	
	313	3	4	1		3412	CV portion
	56	2	1	0		1234	
	434	1	9	3		2334	Testing portion
	123	7	3	1		2786	

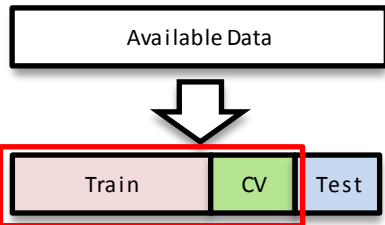
## Evaluating a Model

### Getting More Mileage Out of Your Data

- Linear Regression Using The Normal Equation
- Complex (Non-Linear) Features
- Model Evaluation
  - Testing a Model
  - Comparing Models
  - Getting More Mileage Out of Your Data

## Getting More Mileage Out of Your Data

- Given a limited-size data set, possible to use
- techniques to stretch out “Training” and “CV” portions
- out further:
  - k-Fold Cross-validation
  - Leave-one-out Cross-validation



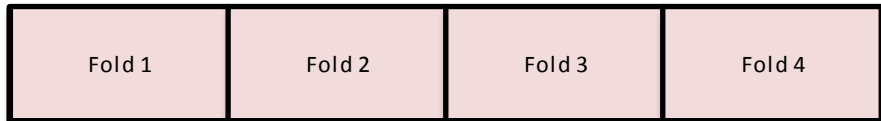
$X =$	460	4	12	2	$y =$	6639	Training portion
	70	1	5	0		1681	
	155	3	8	2		3969	
	429	6	10	3		5095	
	313	3	4	1		3412	CV portion
	56	2	1	0		1234	
	434	1	9	3		2334	Testing portion
	123	7	3	1		2786	

We can stretch out these parts



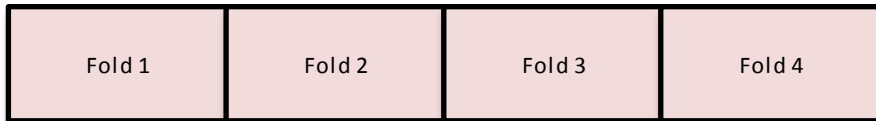
## $k$ -Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below



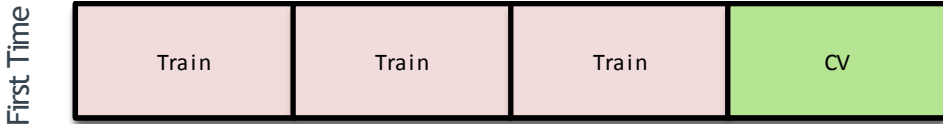
## $k$ -Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets



## k-Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets
    - Train on all the data marked as “Train”; Test on the data marked “CV”
    - Get an evaluation score for the model



## k-Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets
    - Train on all the data marked as “Train”; Test on the data marked “CV”
    - Get an evaluation score for the model

Second Time



## k-Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets
    - Train on all the data marked as “Train”; Test on the data marked “CV”
    - Get an evaluation score for the model



## k-Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets
    - Train on all the data marked as “Train”; Test on the data marked “CV”
    - Get an evaluation score for the model

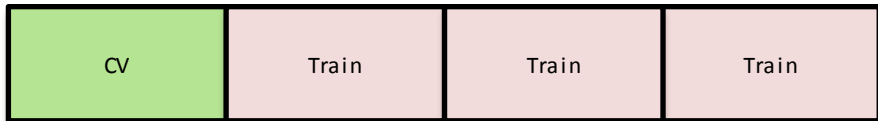
Fourth Time



## k-Fold Cross Validation

- Given the Training + CV parts below
  - Treat them as one big “Training Set”
  - Divide it up into  $k$ -Folds
    - E.g. 4 folds below
  - Train  $k$  times:
    - Each time take 1 fold as the CV set and the remaining folds as training sets
    - Train on all the data marked as “Train”; Test on the data marked “CV”
    - Get an evaluation score for the model

Fourth Time



ood average

- Arguably statistically valid
- Note: Once done, test on the test set (which was set aside)

**THE END**  
Of Linear Regression

---