6] Gradient descent

*) Gradient descent is an optimization algorithm used to minimize a function

*) In machine learning, the function is usually a loss (or cost) function that measures how far off the model's predictions are from the actual data

$$J(\omega, b) \longrightarrow \text{want} \quad \min_{\omega, b} (J(\omega, b))$$

*) start with some $\omega, b$ (set $\omega = 0, b = 0$)
*) keep changing $\omega, b$ to reduce $J(\omega, b)$
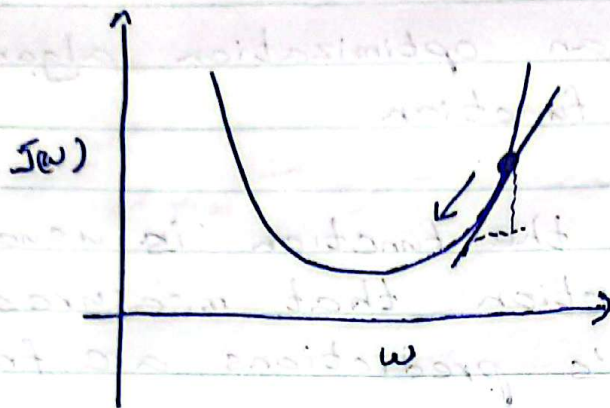*) Until we settle at or near a minimum

◎ If it has more than one minimum it is not linear regression (not squared error cost function)

$$\omega = \omega - \alpha \boxed{\frac{\partial}{\partial \omega} J(\omega, b)}$$

$$b = b - \alpha \boxed{\frac{\partial}{\partial b} J(\omega, b)}$$

$\alpha$ = learning rate
partial derivative

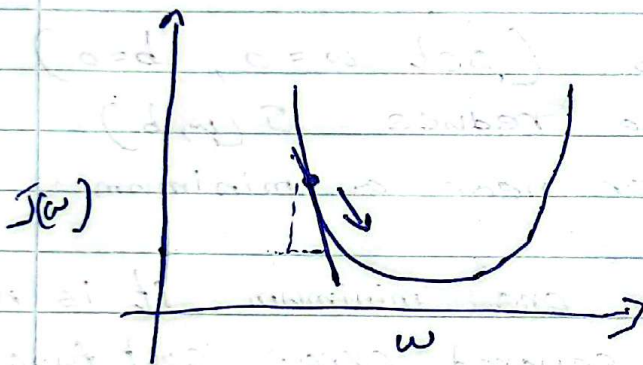*) Simultaneously update $\omega$ and $b$ until Convergence.

Simple case → $J(\omega)$



$$\omega = \omega - a\left(\boxed{\frac{\partial}{\partial \omega} J(\omega)}\right)$$

$$> 0$$

$$\omega = \omega - a \cdot (\text{positive number})$$

~~Derivative~~ decrease
positive slope
positive derivative



$$\omega = \omega - a\left(\boxed{\frac{\partial}{\partial} J\omega}\right)$$

$$< 0$$

$$\omega = \omega - a \, (\text{negative number})$$

$\omega$ = increases
negative slope
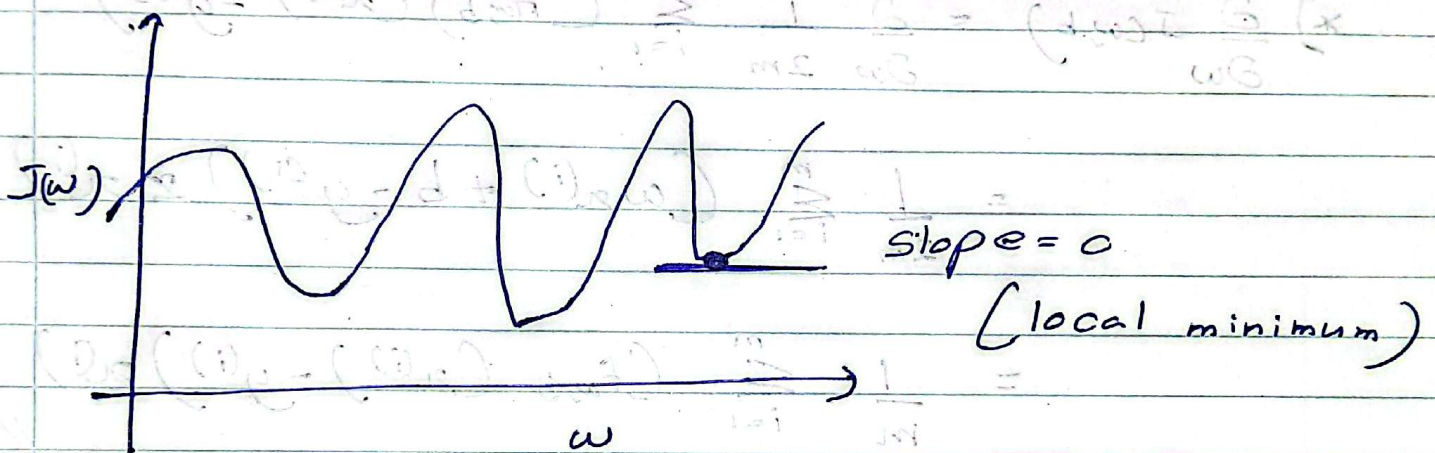negative derivative

7] Learning rate $(a)$

i) If $a$ is too large :
   *) Gradient descent overshoot, never
      reach minimum
   *) Fail to converge, diverge

ii) If a is too small:
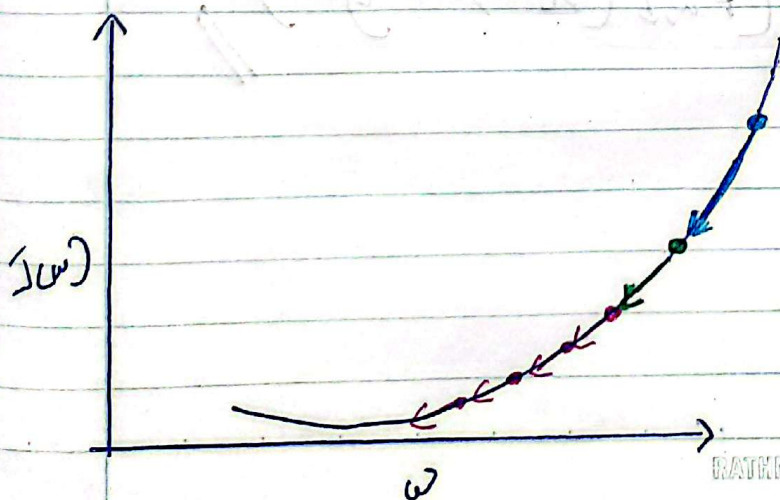
   *) Gradient descent may be slow

Local minimum

*) No matter (what) learning rate you use it will change nothing

$$J(\omega)$$

Slope = 0

(local minimum)

$$\omega$$

$$\omega = \omega - a \left[ \frac{d}{d\omega} (J(\omega)) \right] = 0$$

$$\boxed{\omega = \omega}$$

Fixed learning rate

smaller
(fixed)

not as large

$$\omega = \omega - a \; \frac{d}{d\omega} (J\omega)$$

large

$$J(\omega)$$

Near a local minimum:

*) Derivative is smaller

*) Update steps became smaller

$$\omega$$

RATHNA

# 8) Gradient descent for linear regression

$$f_{w,b}(x) = wx + b$$

$$J(w,b) = \frac{1}{2m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

*) $$\frac{\partial}{\partial w} J(w,b) = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \left( wx^{(i)} + b - y^{(i)} \right) 2\, x^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right) x^{(i)} \quad //$$

*) $$\frac{\partial}{\partial b} J(w,b) = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^{m} \left( wx^{(i)} + b - y^{(i)} \right)^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \left( wx^{(i)} + b - y^{(i)} \right) 2$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right) \quad //$$

# Gradient descent algorithm

repeat until convergence {

$$w = w - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right)$$

}

*) For squared error cost it should only have one global minimum ( ∵ Convex function /bowl )

Batch gradient descent

*) Batch = Each step of gradient descent uses all the training examples

*) There are some methods other than batch gradient descent which use only a subset of training examples.

RATHNA