

Day - 1

* What is Machine Learning?

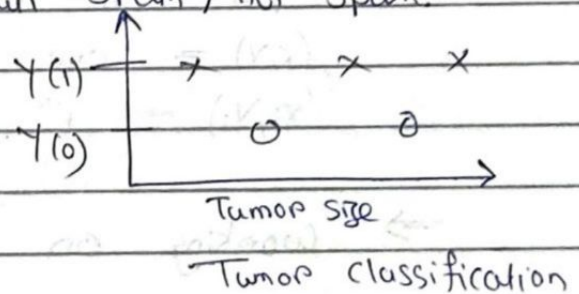
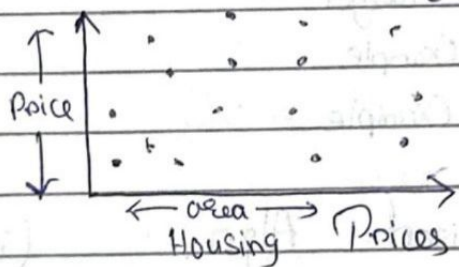
- ① Field of Programming that gives machine the ability to learn without being explicitly Programmed
- ② A Computer Program is said to learn from exp E w.r.t Some task T and Perf. Parameter P. If it's Performance on T as measured by P improved with Experience E.

Supervised Learning

Here the available dataset is such that it consists of Correct answers, and task of this algo to Produce much more right / correct answers.

Regression:- Predicting Continuous valued o/p
e.x. House Prices

Classification:- Predicting discrete valued o/p
e.x. Email Spam / not Spam.



Unsupervised Learning

This type allows us to approach toward solution with very little or no idea about the solution. We can derive structures from data where we don't necessarily know effect of Variable.

NOTE:-

Unsupervised learning don't have flb on Pred Result.

Clustering:- Collection of 100000 balls, we have to find out way to cluster them on the basis of various Parameter.

Non-clustering = Finding structure in chaotic env.

* Linear Regression with one Variable

⇒ Model Representation

Training set of Housing Prices	Size in ft^2 (x)	Price (y)
	2104	460
	1416	232
	1534	315

Notations:-

m = No. of training examples

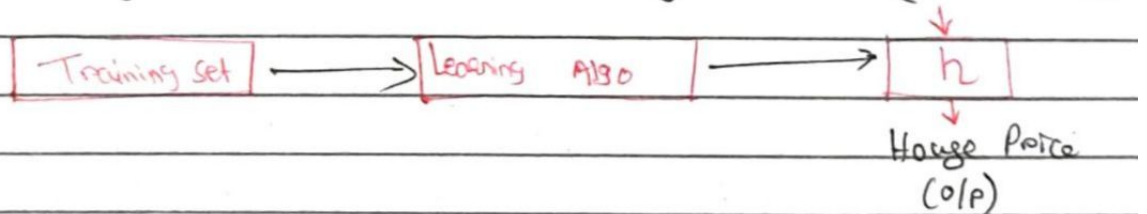
X 's = i/p Variable / feature

Y 's = o/p Variable / Target

(x, y) = one training example

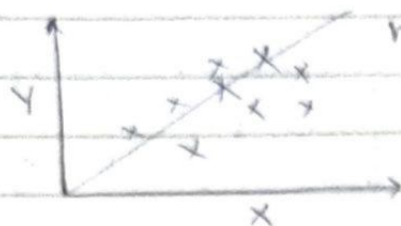
(x_i, y_i) = i^{th} training example

⇒ working on Supervised Algo. (i/p) (House Size)



Hypothesis Matches x 's with y 's.

⇒ Representation of hypothesis.



$$h(x_i) = \theta_0 + \theta_1 x_i$$

Here hypothesis is fitting linearly, so it's called linear Regression with one variable.

⇒ Cost Function.

Here $h(x_i)$ Represents the Predicted response Value for i^{th} obs.

θ_1 = slope ; θ_0 = Y-intercept

We want minimum Value of θ_1 and θ_0 So we will implement Cost

i.e.

Sq. error cost fun $\leftarrow J(\theta_0, \theta_1) = \frac{1}{2m} (h(x_i) - y^{(i)})^2$

Predicted o/p
actual o/p

hence,

$$\theta_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2} = \frac{\sum_{i=1}^m x_i y_i - m \bar{x} \bar{y}}{\sum_{i=1}^m x_i^2 - n(\bar{x})^2}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

⇒ Gradient Descent

Repeat untill convergence {

$$\text{temp } \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x_i) - y^{(i)})$$

$$\text{temp } \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x_i) - y^{(i)}) \cdot x^{(i)}$$

} $\theta_0 = \text{temp } \theta_0$ and $\theta_1 = \text{temp } \theta_1$

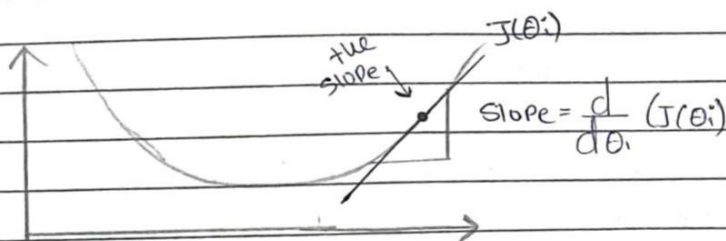
⇒ Gradient Descent intuition

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

\downarrow
 Learning Rate / step size

\rightarrow Derivative

Case I:-



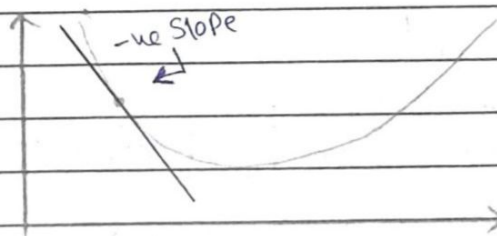
Here to Reach minima theta should decrease so,

$$\theta_j := \theta_j - \alpha \underbrace{\frac{\partial}{\partial \theta_j} J(\theta_j)}_{\geq 0}$$

derivative Part must be greater than zero / Positive

i.e $\theta_j := \theta_j - \alpha (\text{Positive Num})$ [θ_j will \downarrow]

Case II:-

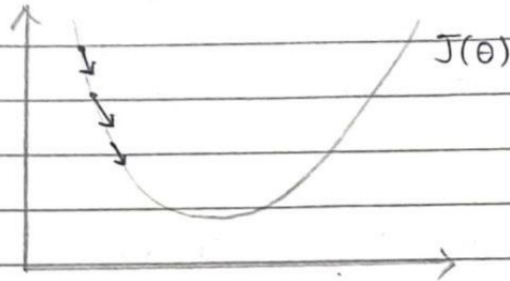


$$\theta_j := \theta_j - \alpha \underbrace{\frac{\partial}{\partial \theta_j} J(\theta_j)}_{\leq 0}$$

here θ_j have to decrease to Reach ~~global~~ minima.

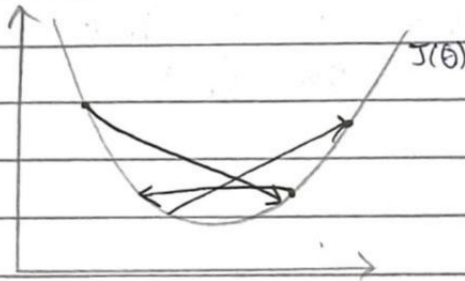
⇒ what α Does?

Case 1 :- α is very Small



If α is very Small
Grad descent will be
Slow

Case 2 :- α is very large



If α is too large Grad
descent will overshoot
local minima

* Linear Regression with Multiple Variables

Sample Data:-

Size (feet ²)	No. of bedroom	No. of floor	Age of home	Price
X_1	X_2	X_3	X_4	(Y)

Notation

n = No. of features

$X^{(i)}$ = i^{th} training ex

$X_j^{(i)}$ = Value of feature j in i^{th} training ex

Hypothesis:-

$$h_\theta(x) = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \theta_3 X_3 + \theta_4 X_4$$

here

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

So,

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

$$= [\theta_0 \ \theta_1 \ \dots \ \theta_n]$$

Parameters vect

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

- feature vect.

* Gradient descent for LR multiple Variable

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

$$\text{Parameter} = \theta_0, \theta_1, \dots, \theta_n$$

Cost fun =

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

Gradient Descent

Repeat {

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

}

for $n=1$

Repeat {

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

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

}

for $n > 1$

Repeat {

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

}

let $n=2$

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_0^{(i)} \quad \text{as } x_0 = 1 \quad \text{①}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_1^{(i)} \quad x_1 = x^{(1)}$$

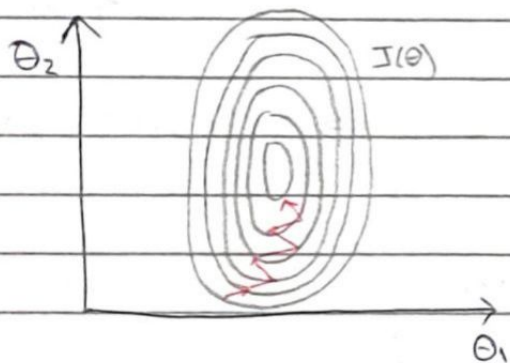
$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_2^{(i)}$$

* Gradient Descent Feature Scaling

Feature Scaling :- Making sure that features are on similar scale

F.g. $x_1 = \text{Size of House (0-2000) } \text{ft}^2$

$x_2 = \text{No. of bedrooms (1-5)}$

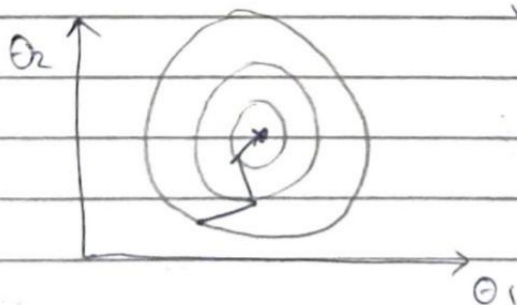


\Rightarrow This will take long time to Reach local minima

So, we Scale feature

$$\text{i.e. } x_1 = \frac{\text{Size (ft}^2\text{)}}{2000} \Rightarrow 0 \leq x_1 \leq 1$$

$$x_2 = \frac{\text{No. of bedrooms}}{5} \Rightarrow 0 \leq x_2 \leq 1$$



\Rightarrow we can Reach local Minima faster

* Feature Scaling

- Get every feature into approx a $-1 \leq x_i \leq 1$ Range

E.g.

$0 \leq x_1 \leq 3$	✓ OK	} Not too large Not too small
$-2 \leq x_2 \leq 0.5$	✓ OK	
$-100 \leq x_3 \leq 100$	X Not OK	
$-0.001 \leq x_4 \leq 0.001$	X Not OK	

→ Generally ⇒ $-3 \leq x_i \leq 3$

* Mean Normalization

Replace x_i with $x_i - \mu_i$ to make feature have approximately zero mean.

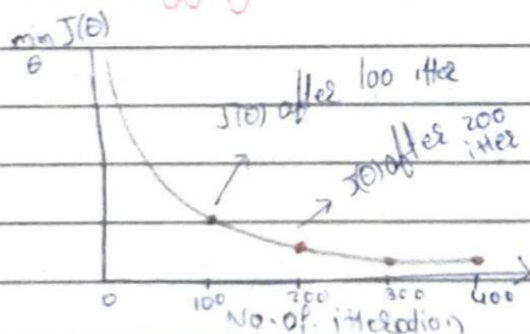
i.e. $x_1 = \frac{\text{Size} - 1000}{2000} \Rightarrow -0.5 \leq x_1 \leq 0.5$

$x_2 = \frac{\text{bedroom} - 2}{5} \Rightarrow -0.5 \leq x_2 \leq 0.5$

Generally,

$x_i \leftarrow x_i - \mu_i$ - Avg value of x training set
 σ_i - Range (max - min) / std. dev

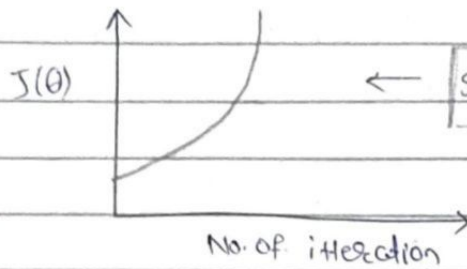
* Debugging Grad. decent



$J(\theta)$ should decrease after every iteration.

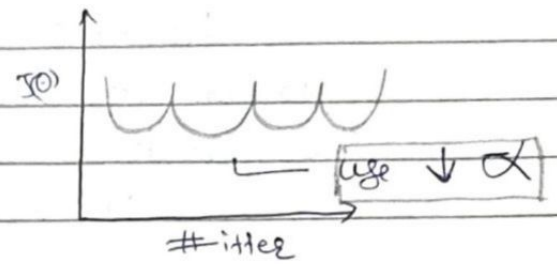
"Debugging" :- Making sure that Grad descent working Properly.

Case I



← Soln: $\downarrow \alpha$

Case II



→ Our α should be suffi small so that $J(\theta)$ will decrease after every iteration

→ If α is too small then, good descent can be slow to Converge.

* Features and Polynomial Regression

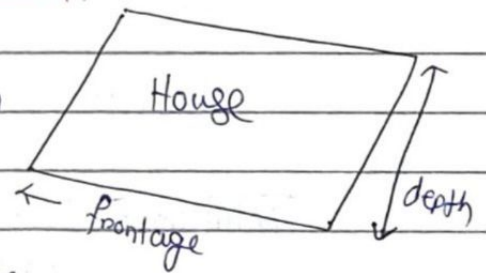
House Price Prediction.

- we are given with two features

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{frontage}) X_1 = \text{frontage} + \theta_2(\text{depth}) X_2 = \text{depth}$$

- but we can create our own feature

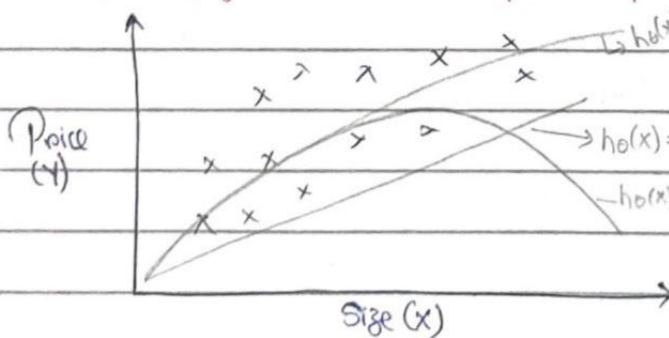
i.e $X_3 = \text{frontage} * \text{depth} = \text{land area}$



Now,

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{land area}) \quad \text{--- New hypo}$$

* Polynomial Regression



• we can see lined model

Can't fit here $\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1x$

• we won't fit Quadratic model as it will come down

$$\hookrightarrow h_{\theta}(x) = \theta_0 + \theta_1x + \theta_2x^2$$

• we are going to fit Cubic model

$$h_{\theta}(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3$$

Now, $h_\theta(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$ — (1)

$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$ — (2)

to map eq. (1) and eq (2)
we can do,

$$\boxed{x_1 = (\text{size})} \quad \boxed{x_2 = (\text{size})^2} \\ \boxed{x_3 = (\text{size})^3}$$

size : 1 - 1000 ; size² : 1 - 1000000 , size³ : 1 - 10⁹

So our features are showing Varied range

So here we will apply

FEATURE SCALING

Other choice,

$$\boxed{h_\theta(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\sqrt{\text{size}})}$$

* Normal Equation

Normal Eq:- Meth to solve θ analytically

Intuition:- If 1D ($\theta \in \mathbb{R}$)

$$J(\theta) = a\theta^2 + b\theta + c$$

To minimize $J(\theta) \Rightarrow \frac{d}{d\theta} J(\theta) = 0$ — solve for θ
Set

Actually

$$\theta \in \mathbb{R}^{n+1}$$

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad \text{for every } j$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$.

n.T.O

Ex. $m = 4$

	Size (feet ²)	# bedroom	# floors	Age	Price
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

Now,

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}_{m \times (n+1)} \quad Y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}_{m \text{-Vec}}$$

$m = \text{no. of Exam/sample}$
 $n = \text{features}$

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

In General

m examples :- $[(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})]$; n -features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(n)})^T \text{---} \end{bmatrix}_{m \times n+1}$$

(design matrix)

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}_m$$

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

• $(X^T X)^{-1}$ is inversed of $X^T X$.

\therefore Octave : $\text{Pinv}(X' * X) * X' * y$

= Set $A = X^T X$
 then $\Theta = A^{-1} X^T y$

Gradient Descent

Disadvantage = Need to choose α
 = Need many iter
 Advan = Work well even when n is large

Normal Equation

Adv = No need of α
 = Don't need to iterate
 Disadv = Need to compute $(X^T X)^{-1}$
 = Slow if n is very large

use Gradient Descent

← beyond this $n = 10000$ can go till

* Normal eqn (Noninvertibility)

What if $(X^T X)$ is non-invertible (singular)
 In octave we have Pinv and inv to calculate inverse

$\text{Pinv} \Rightarrow$ Pseudo Inverse = although the $X^T X$ is singular it will calculate inv

$\text{inv} \Rightarrow$ Real Inv = No ans if $X^T X$ is singular

When $X^T X$ will be non-invertible?

① Redundant features (linearly dependent)

Deleting Redundant feature will solve this problem

let $X_1 = \text{Size (feet)}^2$

$1\text{m} = 3.28\text{ feet}$

$X_2 = \text{Size (m)}^2$

$\therefore X_1 = (3.28)^2 X_2$

= if X_1 and X_2 are linearly dependent then $X^T X$ is non-invertible

② Too many features (e.g. $m \leq n$)

$m = \text{training ex}$
 $n = \text{features}$

— Delete some feature / Regularization

→ This will solve non-invertibility