

AI & ML

1 row \rightarrow sample

1 column \rightarrow Criteria under that sample

Sample = ~~Attribute~~ Instance

Features = Attribute

In a row:

$n-1$ = Features

$[-1]$ = Output

Classification

Discrete Set of Output.

*** A classification Algo may predict a continuous value but the continuous value is in form of a probability for a class label.

Regression

Continuous Set of Output.

*** A regression Algo may predict a discrete value, but the discrete value in the form of integer quantity.

Paloxi
Palonosetron INN

Just to know!

Human learning vs ML ^{Algo}

Human - brain stimulation.

ML - parameter.

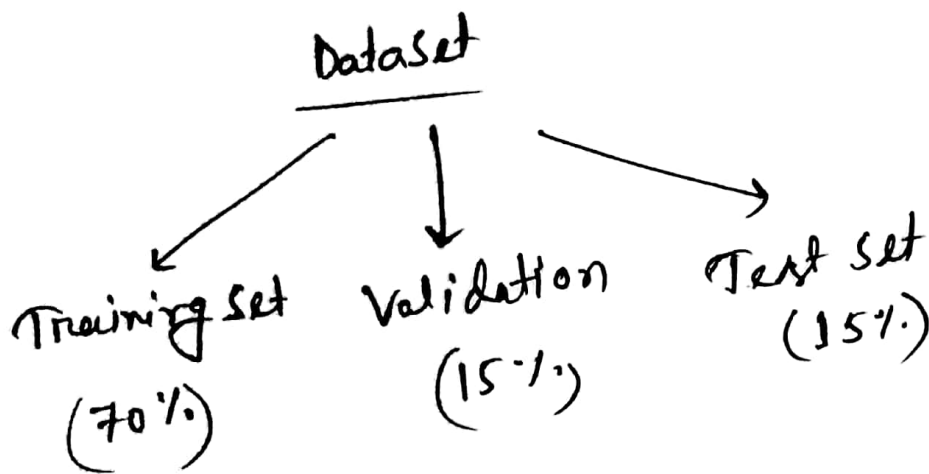
parameter vs Hyperparameter (Not yet)

* Get idea of Training Set, Validation Set & Test Set.

By Example! SSC Exam:

Training Set: Given that u can learn from.

- i) Class, coaching, private tuition: Training Set.
- ii) Model Test, Class Test, private tutor Exam: Validation Set.
- iii) SSC Final Exam: Test Set (Pass, Fail)



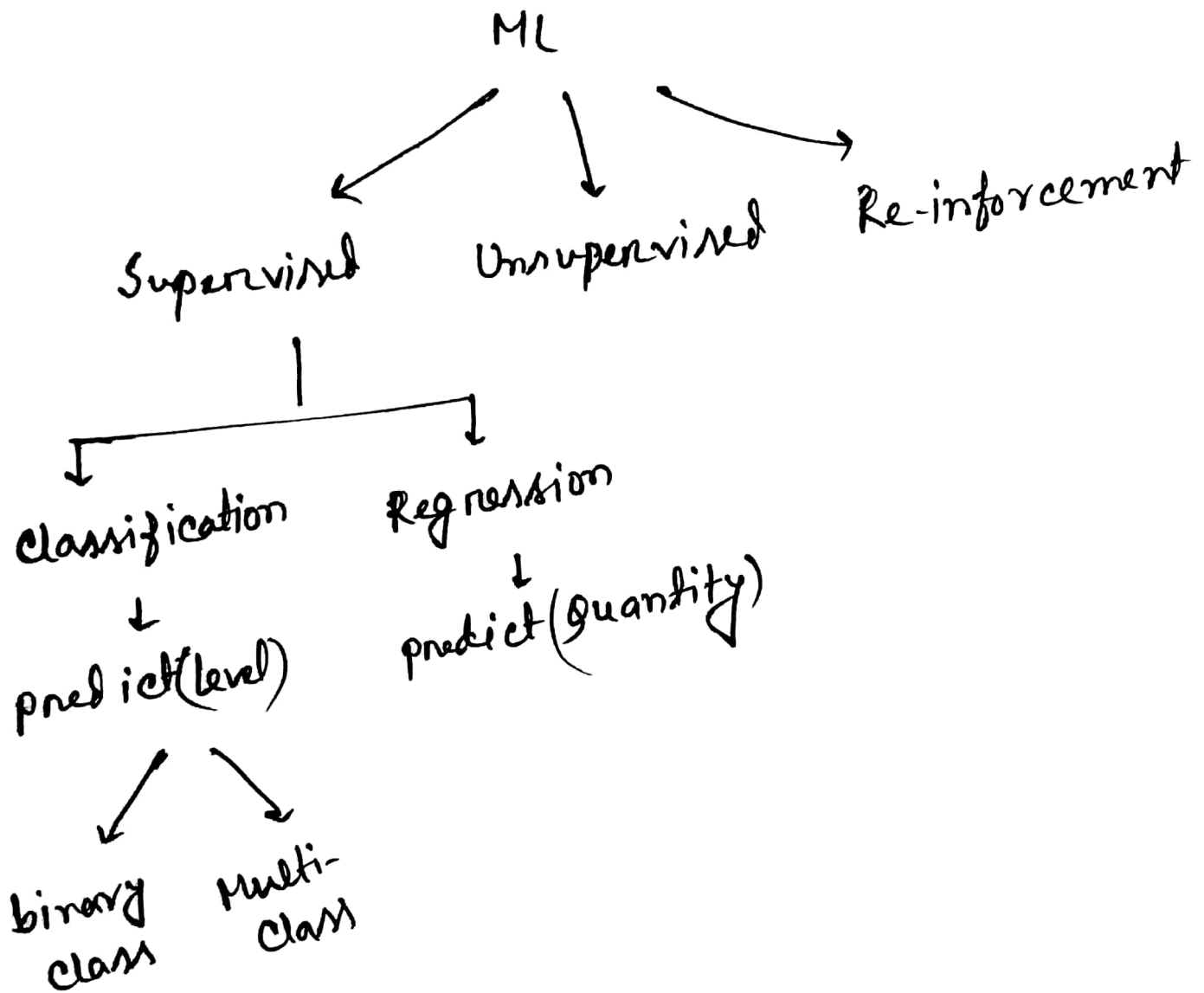
All are disjoint.

Training Set: to fit the ~~para~~ parameters.

Validation Set: to tune the parameters.
[Validation Set is used to estimate prediction error for model ~~selection~~ selection.]

Test Set: to assess the performance. [The test set is used for assessment of the generalization error of the
(judgemental)]

final chosen model.



* What is Linear Regression?

⇒ Linear regression is a method to predict dependent variable (Y) based on values of independent variable (X). It can be used for the cases where we want to predict some continuous quantity.

• 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
...	...

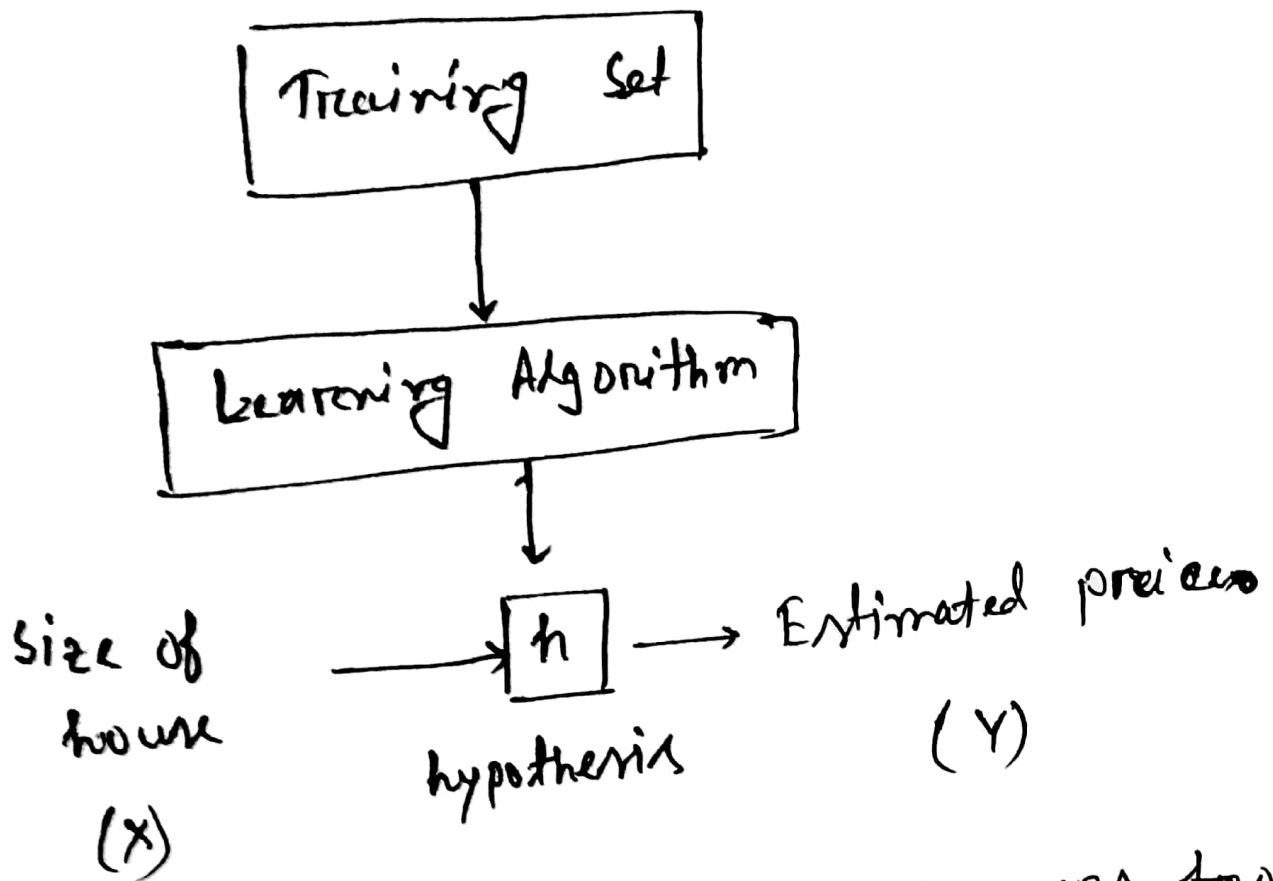
} m = 47

Notation! m = Number of training Examples

X_i = input var / features

y_i = output var / target var

Paloxi
Palonosetron INN

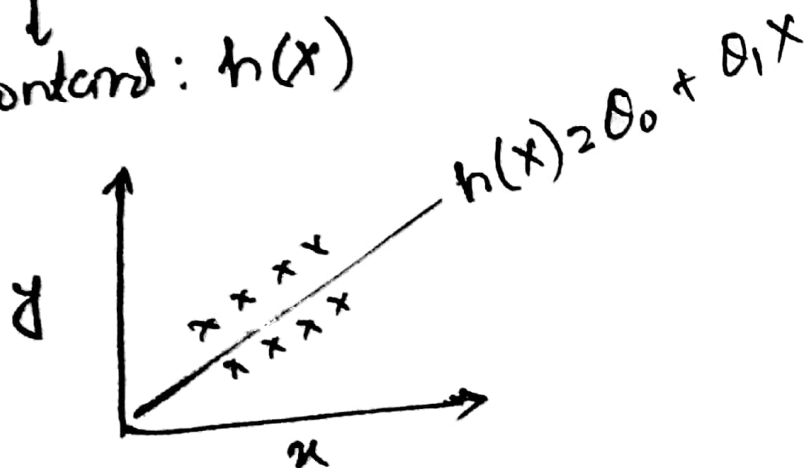


So, h is a function that maps from x 's to y 's.

How do we represent h ?

$$h_0(x) = \theta_0 + \theta_1 x$$

↓
Shortened: $h(x)$



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

θ_i 's : Parameters

How do we choose parameters?

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

* Minimizing the problem;

Minimize
 $\theta_0, \theta_1 \quad (h_{\theta}(x) - y)^2$

for (x_i, y_i) : \rightarrow # training examples.

minimize
 $\theta_0, \theta_1 \quad \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$

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

It means you'll find the values of θ_0 & θ_1 that

causes this expression to be minimized. And this expr depends on θ_0 & θ_1 .

Paloxi
Palonosetron INN

By convention:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

Minimize
 θ_0, θ_1

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

cost function is also
called ~~of~~ Squared
Error ~~cost~~ function.

→ In this case, Not always.

Just Summarize in this particular Area:

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

Parameters: θ_0, θ_1

Cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

Goal: minimize $J(\theta_0, \theta_1)$ to get optimal.
 θ_0, θ_1

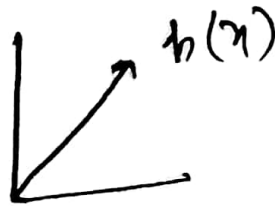
Paloxi
Palonosetron INN

Simplified

$$h_0(x) = \theta_0 + \theta_1 x$$

$$\theta_0 = 0$$

$$h_0(x) = \theta_1 x$$

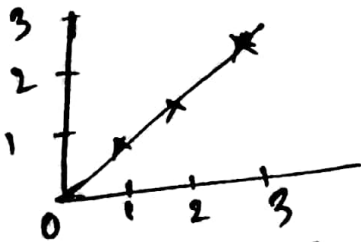


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(h_0(x_i) - y_i)^2}_{\theta_1 x_i}$$

hypothesis func:

$$h_0(x)$$

for fixed θ_1 , this is
function of x



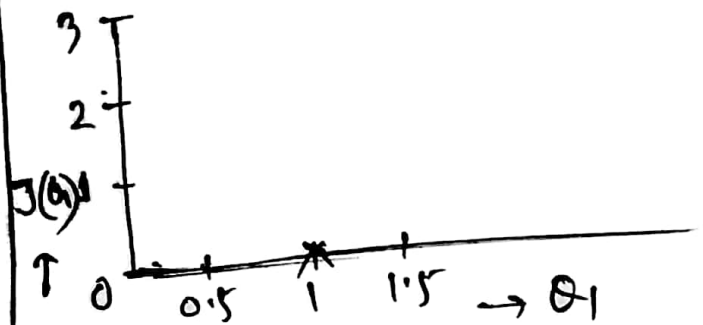
$$\theta_1 = 1$$

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

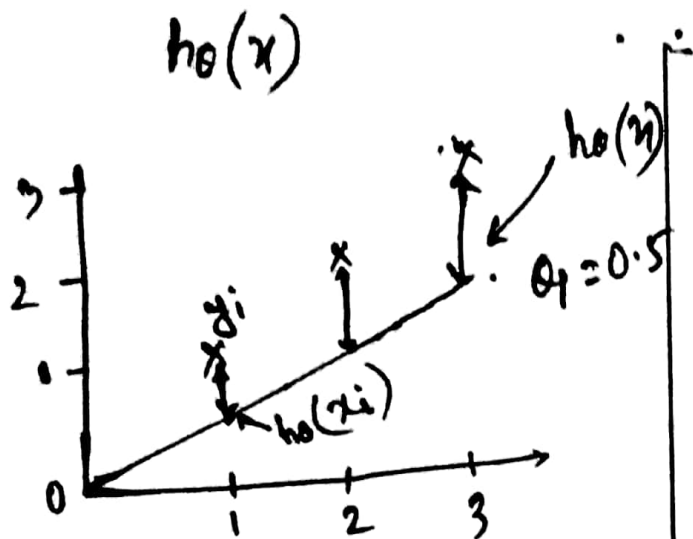
$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$

cost func:
 $J(\theta_1)$

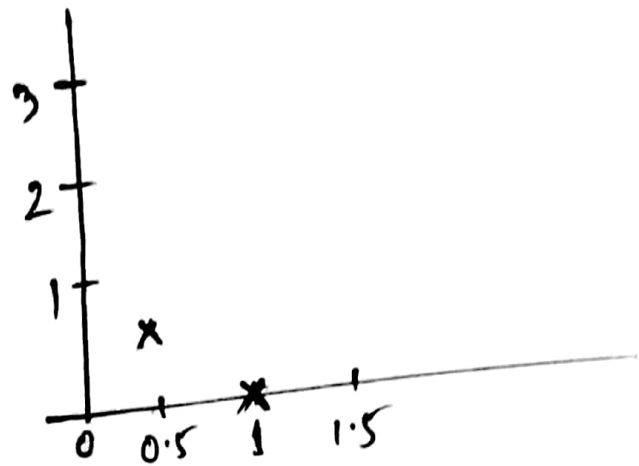
function of the parameter θ_1



$$J(1) = 0$$



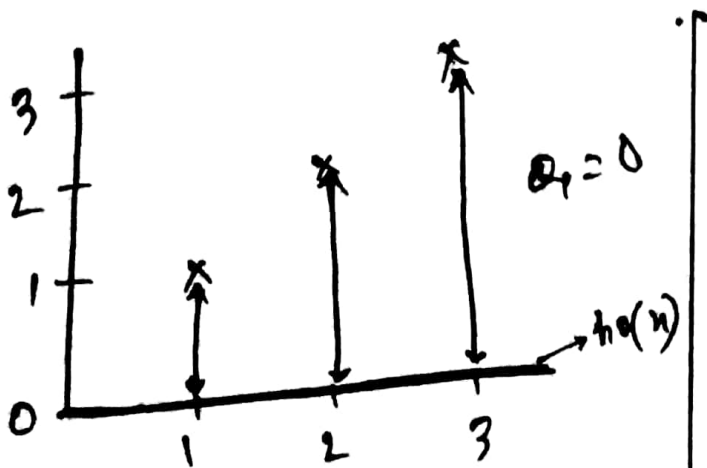
$J(\theta_1)$



$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

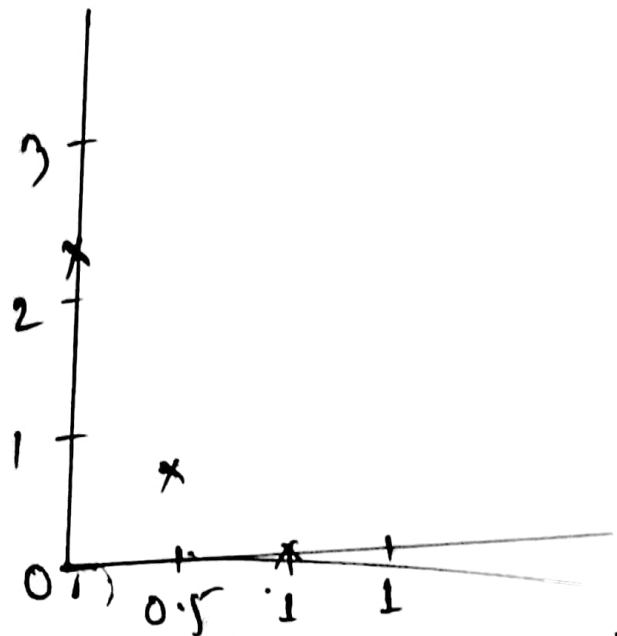
$J(0.68)$

≈ 0.58



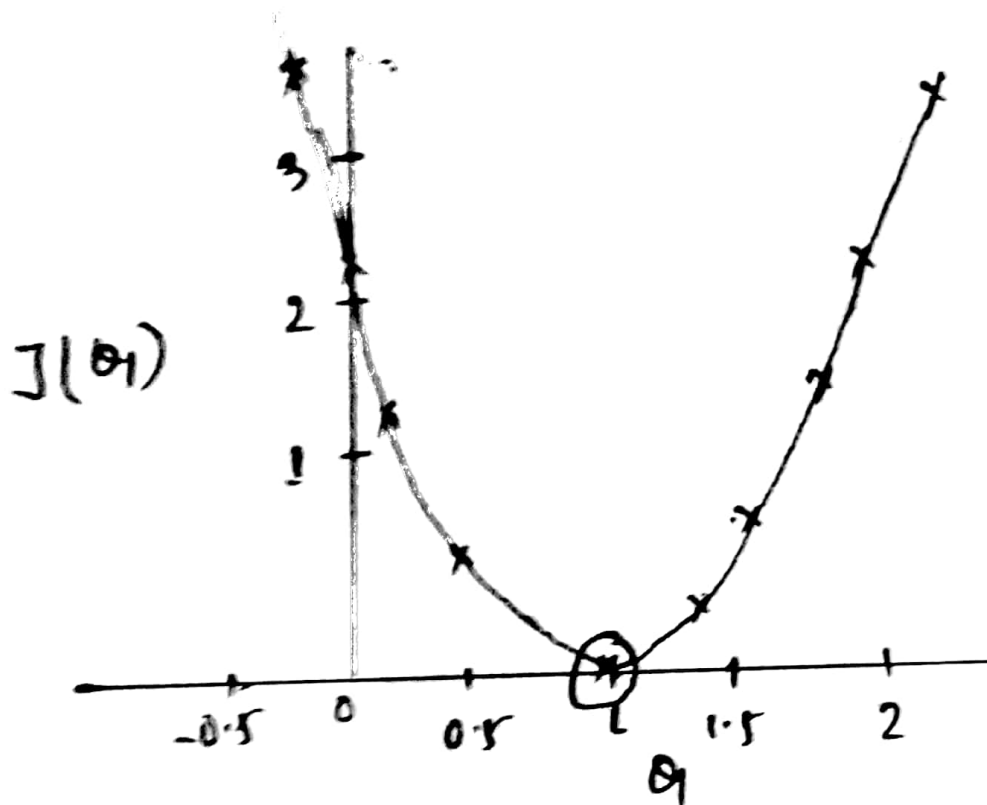
$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2)$$

$$= \frac{1}{6} \cdot 14 \approx 2.3$$



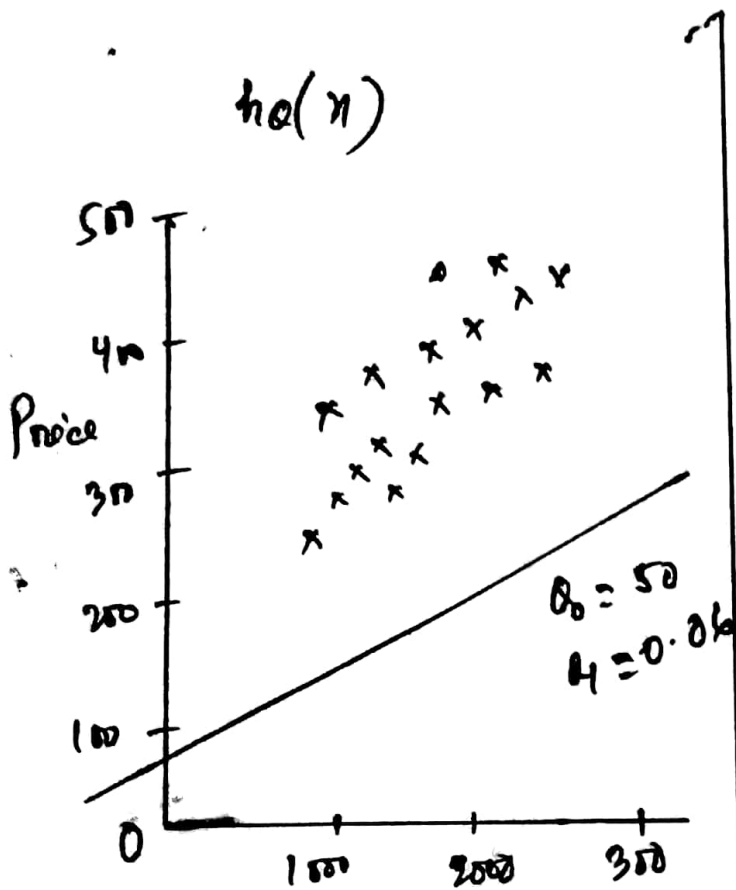
Paloxi
Palonsetron INN

$$J(\theta_1)$$



$$: J(\theta_1)$$

Minimize $J(\theta_1) \hat{=} \theta_1 = 1$



$$h_0(x) = 50 + 0.06x$$

Learn to know!
Contour plots
or Contour figures.

Gradient descent

is an Algorithm to minimize the cost function $J(\theta_0, \theta_1)$. It's not only used in Linear Regression, it's generally used in almost of ML places.

Have some function $J(\theta_0, \theta_1)$

$$\text{Want } \min_{\theta_0, \theta_1} J(\theta_0, \theta_1) \quad \left| \quad \begin{array}{l} J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) \\ \min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n) \end{array} \right.$$

Outline:

- Start with some θ_0, θ_1
- keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum.

:= Assignment

repeat until convergence {

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

learning rate

}
(for $j=0$ and $j=1$)
Simultaneous update
 θ_0 and θ_1

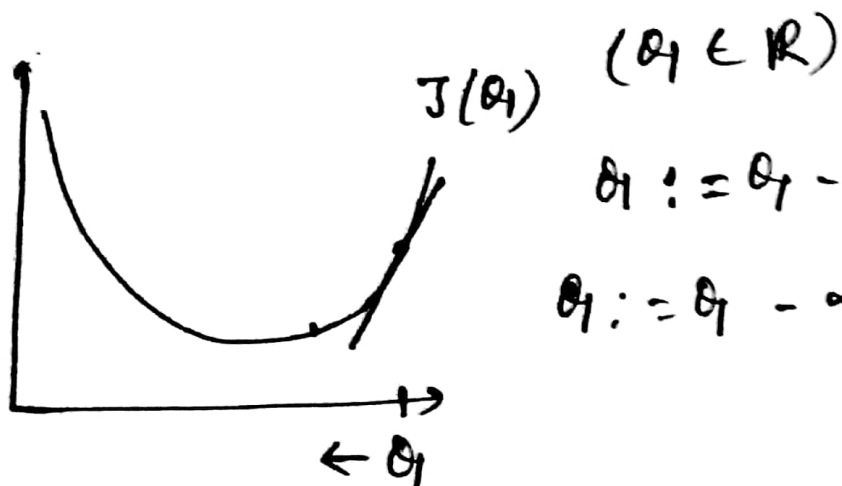
Correct: Simultaneous Update:

$$\Rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\Rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

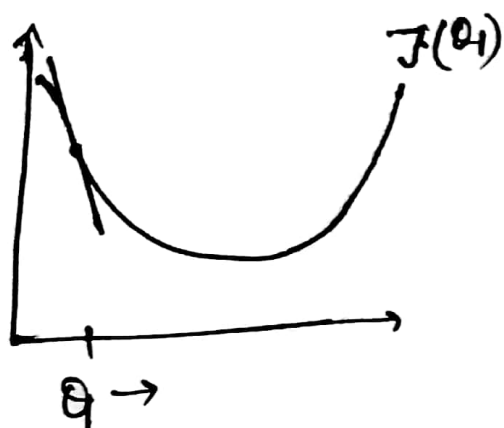
$$\Rightarrow \theta_0 := \text{temp0}$$

$$\Rightarrow \theta_1 := \text{temp1}$$



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

$$\alpha := \alpha \cdot (\text{positive no})$$



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

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

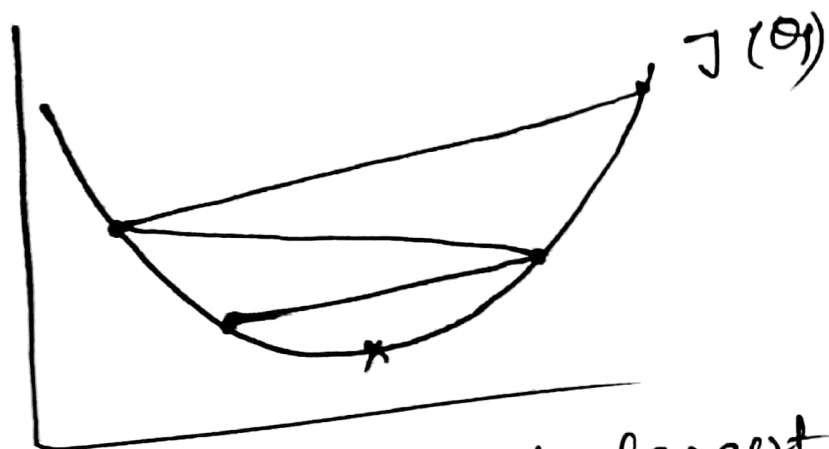
so, if α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

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

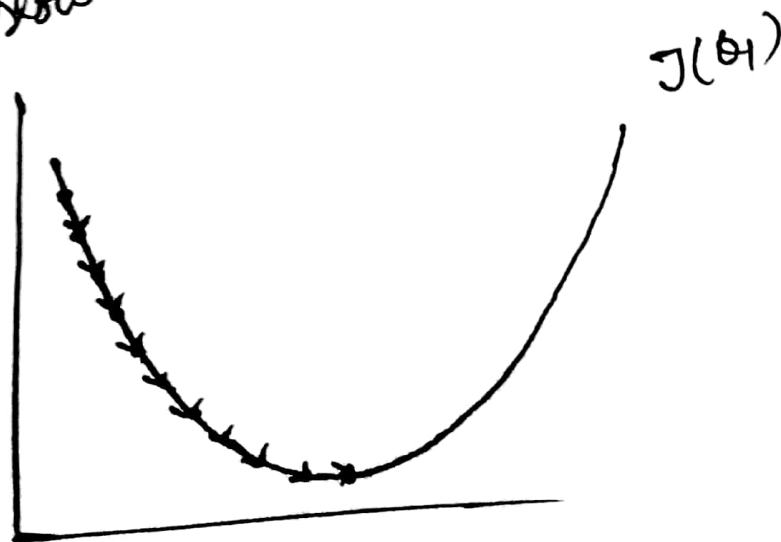
~~$\theta_1 := \theta_0$~~

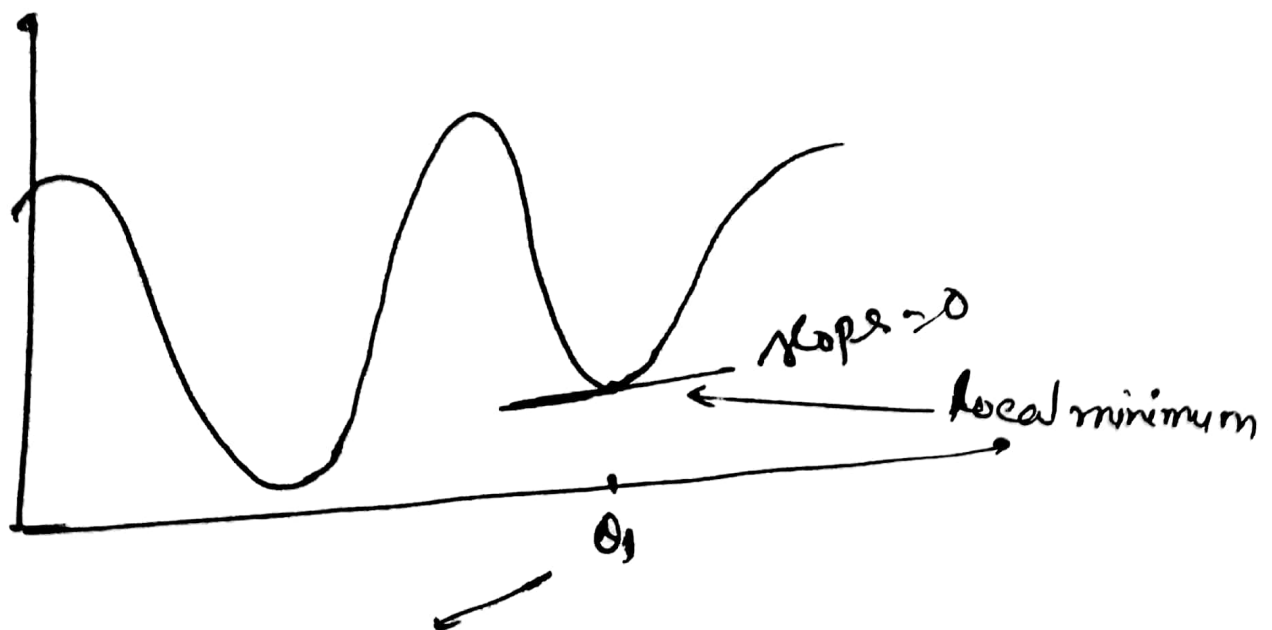
Let's see example:

α is too large means learning steps are very big.



If α is too small, gradient descent can be slow





current
value of θ_1

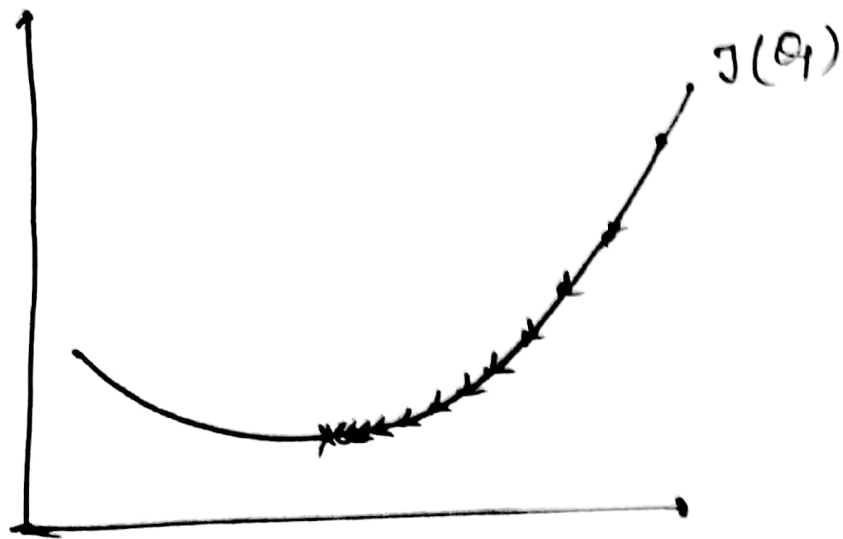
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha \cdot 0$$

$\theta_1 := \theta_1$ remain unchanged.

➤ Gradient descent can converge to a local minimum, even with the learning rate α fixed.

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



As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.

Gradient descent Algo:

Linear Regression Model

repeat until convergence

$$\left\{ \begin{array}{l} \theta_j: \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ \text{for } (j=1 \text{ \& } j=0) \end{array} \right\}$$

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

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

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

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2$$

Now,

$$\theta_0 \mid j=0: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$

$$\theta_1 \mid j=1: \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_i$$

That's final Gradient descent Algo for Linear Reg.

Paloxi
Paloxi.com

A 3D bowl shape function represent
that Linear Regression cost function.
Another called: "convex function".

"Batch" Gradient Descent.

"Batch": Each step of gradient descent
usage uses all the training example.

$$\rightarrow \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$