

REGRESSION ANALYSIS : LINEAR

BY – MAUJAMA FIRDAUS & TULIKA SAHA

MACHINE LEARNING

- It is the science of getting computer to learn without being explicitly programmed.
- Machine learning is an area of artificial intelligence concerned with the development of techniques which allow computers to "learn".
- More specifically, machine learning is a method for creating computer programs by the analysis of data sets.

APPLICATIONS

- Search Engines like Google, Bing etc.
- Facebook photo tagging application.
- Self Customizing Programs and many more.

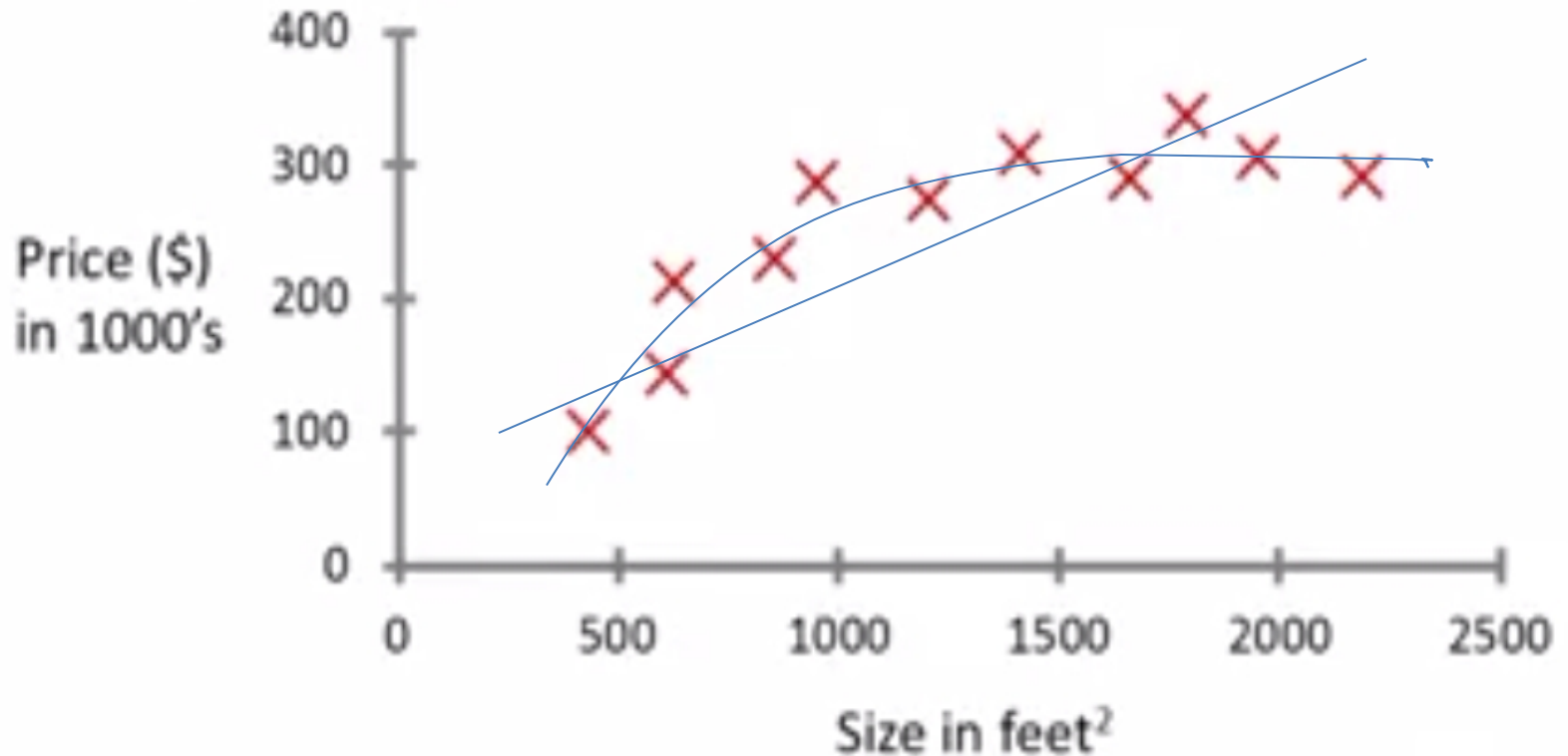
SUPERVISED LEARNING

- Given the “right answer” for each of our examples in the training set.
- The task of the algorithm is to find many more such right answers for new examples.
- Consists of two problems-
 - Regression Problem
 - Classification Problem

REGRESSION PROBLEM

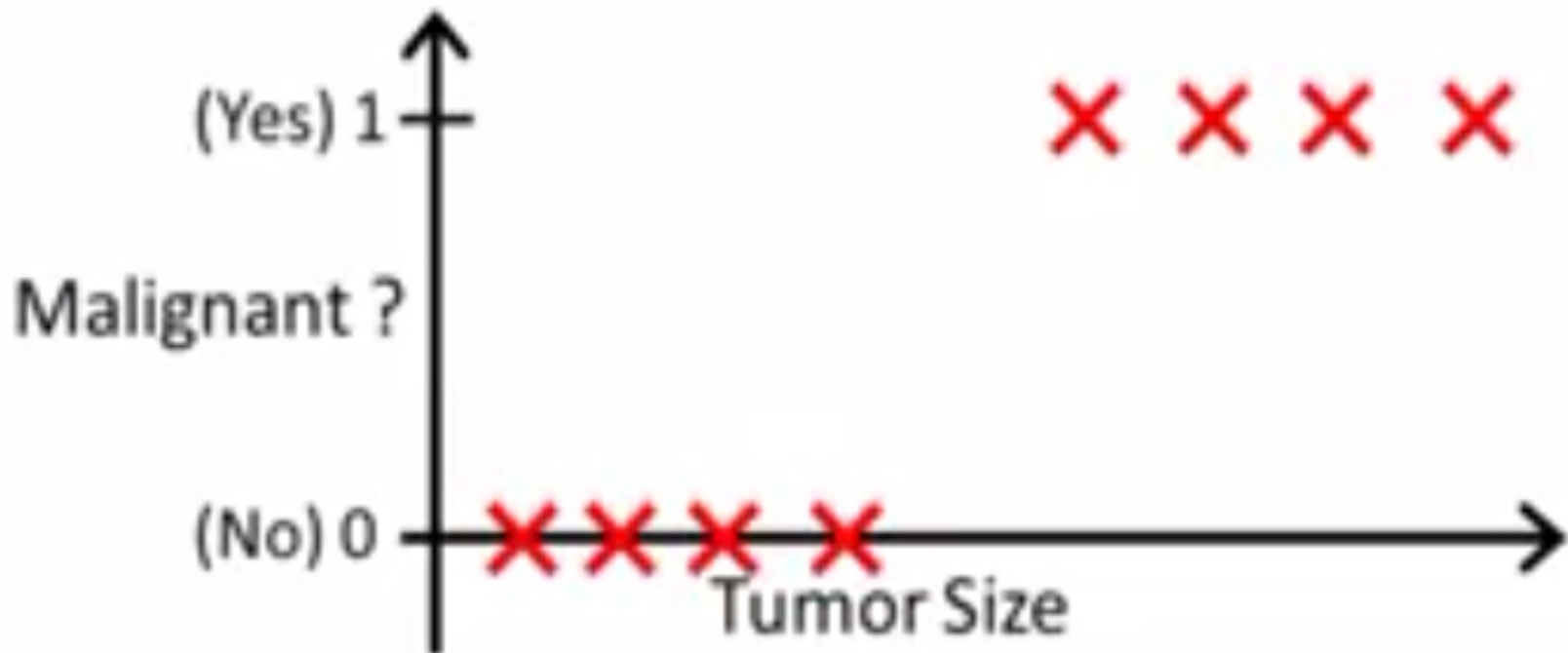
- The term “regression” refers to the fact that we are trying to predict a continuous-valued output.

Housing price prediction.



CLASSIFICATION PROBLEM

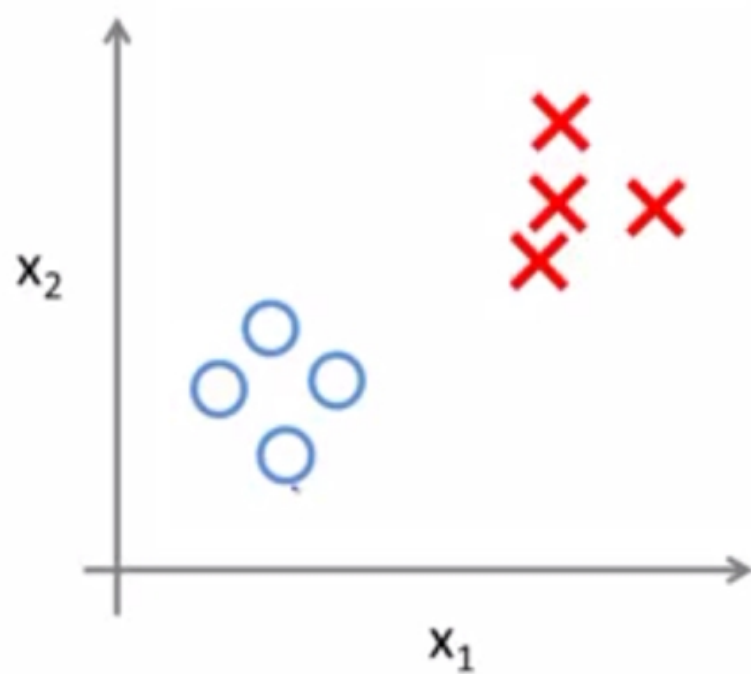
- It refers to the fact that we are trying to predict discrete valued outputs.



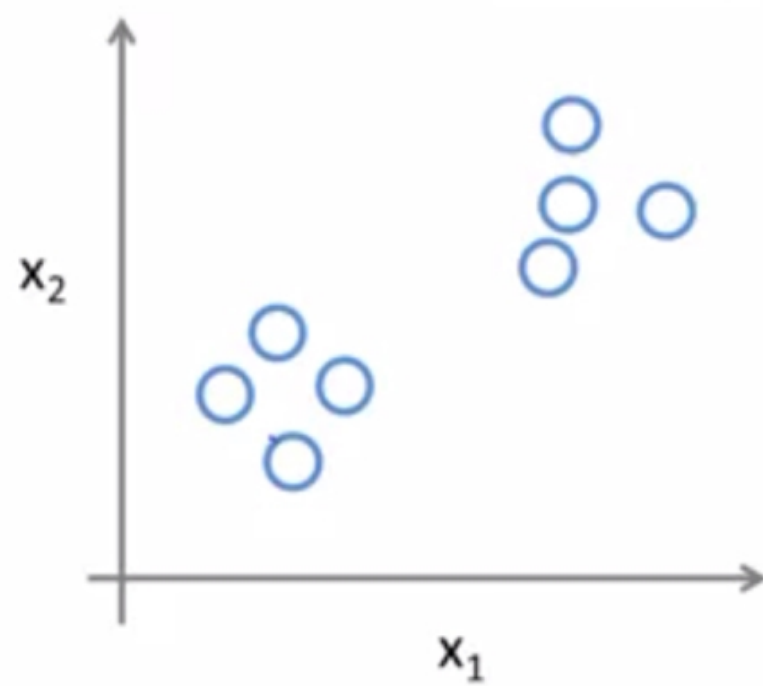
UNSUPERVISED LEARNING

- Given a set of data or examples in the training set without prior information of what the data is all about.
- The task of the algorithm is to find the structure among the data.
- Clustering Problem is an example of Unsupervised Learning.

Supervised Learning



Unsupervised Learning



LINEAR REGRESSION

Why are we looking at this algorithm ?

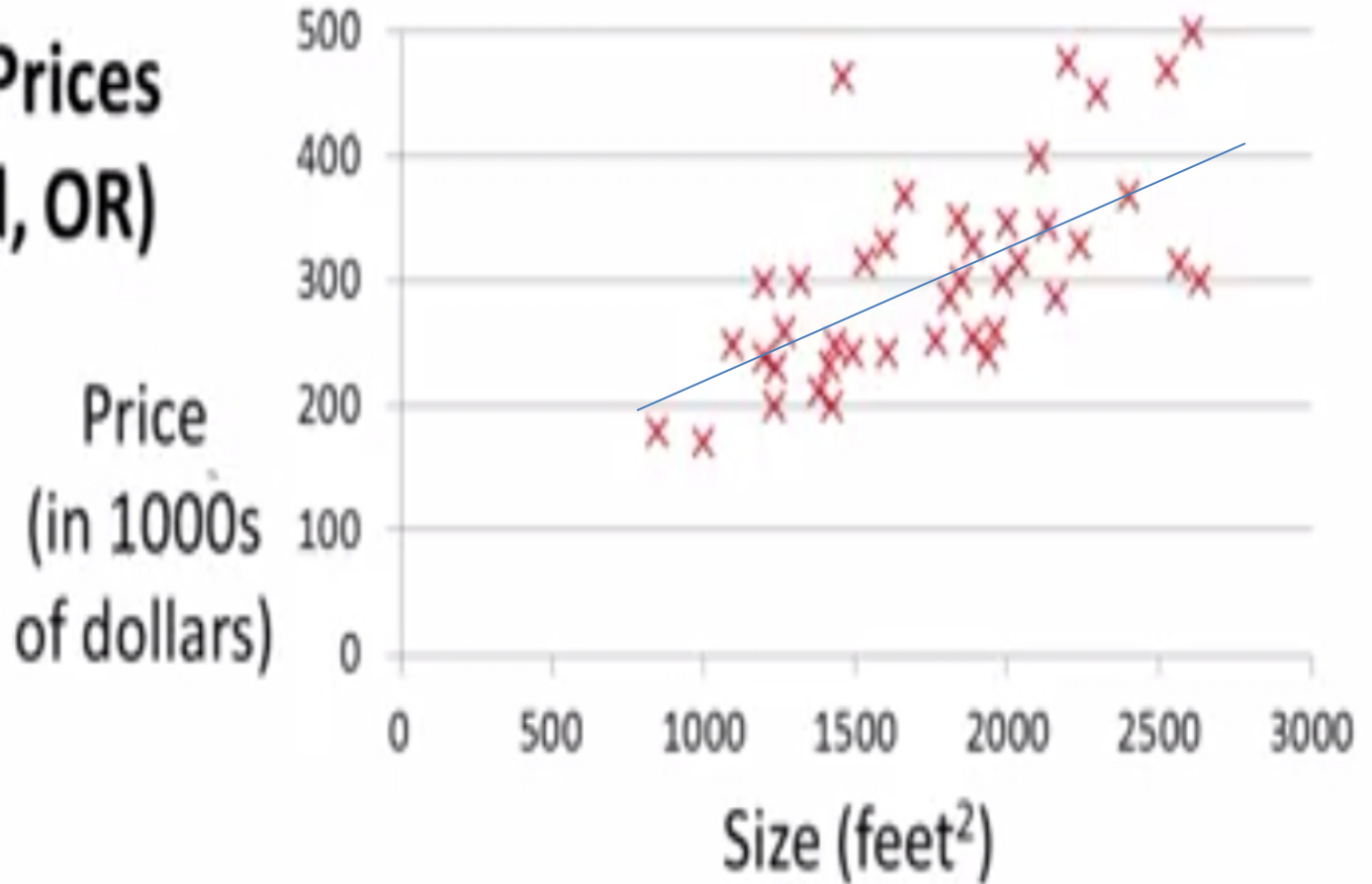
- Machine Learning – a field of predictive modelling
- Linear Regression
 - Developed in the field of statistic
 - A model that attempts to show the relationship between two variables with a linear equation
 - Borrowed by Machine Learning
- Finds it's application in
 - Evaluating Trends and Sales Estimates
 - Analyzing the Impact of Price Change and many more

LINEAR REGRESSION

- Basic Framework
 - Dependent variable \mathbf{y} , also called the output/explained variable which is to be predicted
 - Independent variables \mathbf{x}_i , also called the input/explanatory variables that is to be used for making predictions.
- Regression is the general task of attempting to predict the value of output variable \mathbf{y} from the input variables \mathbf{x}_i .

MODEL REPRESENTATION

Housing Prices (Portland, OR)



MODEL REPRESENTATION

Training set of housing prices (Portland, OR)	Size in feet ² (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178

Notation:

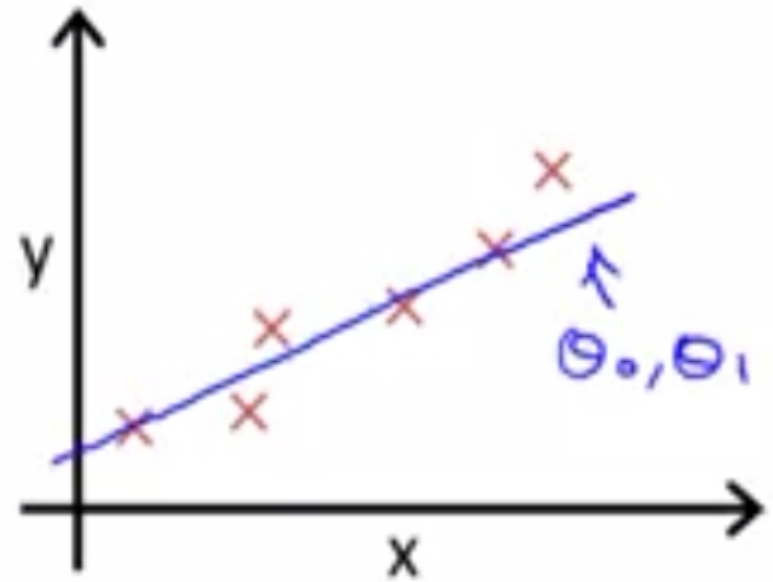
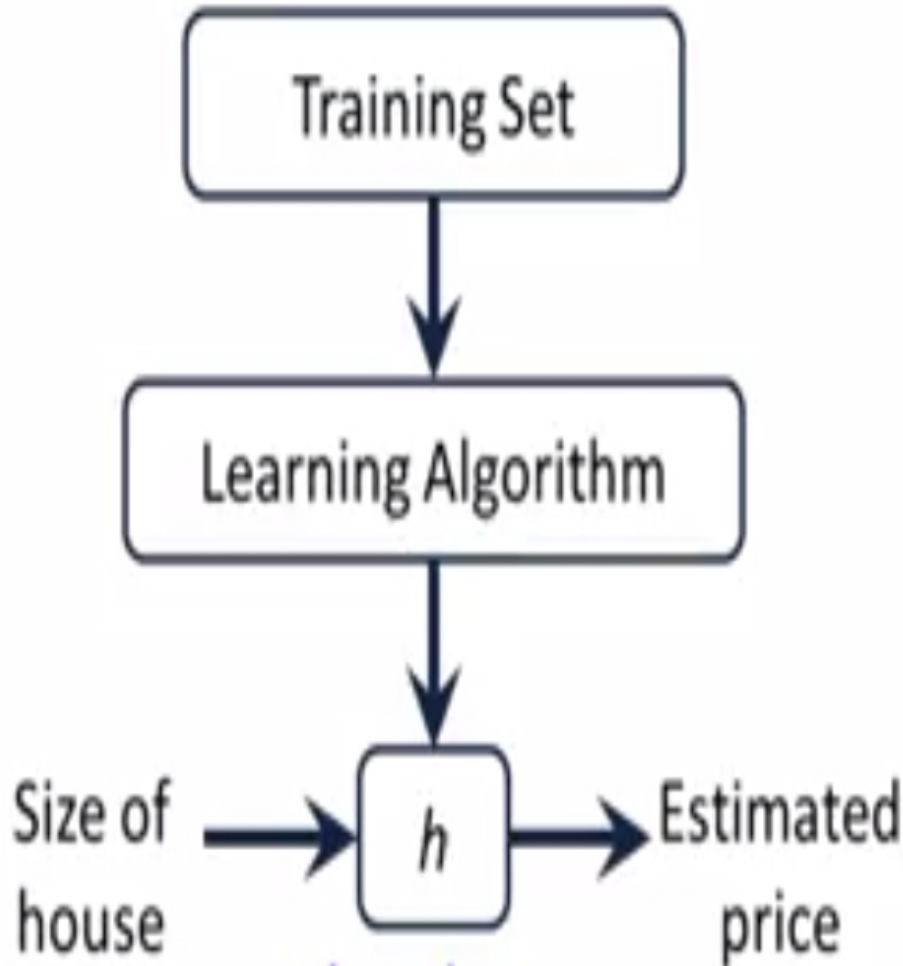
m = Number of training examples

x's = "input" variable / features

y's = "output" variable / "target" variable

MODEL REPRESENTATION

How do we represent h ?



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

MODEL REPRESENTATION

Given a training set ,

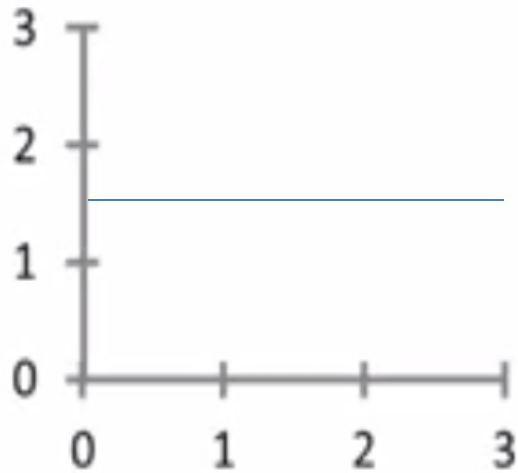
Hypothesis : $h_{\theta}(x) = \theta_0 + \theta_1 x$

where θ_i 's : Parameter

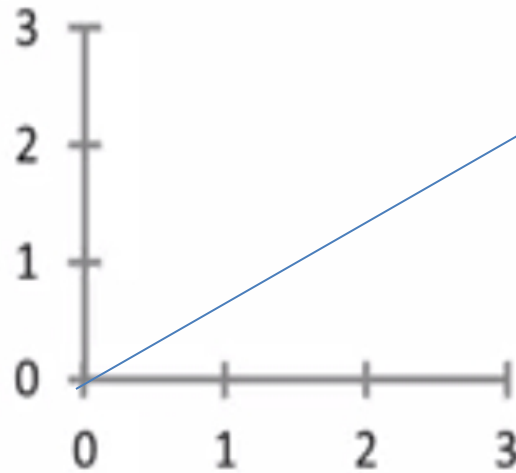
How to choose θ_i 's

MODEL REPRESENTATION

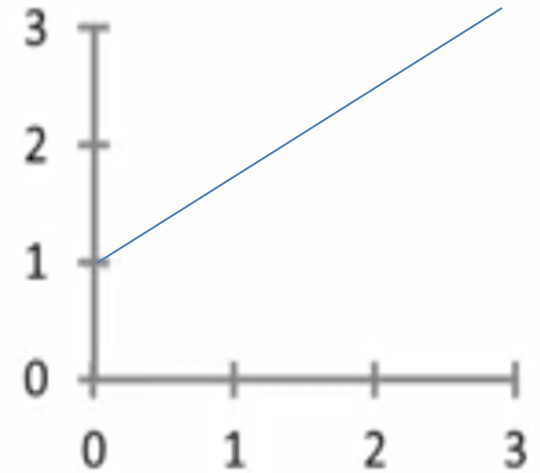
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



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

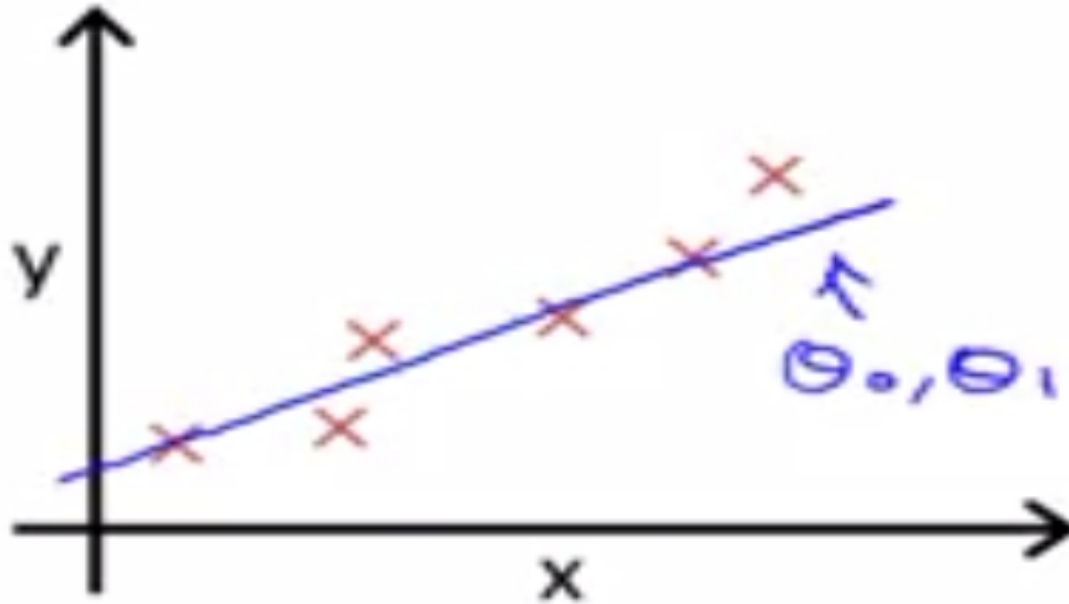


$$\theta_0 = 0$$
$$\theta_1 = 0.5$$



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

COST FUNCTION



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

COST FUNCTION

So going with the idea ,

If hypothesis : $\mathbf{h}_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 \mathbf{x}$

We want to minimise the squared error function

i.e. $(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$ i.e the difference between the predicted value ($\mathbf{h}_{\theta}(\mathbf{x})$) and the actual label (\mathbf{y}) should be as minimum as possible.

COST FUNCTION

Therefore, the cost function becomes

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

where summation(Σ) represents sum over my training set from i to m .

and where ($\frac{1}{2m}$) represents that minimising one half of something shall give us same values of θ_0 and θ_1 as minimising the entire thing.

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

So, overall we want to find the value of θ_0 & θ_1 that minimises the entire expression.

GRADIENT DESCENT ALGORITHM

- Gradient Descent algorithm is basically used to minimise some function J (say a cost function).
- Problem Setup
 - Have some function $J(\theta_0, \theta_1)$
 - Want to $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$
- Outline
 - Start with some θ_0 & θ_1 (say $\theta_0 = 0, \theta_1 = 0$)
 - Keep changing θ_0 & θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum θ_0, θ_1

GRADIENT DESCENT ALGORITHM

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1) \quad , \quad (\text{for } j = 0 \ \& \ j = 1) \quad \}$$

where $:=$ is the assignment operator

α is the learning rate that controls how big a step we can take downhill while creating descent

- Simultaneous update

$$\text{temp 0} := \theta_0 - \alpha \frac{\delta}{\delta \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp 1} := \theta_1 - \alpha \frac{\delta}{\delta \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp 0}$$

$$\theta_1 := \text{temp 1}$$

GRADIENT DESCENT INTUITION

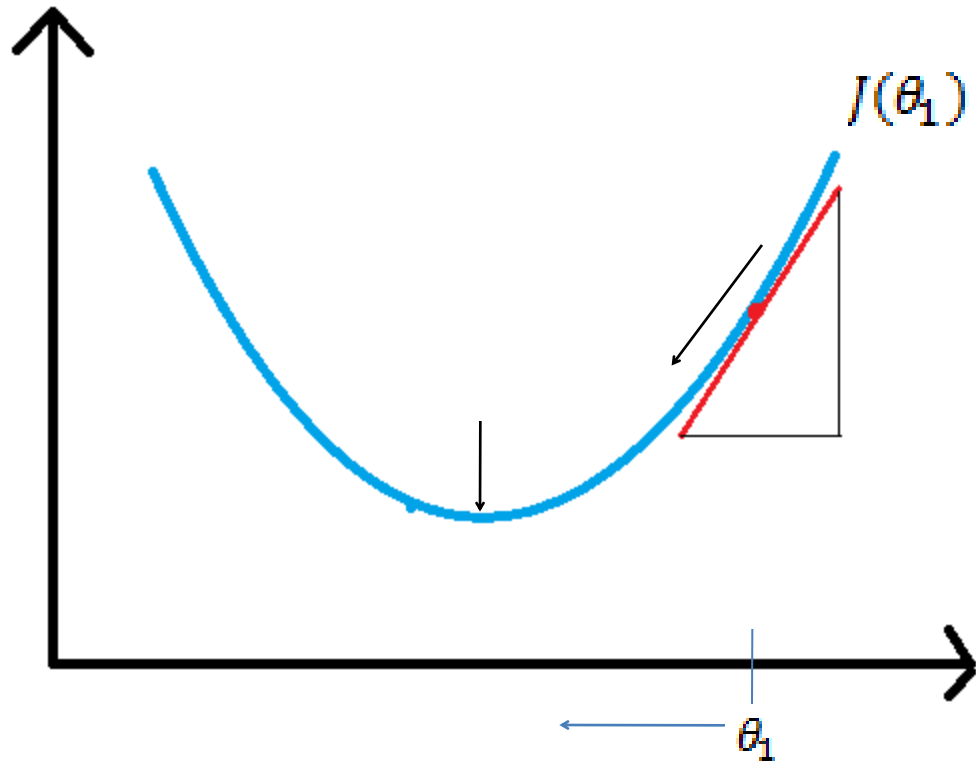
For eg :) say we have only one parameter, so,

$\min_{\theta_1} J(\theta_1)$ & $\theta_1 \in \mathbb{R}$, hence the GD equation becomes

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

$\frac{d}{d\theta_1} J(\theta_1)$ is the derivative term which basically denotes the slope of the line that is just tangent to the function $J(\theta_1)$ at the point θ_1

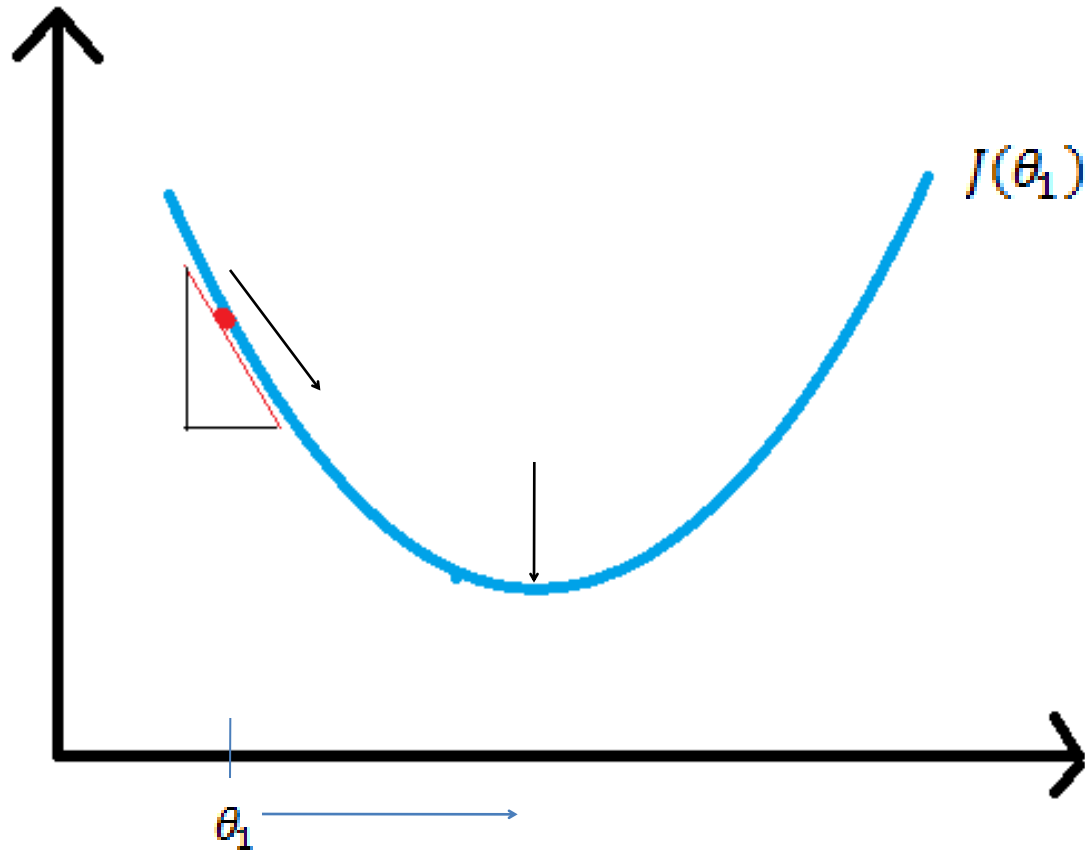
GRADIENT DESCENT INTUITION



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

$$\theta_1 := \theta_1 - \alpha \text{ (positive no.)}$$

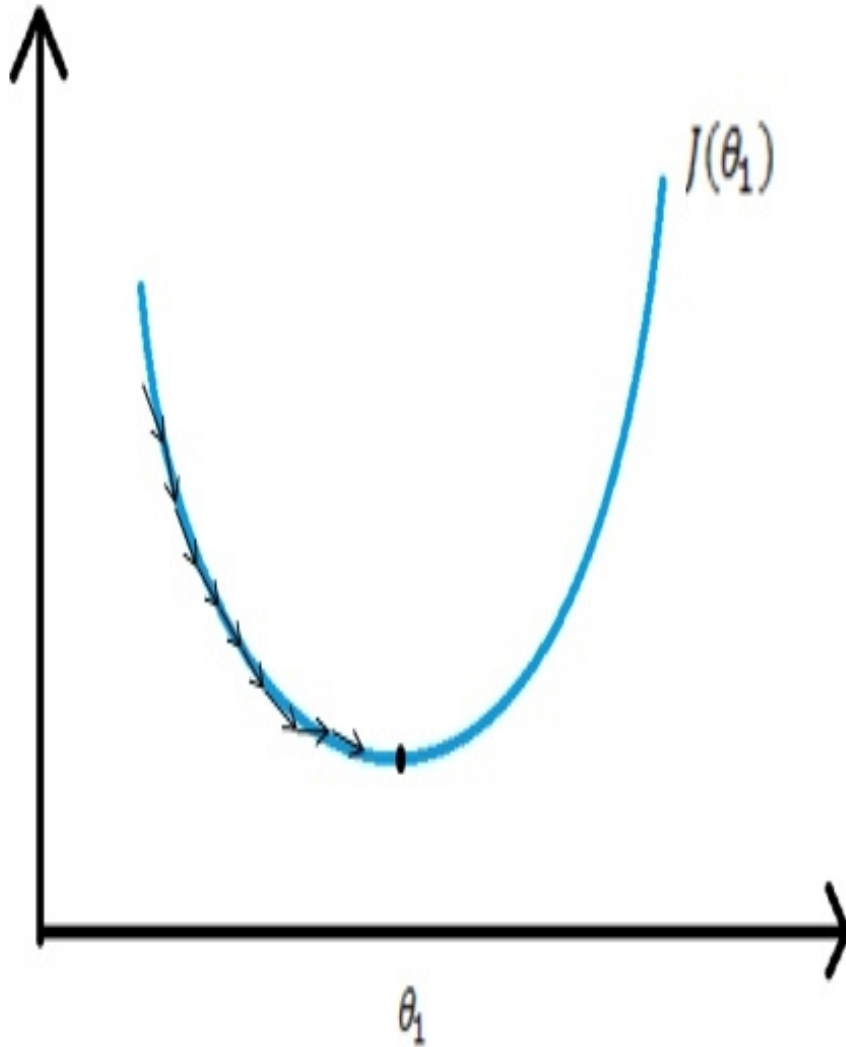
GRADIENT DESCENT INTUITION



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

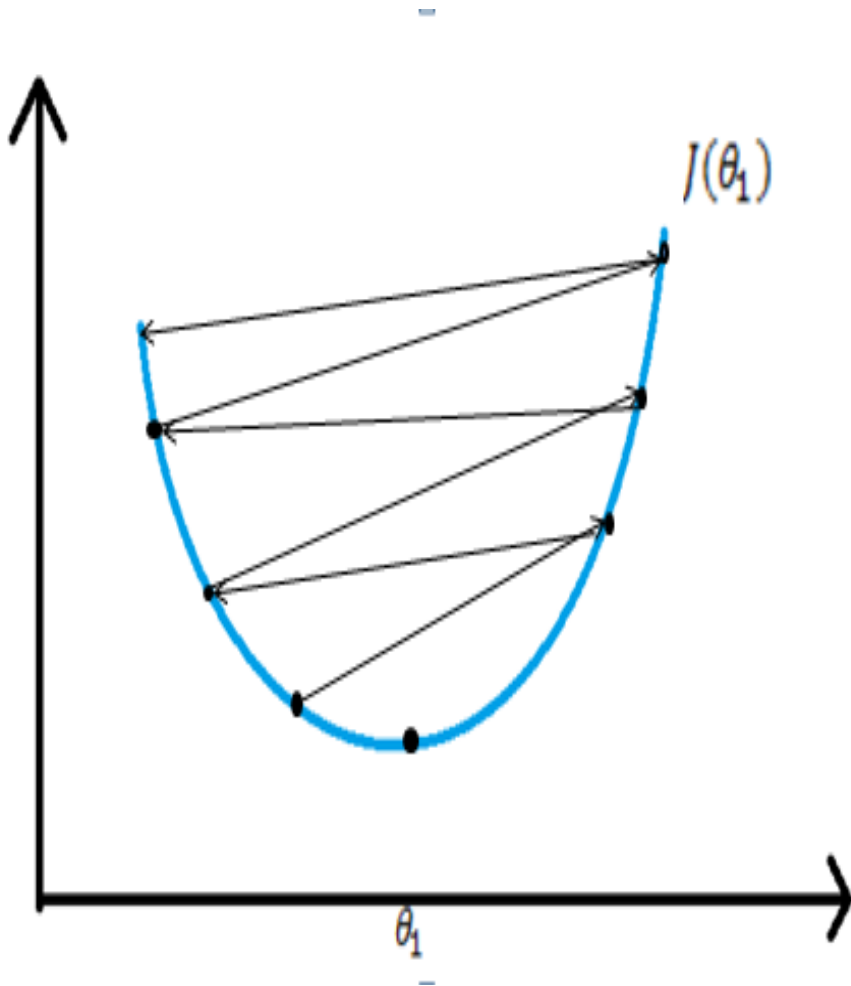
$$\theta_1 := \theta_1 - \alpha (\text{negative no.})$$

GRADIENT DESCENT INTUITION



If α is too small, then gradient descent will converge slowly because very small and hence lot of steps are being taken before it gets anywhere close to the global minimum.

GRADIENT DESCENT INTUITION



If α is too large, then gradient descent will overshoot the minimum because the steps taken are huge . It may fail to converge .

GRADIENT DESCENT FOR LINEAR REGRESSION

So, applying GD algorithm to minimise the squared error cost function i.e. $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

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

$$\mathbf{h}_\theta(\mathbf{x}) = \theta_0 + \theta_1 \mathbf{x}$$

&

Calculating the derivative part,

$$\frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1) = \frac{\delta}{\delta \theta_j} \frac{1}{2m} \sum_{i=1}^m (\mathbf{h}_\theta(\mathbf{x}^{(i)}) - y^{(i)})^2$$

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

GRADIENT DESCENT FOR LINEAR REGRESSION

So for θ_0 , or $j = 0$:

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

for θ_1 , or $j = 1$:

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

GRADIENT DESCENT FOR LINEAR REGRESSION

Now , by applying GD algorithm to the previous equation(derivative part) ,

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\mathbf{h}_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\mathbf{h}_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) \cdot \mathbf{x}^{(i)}$$

} update θ_0 & θ_1 simultaneously

FEATURE SCALING

- If one feature has a range say 0-2000 & another say 0-5 i.e. a 2000 to 5 ratio, then the contours of the cost function takes up a very skewed elliptical shape & the GD may end up taking a long time & can oscillate back & forth before it can finally find its way to the global minimum.
- A useful thing to do is scale the features i.e. the different features take on similar range of values.
- So, by feature scaling we get the features in a range of **-1** to **+1** .

LINEAR REGRESSION WITH MULTIPLE VARIABLES

for multiple features ,

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

For convenience of notation , define $\mathbf{x}_0 = 1$

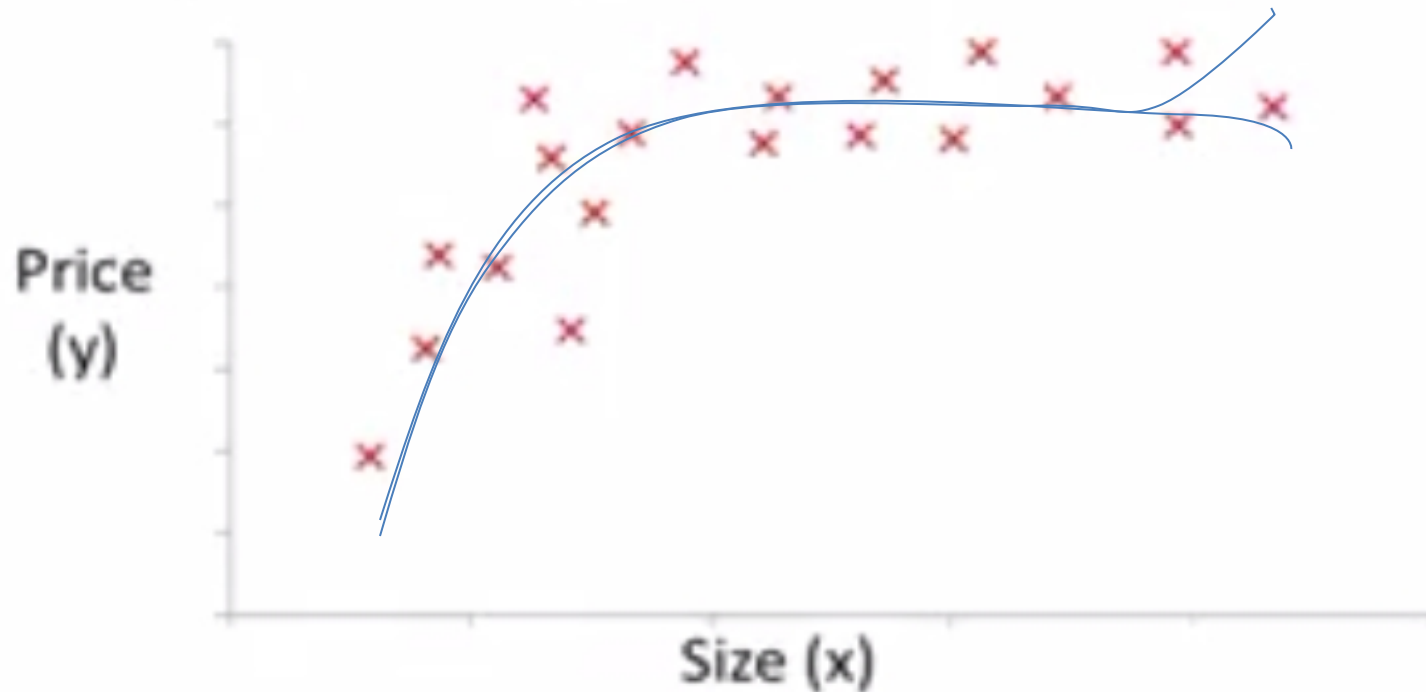
$$\mathbf{X} = \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} , \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \cdot \\ \cdot \\ \cdot \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\text{So, } \mathbf{h}_\theta(\mathbf{x}) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\therefore \mathbf{h}_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$$

POLYNOMIAL REGRESSION

- It allows to use the machinery of linear regression to fit very complicated even non-linear function to the data.



Quadratic model : $\theta_0 + \theta_1x + \theta_2x^2$

Cubic model : $\theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3$

POLYNOMIAL REGRESSION

The form of hypothesis is :

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

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

where $x_1 = \text{size}$, $x_2 = (\text{size})^2$, $x_3 = (\text{size})^3$

- Feature Scaling is required here

NORMAL EQUATION

- GD is basically an algorithm for linear regression that takes multiple iterations to reach the global minimum or to minimise a cost function $J(\theta)$.
- Normal Equation is a method to solve for θ analytically i.e. in one step we get the optimum value .
- Intuition : If 1D , ($\theta \in \mathbb{R}$)
$$J(\theta) = a\theta^2 + b\theta + c$$

To minimise the above quadratic function, we calculate the derivative of the function and set it to 0 i.e. $\frac{d}{d\theta} J(\theta) = 0$ & solve for θ .

NORMAL EQUATION

Examples: $m = 4$.

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

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

GRADIENT DESCENT Vs NORMAL EQUATION

GRADIENT DESCENT

- Need to choose α .
- Needs many iterations
- Works well even if n is large where n is the no. of features.

NORMAL EQUATION

- No need to choose α .
- Needs no iteration
- Slow if n is large because it needs to compute $(X^T X)^{-1}$ which is $n*n$ matrix .