

# COURSE MACHINE LEARNING

BY  
ANDREW NG,  
STANFORD UNIVERSITY



NOTES BY RUTHUPARNA K (EC dept.)

# MACHINE LEARNING

BY ANDREW NG

## WEEK 1

### INTRODUCTION

What is Machine Learning?

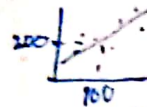
A computer program is said to learn from a experience  $E$  with respect to some task  $T$  and some performance measure  $P$ , if its performance on  $T$ , as measured by  $P$ , improves with experience  $E$ ."

- Tom Mitchell

Supervised Learning.

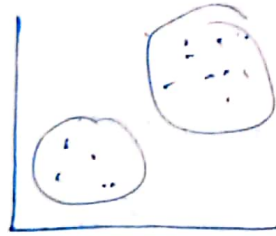
In supervised learning, we are given a labeled data set and already know what our correct output should look like, having the idea that there is a relationship between the input and output.

- Regression: Mapping input variable to some continuous function
- Classification: Mapping input variable to discrete categories



### 3) Unsupervised Learning

When the data isn't labelled, the algorithm finds pattern and clusters the data.



#### b) Notation:

$m$  - number of training examples

$x$  - feature - input variable

$y$  - target variable - output variable

$(x^{(i)}, y^{(i)})$  -  $i$ th training example where  
 $i = 1, 2, \dots, m$

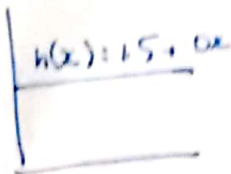


# Univariate Linear Regression

Cost Function

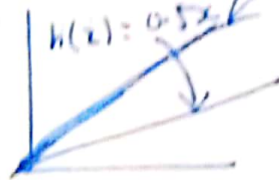
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

hypothesis function



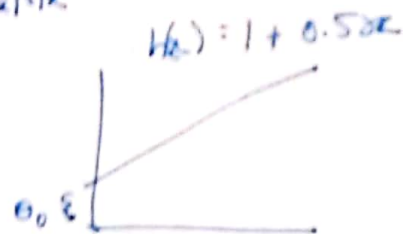
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

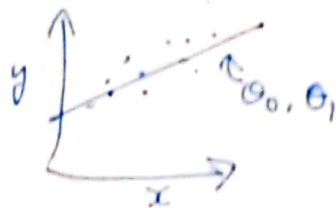
$$\theta_1 = 0.5$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

Aim



Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training example  $(x, y)$

Cost Function: - also called squared error function

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

↑

hypothesis function

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

(prediction)

↑

value

of  $y$

(actual)

∴ We minimise  $J(\theta_0, \theta_1)$   
 $\theta_0, \theta_1$

## Understanding cost function (Intuition I)

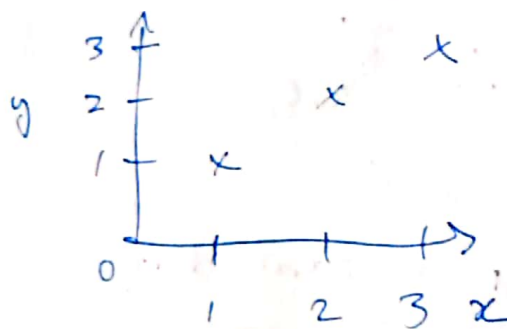
Let  $\theta_0 = 0$

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

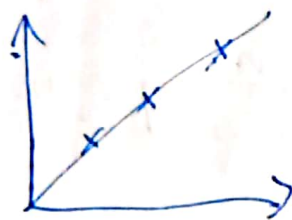
$$\therefore J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m \left( \underset{\substack{\uparrow \\ \theta_1 x^{(i)}}}{h_{\theta}(x^{(i)})} - y^{(i)} \right)^2$$

$\therefore$  We minimise  $J(\theta_1)$   
 ~~$\theta_1$~~

consider given data set

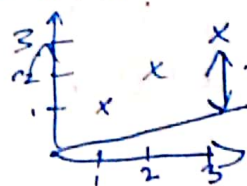


① Let  $\theta_1 = 1$ :



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

② Let  $\theta_1 = 0.5$ :



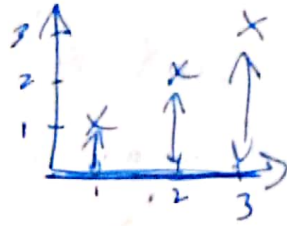
Difference between predicted & actual value  
 $h_{\theta}(x^{(3)}) - y^{(3)}$

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

$\uparrow$   
3

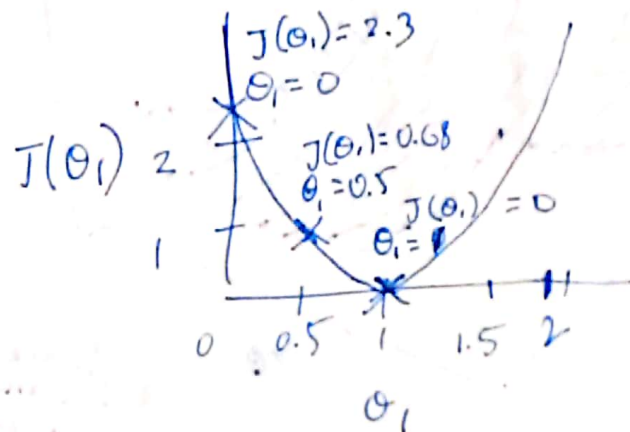
$$= 0.68$$

③ Let  $\theta_1 = 0$ :

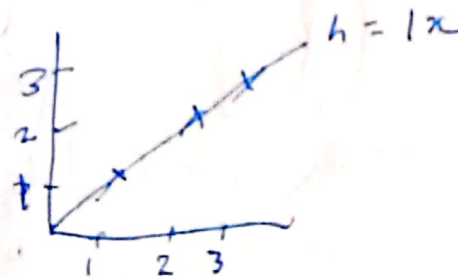


$$J(0) = \frac{1}{2 \times 3} (1^2 + 2^2 + 3^2) = \frac{14}{6} = 2.3$$

By plotting different values of  $\theta_1$ , we get:

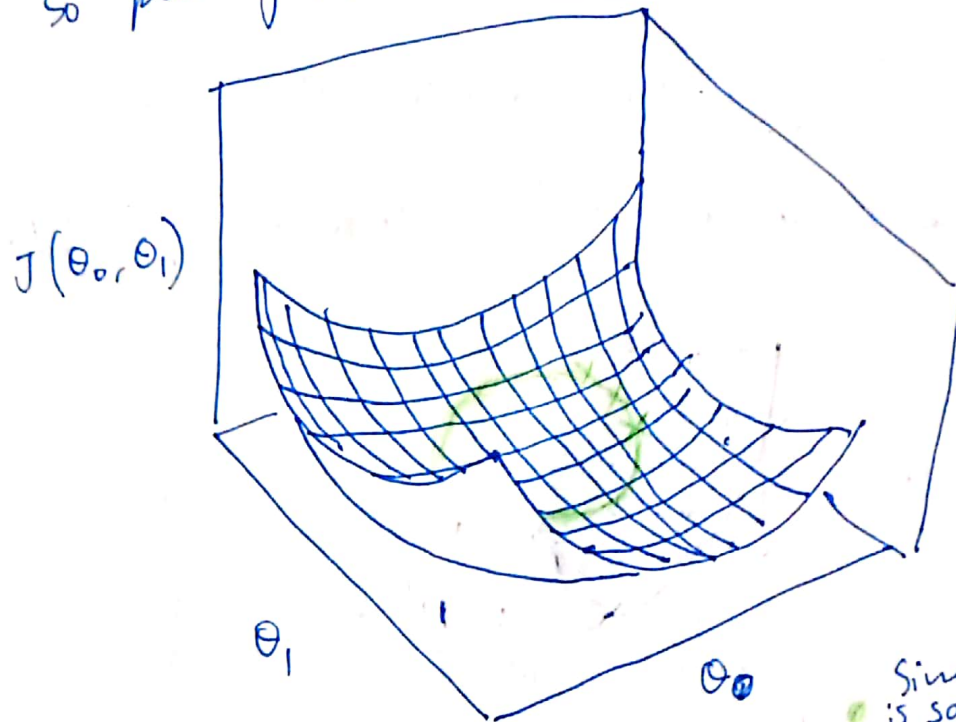


Hence when  $\theta_1 = 1$ , we get minimum value of  $J(\theta_1)$ . This is the best fit of the data



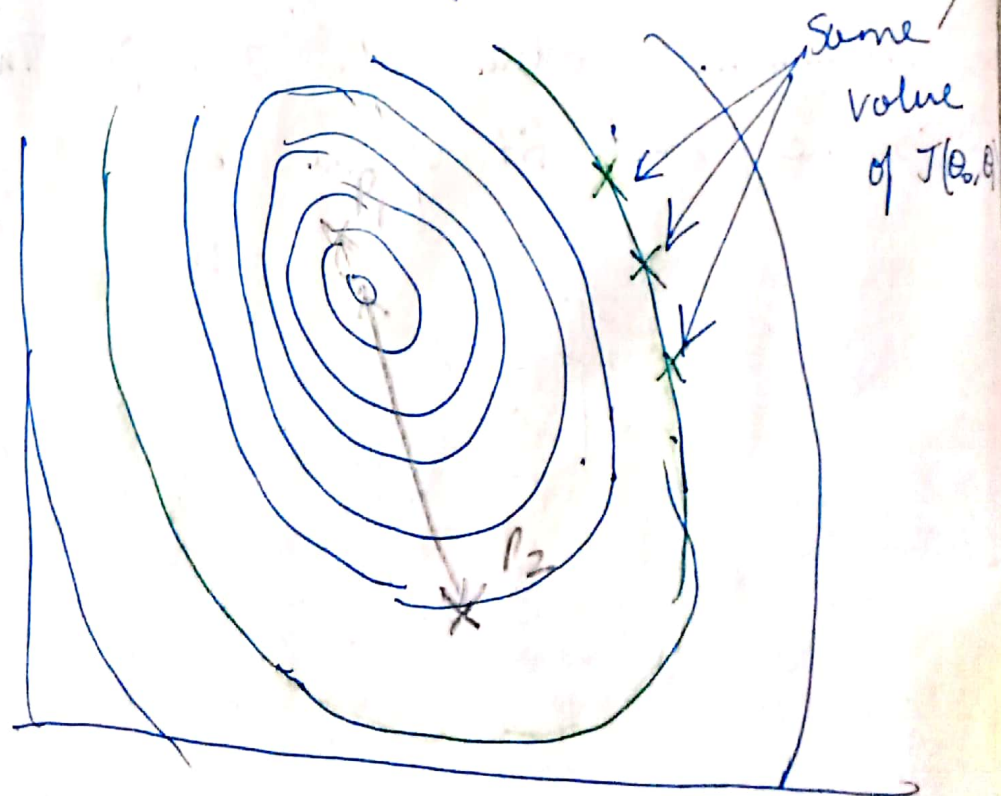
## Intuition II

Now we consider both  $\theta_0$  &  $\theta_1$   
So plotting  $J(\theta_0, \theta_1)$  similar to previous example.



Since height is same along that circle

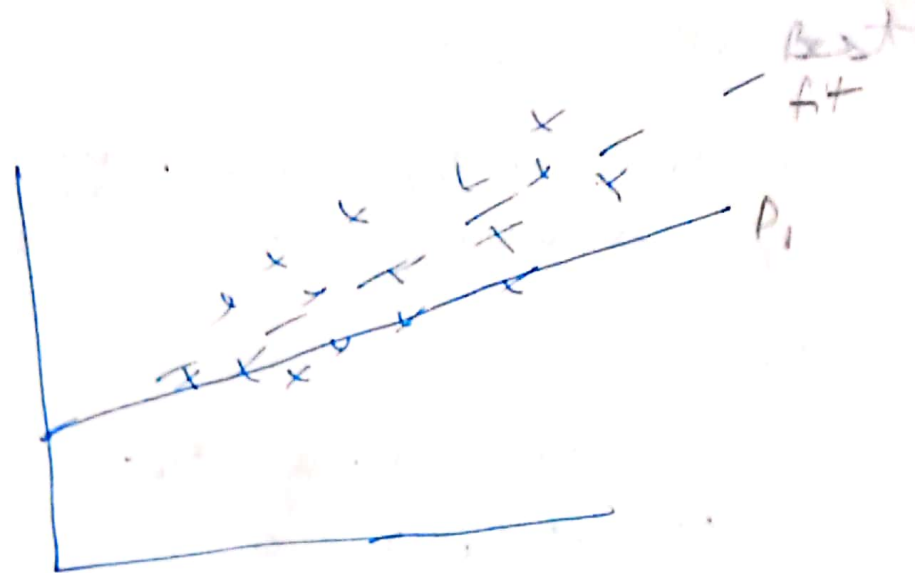
We can represent this in contour plot.



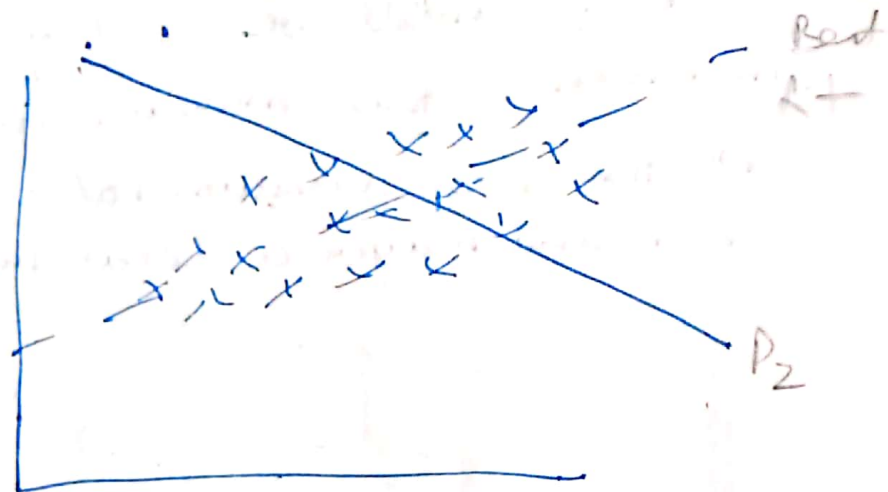
This is the top view of the 3D graph



$P_1$



$P_2$



$\therefore$  Closer to the centre, better the fit. So we need an algorithm that helps us find the centre. Therefore we use GRADIENT DESCENT



## Gradient Descent

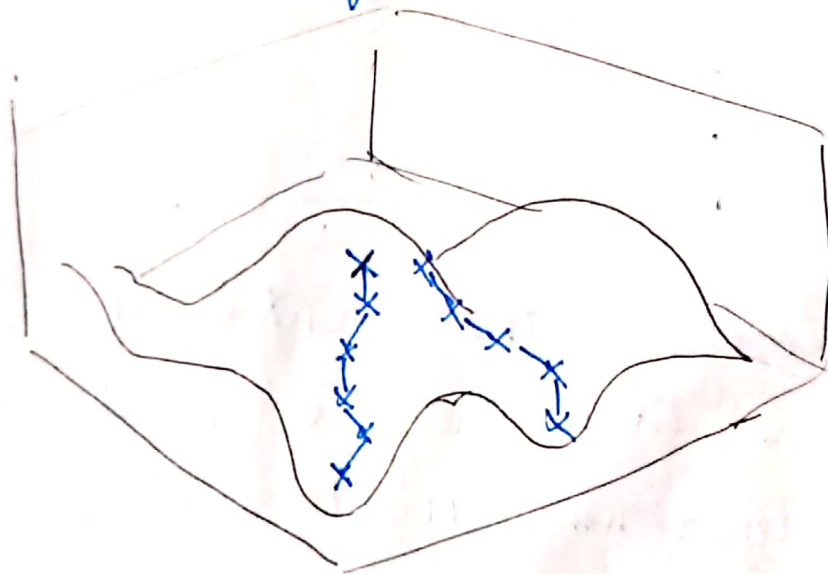
- Can solve function with a number of  $\theta$  parameters  $J(\theta_1, \theta_2, \dots, \theta_n)$

Logic:

- Start with some  $\theta_0, \theta_1$

Ex: let  $\theta_0 = 0$  &  $\theta_1 = 0$

- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we end up at a minimum (this minimum may not be the most minimum value in some cases and maybe a local minimum)



- Start
- Look around, which is lower point?
- Take a small step towards the point of maximum descent

## Gradient Descent Algorithm

When the value starts changing by very little with every iteration

repeat until convergence

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

}

assignment operator

learning rate - the size of the step taken

Here  $j=0$  and  $j=1$

$\therefore$  Simultaneously update  $\theta_0$  &  $\theta_1$

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

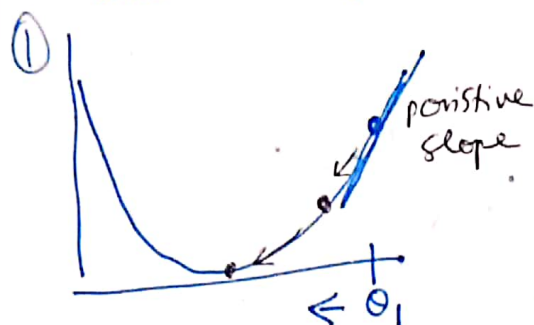
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

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

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

## Intuition

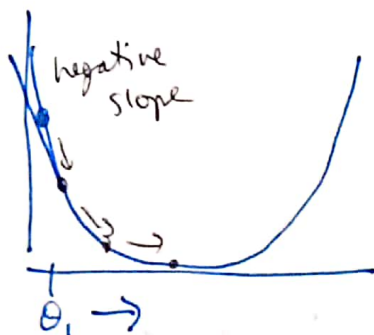
consider only  $\theta_1 \therefore J(\theta_1)$



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

slope

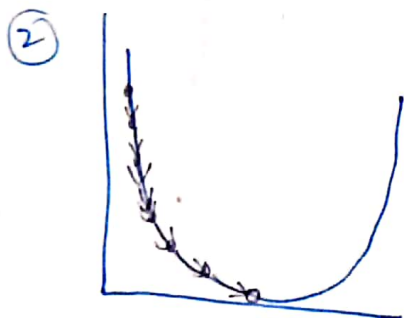
$\therefore \theta_1 = \theta_1 - \alpha (\text{positive value})$   
 $\therefore \theta_1$  is decreasing



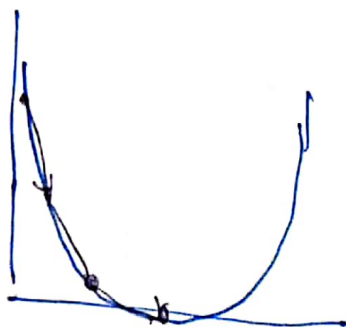
$$\theta_1 = \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$\therefore \theta_1 = \theta_1 - \alpha (\text{negative value})$   
 $= \theta_1 + \alpha (\text{positive value})$   
 $\therefore \theta_1$  is increasing

$\therefore \frac{d}{d\theta_1} J(\theta_1)$  controls where  $\theta_1$  changes (direction of descent)

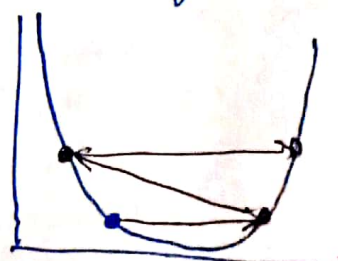


For small  $\alpha$   
 - gradient descent is slow



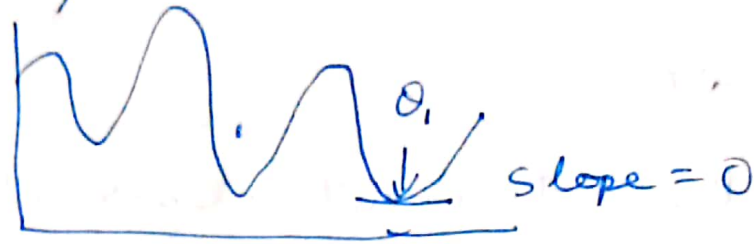
For large  $\alpha$   
 - gradient descent is fast

For too large  $\alpha$



- the steps are too big so it keeps overshooting & never converges

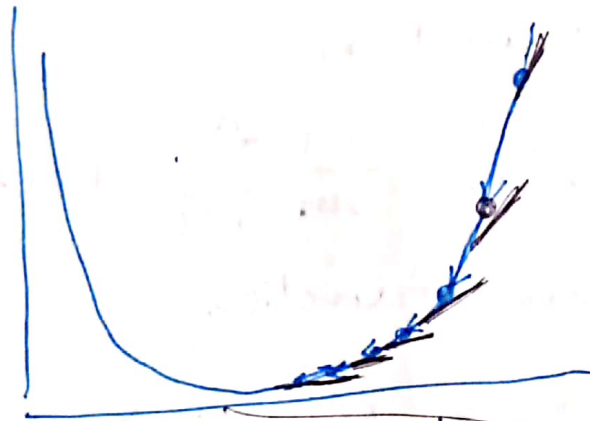
③ If  $\theta_1$  is already at local minimum,



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

$\therefore \theta_1$  doesn't change

④ As gradient descent occurs, the slope gradually decreases, so the step size also decreases



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

• This slope reduces with descent so overall step size decreases as we approach minimum



# Gradient Descent for Linear Regression

So far,

- Gradient Descent algorithm,  
repeat until convergence

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

(for  $j=1 \leftarrow j=0$ )

}

- Cost function

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

- Hypothesis function

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

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

Simplifying this for  $\theta_0$  &  $\theta_1$ ,

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

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

$\therefore$  Gradient descent algorithm for linear regression (univariate)

repeat until convergence  $\{$

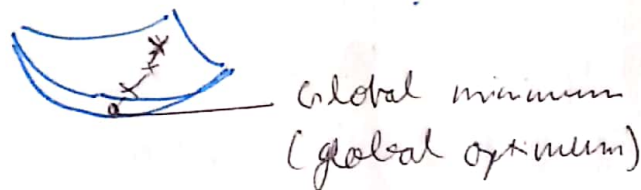
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

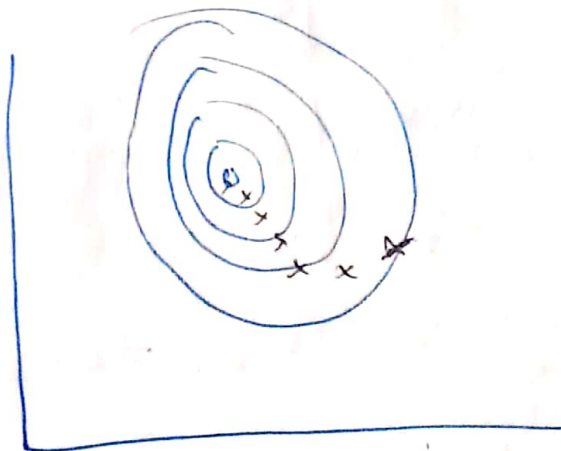
$\}$

Note: update  $\theta_0$  &  $\theta_1$  simultaneously

Note: we had a problem where gradient descent would go to local ~~maximum~~<sup>minimum</sup>, but in linear regression, we get a 'convex function' that has a global ~~max~~ minimum



Example



This is also called Batch gradient descent  
Batch: Each step of gradient descent uses all the training examples (entire batch)