# * Gradient Descent

## Have some function $J(w,b)$

## Want $\min\limits_{w,b} J(w,b)$

## Outline:

(Random initialization)

- Start with some $w, b$   (Set $w=0, b=0$)
- Keep changing $w, b$ to reduce $J(w,b)$
- until we settle at or near a minimum

→ Repeat till Convergence   → learning rate (usually $0 \to 1$)
                              ↳ represents the Cotol of steps of weights update
                                                              x

$$w = w - \alpha \frac{\partial J}{\partial w}$$

$$b = b - \alpha \frac{\partial J}{\partial b}$$

(Simultaneously update w and b)

# * Gradient Descent:

→ general algorithm to optimize nearly any function (find Min.)

→ It is an algorithm to find the parameters $\theta_0, \theta_1$ of the Cost function $J(\theta_0, \theta_1)$ that makes the Cost function as minimum as possible.

→ The idea behind this algorithm is by getting the derivative of the Cost function (tangent line to the function) and it will give us a direction to move towards. We make steps down the Cost function in the direction with the steepest descent.

The gradient descent algorithm: → Repeat until Convergence

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

where  $j = 0, 1$   represents the feature index number

$\qquad\qquad \alpha = $ learning rate (the size of each step)

At each iteration, $(\theta_0, \theta_1)$ must be updated simultaneously as:

$\quad$ tempo $= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

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

$\quad$ then:

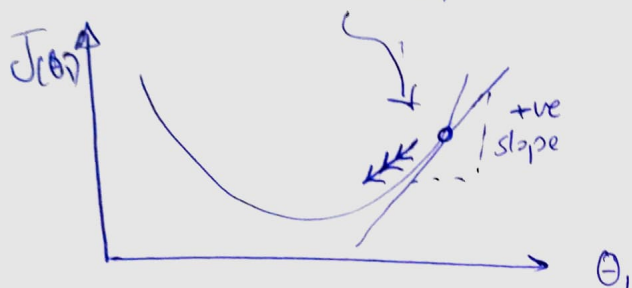$\qquad \theta_0 = $ tempo0
$\qquad \theta_1 = $ temp1

Note
- If $\alpha$ is very small, it causes small steps
- If $\alpha$ is very large, it causes large steps

Let $\theta_0 = 0$ for simplicity:

$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1) \quad [\alpha \text{ is always positive}]$$

If we started at $\theta_1$ larger than that makes $J(\theta_1)$ minimum
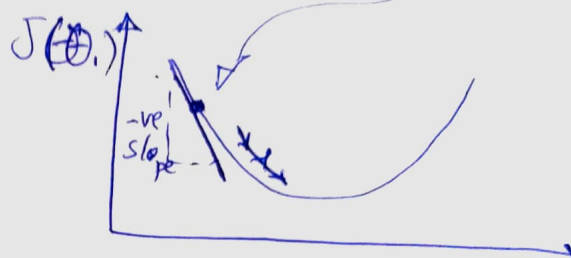


the derivative $\frac{\partial}{\partial \theta_1} J(\theta_1)$, which is the **slope** of tangent line at this point, will be positive

So, $\theta_1 = \theta_1 - (\text{positive value})$

So, $\theta_1$ will decrease

If we started at $\theta_1$ smaller than that make $J(\theta_1)$ minimum



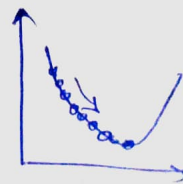the derivative $\frac{\partial}{\partial \theta_1}$ will be negative

So, $\theta_1 = \theta_1 - (\text{negative}) \text{ value}$
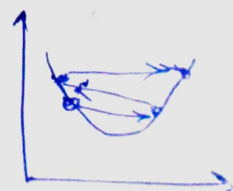
So, $\theta_1$ will increase

So, $\theta_1$ will eventually Converges to the minimum value of $J(\theta_1)$

Note we should adjust out parameter ($\alpha$) to ensure that the gradient descent will Converge in a reasonable time

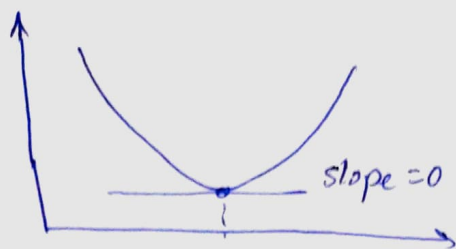- If $\alpha$ is too small: gradient descent can be slow



- If $\alpha$ is too large: gradient descent can overshoot the minimum. It may fail to Converge, or may even diverge

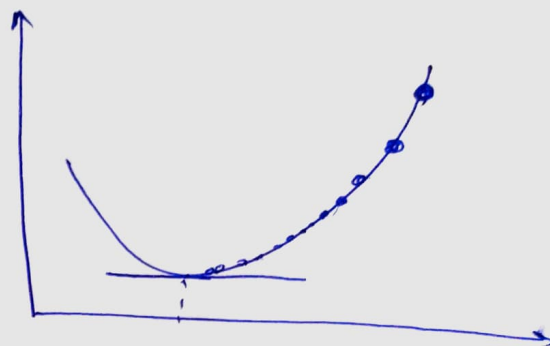→ What if we started at $\overline{\theta}(\theta_1)$ is minimum already?



slope = 0

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

so, $\theta_1$ will be unchanged

→ The steps will automatically be smaller because the slope decreases till it equals to zero



• Gradient Descent for Linear Regression :

**Gradient Descent Algorithm**

repeat until convergence {

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

for $j = 0, 1$

}

**Linear Regression Model**

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

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

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

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

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

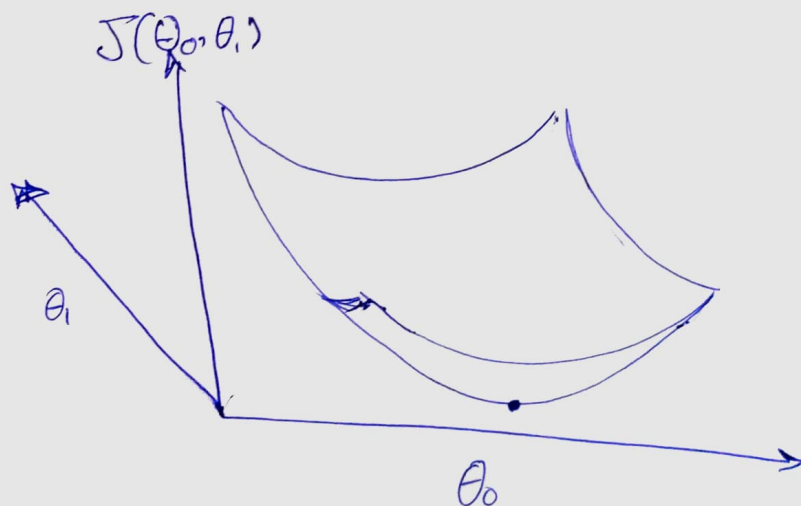∴ So, Gradient descent for linear regression:

repeat until convergence {

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

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

} update $\theta_0, \theta_1$ Simulataneously

• The point of all this that if we start with a guess for our hypothesis (a guess of $\theta_0, \theta_1$) and then we repeatedly apply these gradient descent equations, our hypothesis will become more and more accurate.



Convex function

↓

one global minimum

→ This gradient descent called "Batch" Gradient Descent as Each step of gradient descent uses all the training examples

→ There are other versions of gradient descent (Stochastic, — mini-batch)