Numerical Optimization for ML & DS.

Mansoura, AI46, Session 2 ——

1) Review of ~~Probal~~ differentiation.

⇒ Gradient $\nabla$ → del / nabla

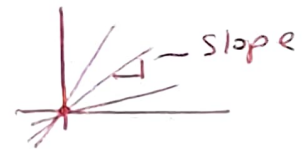2) Contour plot ⟵

⇒ Objective function (Review)

3) Gradient of a cost function

(Linear regression / Logistic Reg.)

↳ 1 - parameter
↳ 2 - parameter
↳ ⋮
↳ n - parameter

— slope

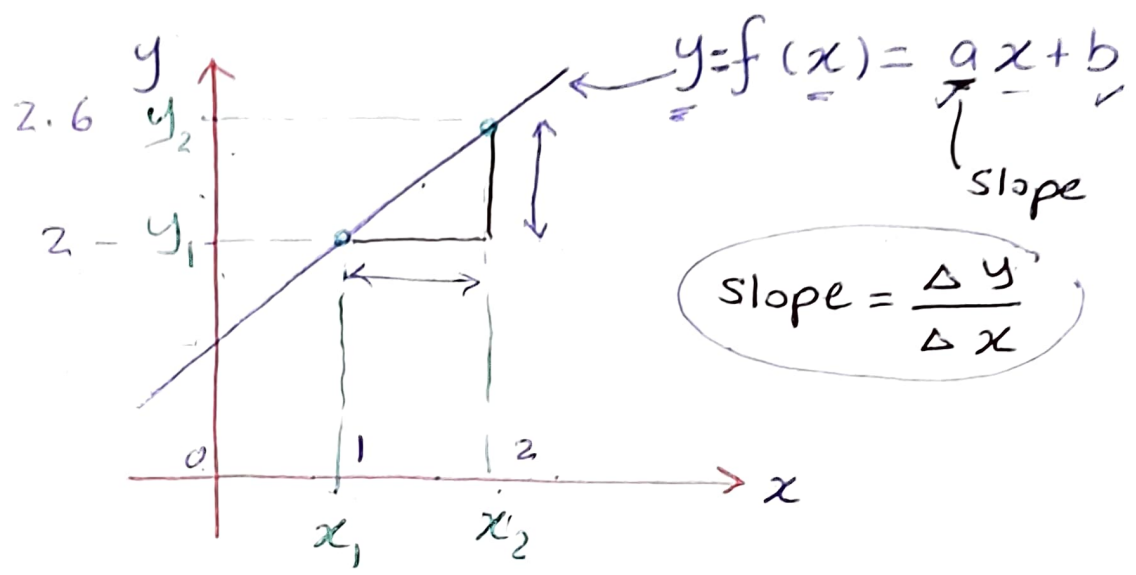4) Problems with "Vanilla Gradient Descent" algorithm.

5) feature scaling

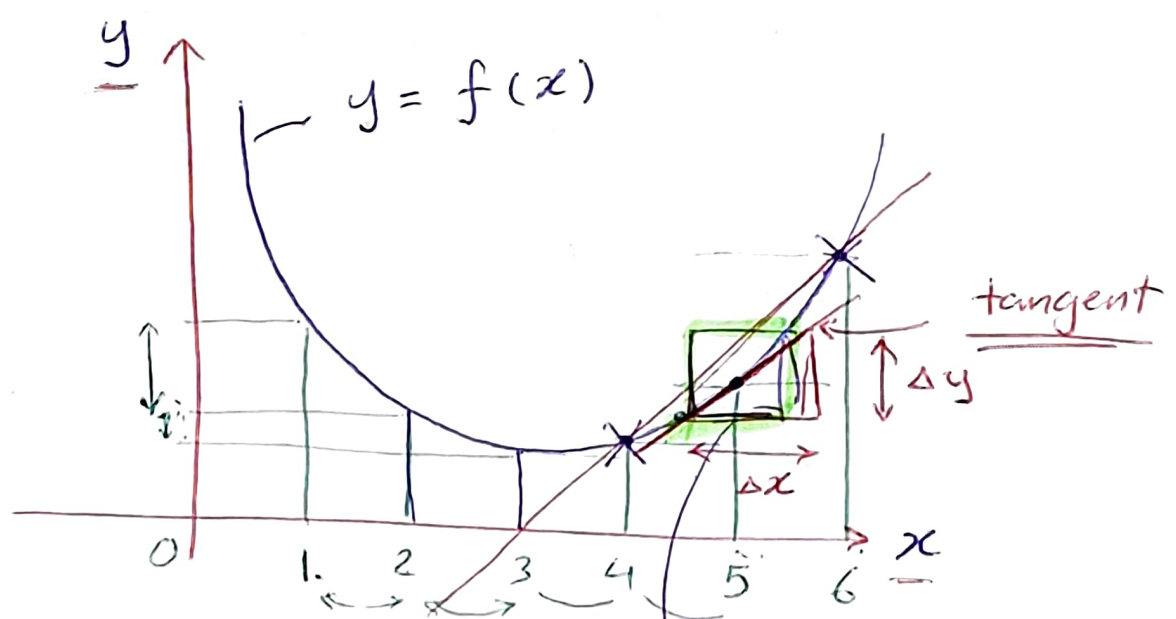↓

6) Variants of GD algorithm.
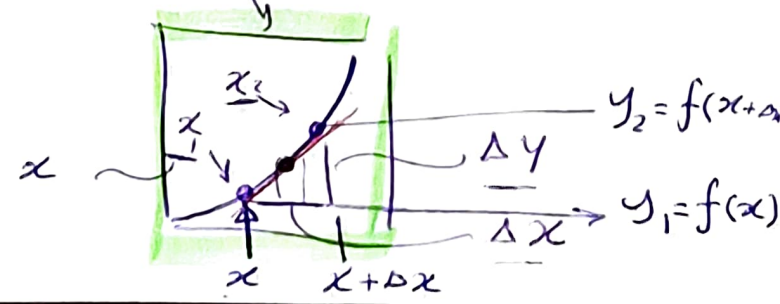
~ derivative of a function

"defferintiation"



$y = f(x) = ax + b$

slope

$$\text{slope} = \frac{\Delta y}{\Delta x}$$

ex.   slope $= \dfrac{y_2 - y_1}{x_2 - x_1} = \dfrac{2.6 - 2}{2 - 1} = \dfrac{0.6}{1} = 0.6$

rate of change $= 0.6$

↳ of y value w.r.t. x value



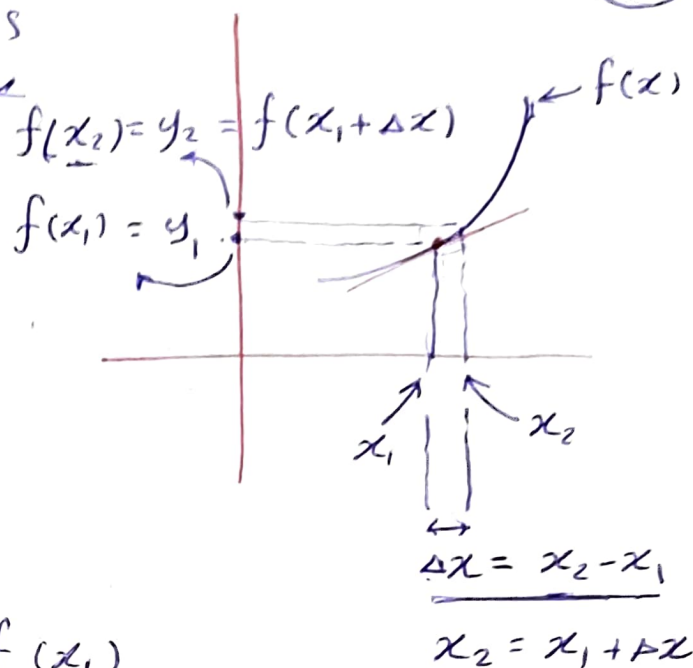$y = f(x)$

tangent

$$\left( \frac{\Delta y}{\Delta x} \right)_{\Delta x \to 0}$$

$\Delta x \to 0$ ;  $dx$

$y_2 = f(x + \Delta x)$

$y_1 = f(x)$

$x + \Delta x$

→ for continuous functions

$$\frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1}$$

$$= \frac{f(x_2) - f(x_1)}{\Delta x}$$

$$f(x_2) = y_2 = f(x_1 + \Delta x)$$

$$f(x_1) = y_1$$

← f(x)

$x_1$    $x_2$

$\Delta x = x_2 - x_1$

$x_2 = x_1 + \Delta x$
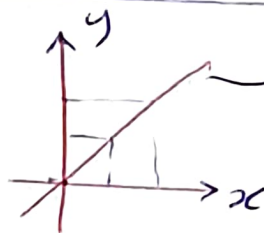
$$\text{slope}\Big|_{x = x_1} \approx \frac{f(x_1 + \Delta x) - f(x_1)}{\Delta x}$$

$$\frac{dy}{dx} = \lim_{\Delta x \to 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

$$dy \stackrel{\circ}{=} f(x + \Delta x) - f(x), \; \Delta x \to 0$$

$$dx \stackrel{\circ}{=} \Delta x \quad, \; \Delta x \to 0$$

---

$$y = x$$

$$\frac{dy}{dx} = 1$$

$$y = f(x) = x$$

$$y = a x$$

$$\frac{dy}{dx} = a$$

$$y = a$$

$$\frac{dy}{dx} = 0$$

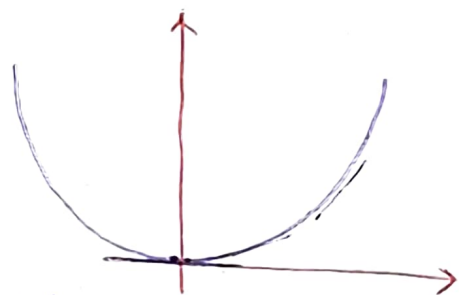| $y = f(x)$ | $\dfrac{df(x)}{dx} \equiv \dfrac{dy}{dx}$ |
|---|---|
| $y = $ const. | $\dfrac{dy}{dx} = 0$ |
| $y = ax$ | $\dfrac{dy}{dx} = a$ |
| $y = ax + b$ | $\dfrac{dy}{dx} = \dfrac{d(ax)}{dx} + \dfrac{d(b)}{dx}$ <br> $\dfrac{dy}{dx} = a$ |
| $y = ax^2$ | $\dfrac{dy}{dx} = a(2x)$ |
| $y = b.x^K$ | $\dfrac{dy}{dx} = b.K\, x^{K-1}$ |
| $y = f(g(x))$ | $\dfrac{dy}{dx} \equiv \dfrac{df}{dx} \cdot \dfrac{dg(x)}{dx}$ |
| ~~$y = f(ax)$~~ | |
| $y = (ax+b)^2$ | $\dfrac{dy}{dx} = 2(ax+b) \cdot \dfrac{d(ax+b)}{dx}$ |
| $y = |x|$ | |

Piecewise continuous

derivative is not defined ← 

← discontinuity
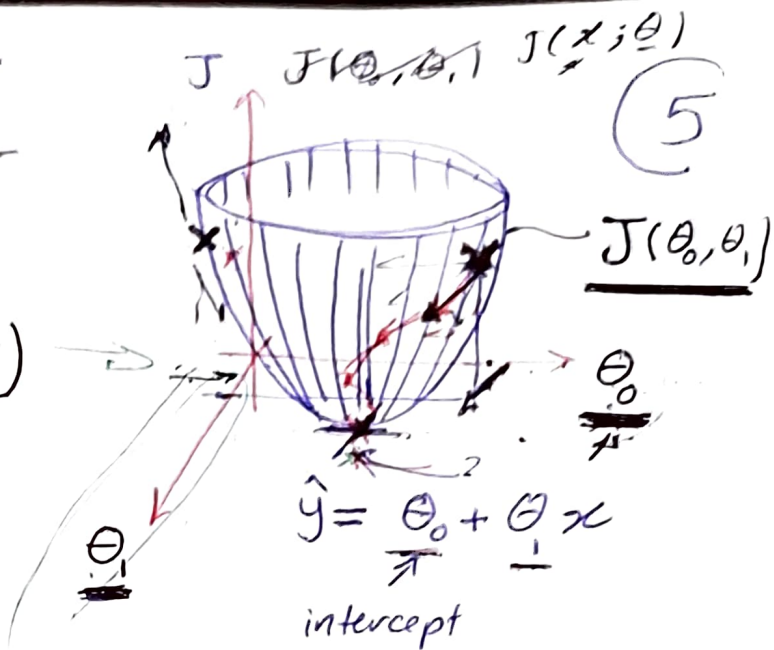
# Partial Derivatives
## ~~Gradient~~ $\nabla J$

Partial derivatives

$$\frac{\partial J}{\partial \theta_0}$$

$$\frac{\partial J}{\partial \theta_1}$$
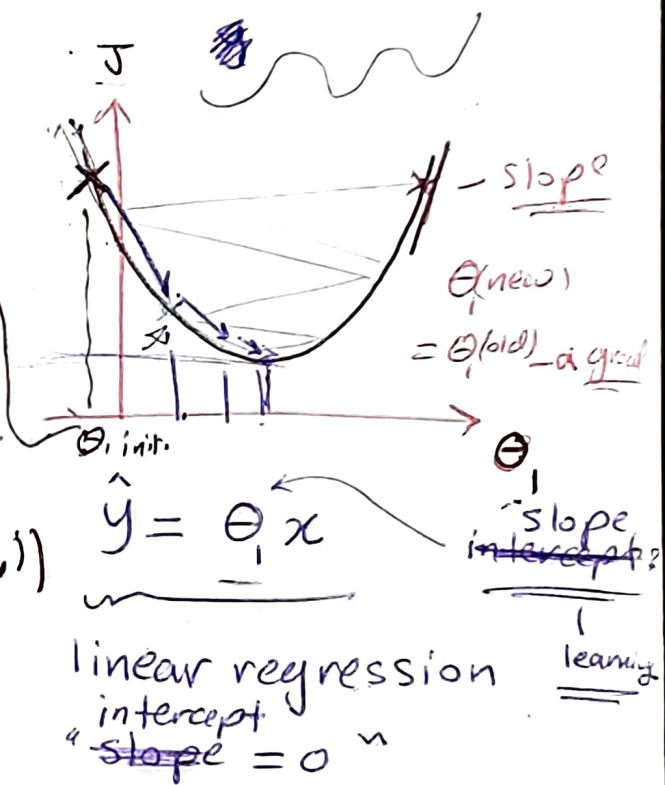
$$J(\theta_0, \theta_1)$$

$J \qquad J(\theta_0, \theta_1) \qquad J(x; \theta)$

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

$$\hat{y} = \theta_0 + \theta_1 x$$

intercept

---

→ initialize $\theta_1$ (initial)

→ find gradient at $\theta_1$
of $J(\theta_1)$

→ update ~~$\theta_1^{(new)} = \theta_1^{(old)} \text{ tangently}$~~

$$\theta_1^{(new)} = \theta_1^{(old)} - \alpha \text{ gradient} (J(\theta_{old}))$$

learning rate

→ repeat until convergence

$J$

- slope

$\theta_1^{(new)} = \theta_1^{(old)} - \alpha$ grad

$\theta_1$ init.

$$\hat{y} = \theta_1 x$$

$\theta_1$

"slope intercept?"

linear regression
intercept
"~~slope~~ = 0"

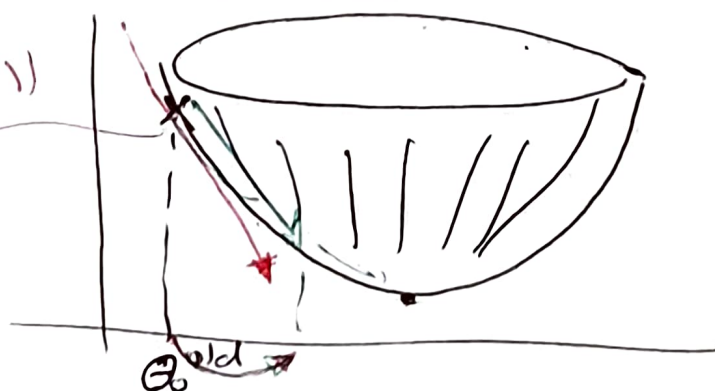learning

---

$$\theta_{new} = \theta_{old} + \text{positive value}$$

$$\theta_{new} = \theta_0^{old} - \alpha \text{ gradient} (J(\theta_0^{old}_1))$$

$\rightarrow$ -ve.

grad $(J(\theta_0^{old}))$
= -ve, large

$\theta_0^{old}$
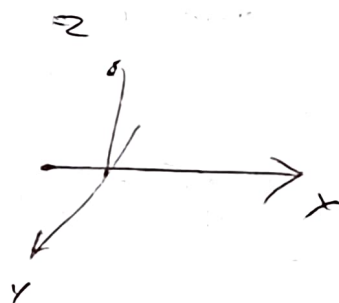
$\nabla$ : Del, Nabla operator

$$\cdot\nabla = \begin{bmatrix} \dfrac{\partial}{\partial\Theta_0} \\[6pt] \dfrac{\partial}{\partial\Theta_1} \\[6pt] \vdots \\[6pt] \dfrac{\partial}{\partial\Theta_n} \end{bmatrix}$$

---

e.g., 3D problem

$$f(x,y,z) = 2x + 3y^2 - z^3$$

$$\frac{\partial f}{\partial x}, \quad \frac{\partial f}{\partial y}, \quad \frac{\partial f}{\partial z}$$

$$\nabla f \equiv \text{vector}$$

gradient → "Vector" | Scalar function

$$\vec{\nabla} = \begin{bmatrix} \partial/\partial x \\ \partial/\partial y \\ \partial/\partial z \end{bmatrix} \Rightarrow \nabla f = \begin{bmatrix} \partial f/\partial x \\ \partial f/\partial y \\ \partial f/\partial z \end{bmatrix}$$

gradient (f)

grad (f)

$$\vec{\nabla} f = \begin{bmatrix} \partial/\partial x \, (f) \\ \partial/\partial y \, (f) \\ \partial/\partial z \, (f) \end{bmatrix}$$

gradient

$$f = 2x + 3y^2 - z^3$$

$$= \begin{bmatrix} \dfrac{\partial}{\partial x} (2x + 3y^2 - z^3) \\ \dfrac{\partial}{\partial y} ( \qquad ) \\ \partial/\partial z \, ( \qquad ) \end{bmatrix}$$
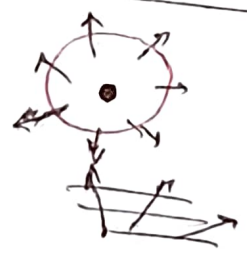
$$\vec{\nabla}(f(x,y,z)) = \begin{bmatrix} 2 \\ 6y \\ -3z^3 \end{bmatrix}$$

gradient

note

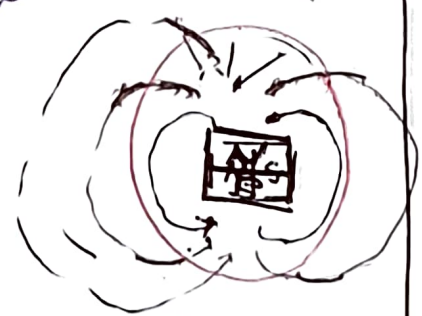$$\vec{\nabla} \cdot \vec{v} = \text{scalar}$$

divergence of $\vec{v}$

cross product

$$\vec{\nabla} \times \vec{u} = \text{vector}$$
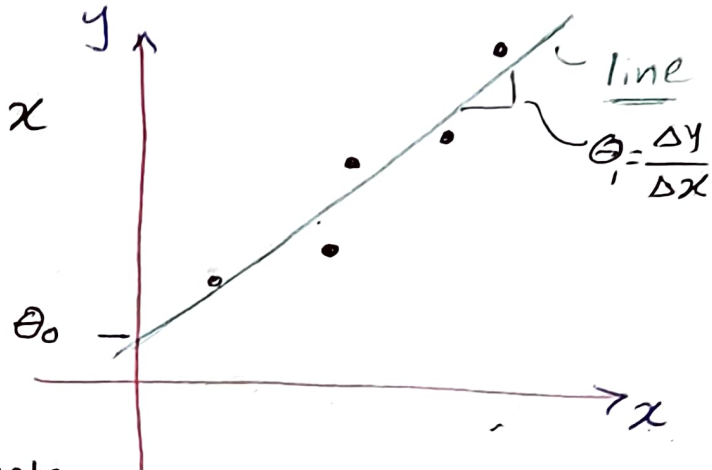
curl of $\vec{u}$

magnetic monopoles?

are they used in ML?

→ Linear Regression    (two parameters)

$$f(x; \theta_0, \theta_1) \equiv h_\theta(x) = \hat{y} = \theta_0 + \theta_1 x$$

hypothesis
$\underset{\sim \text{model}}{}$

predicted y



line

$\theta_1 = \frac{\Delta y}{\Delta x}$

$\theta_0$

→ loss: $\ell(\cdot)$ ;  for a single ~~certain~~ example  (1-data points)

→ cost; $\sum\limits_{i=1}^{m} \ell(\cdot)$  $=$  $\underline{J(\theta)}$

— sum of losses  ↳?

objective :   minimize $J(\theta)$
$\underline{\theta_0, \theta_1}$ ~?

→ let  ;   MSE

m-data
points

m-examples

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum\limits_{i=1}^{m} (\hat{y}_i - y_i)^2$$

?

RSS

$$= \frac{1}{2m} \sum\limits_{i=1}^{m} (y_i - \hat{y}_i)^2$$

$$\frac{1}{m} \sqrt{\sum\limits_{i=1}^{m} (y_i - \hat{y}_i)^2} = (\| \vec{y} - \vec{\hat{y}} \|)^2 / m$$

ex

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\vec{\hat{y}} = \theta_0 + \theta_1 \vec{x}$$

$$\vec{\Theta} := \vec{\Theta} - \alpha \vec{\nabla}(J(\vec{\Theta}))$$

<u>update</u>

learning rate

"or $\eta$"

$$\vec{\Theta} = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_n \end{bmatrix}$$

$$\vec{\Theta}_{\xi}^{(new)} = \vec{\Theta}^{(old)} - \alpha \vec{\nabla}(J(\vec{\Theta}^{(old)}))$$
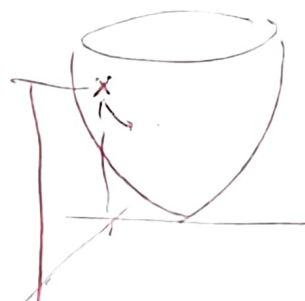
repeat until convergence

$$\vec{\Theta}(K+1) = \vec{\Theta}(K) - \alpha \vec{\nabla}(J(\Theta(K)))$$

time index / step index

$$\Theta_0^{new} := \Theta_0^{old} - \alpha \frac{\partial}{\partial \Theta_0} J(\Theta^{old})$$

and

$$\Theta_1^{new} := \Theta_1^{old} - \alpha \left\{ \frac{\partial}{\partial \Theta_1} J(\Theta^{old}) \right.$$

↳ all parameters $(\Theta_0, \Theta_1, \dots, \Theta_n)$ are updated <u>simultaneously</u> !

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

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

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

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

if $\quad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}_i - y_i)^2$

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

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{m} (\theta_0 + \theta_1 x_i - y_i)(+1)$$

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

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

$$\nabla J(\theta) = \begin{bmatrix} \dfrac{\partial J(\theta)}{\partial \theta_0} \\[2mm] \dfrac{\partial J(\theta)}{\partial \theta_1} \end{bmatrix}$$
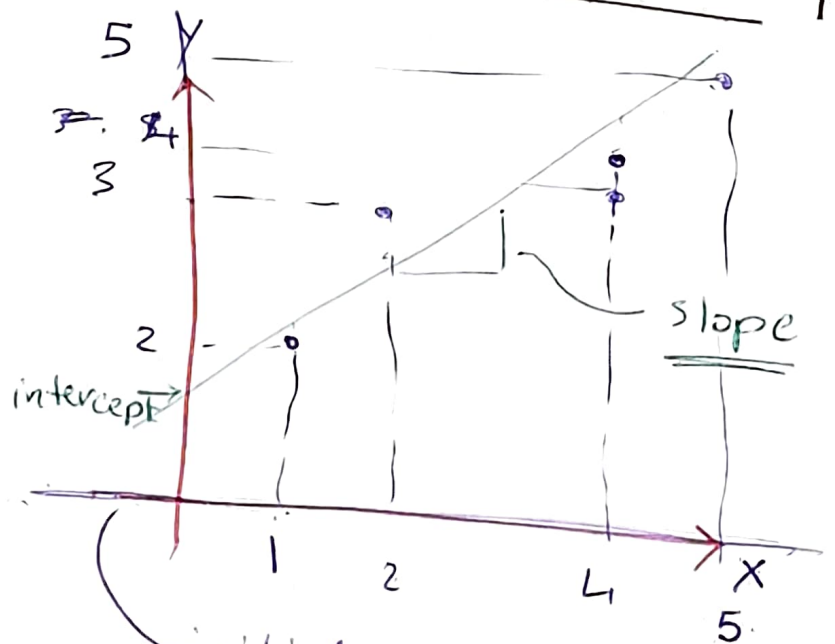
| $i$ | $x$ | $y$ |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 2 | 3 |
| 3 | 4 | 4 |
| 4 | 5 | 5 |

$\boxed{m = 4}$ : 4 - examples,



initial

slope = 0 = $\theta_1$

intercept = 0 = $\theta_0$

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \qquad \vec{\theta}_{(init)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\hat{\vec{y}} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \qquad \vec{y} = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$$

$$\overrightarrow{error} : \begin{bmatrix} -2 \\ -3 \\ -4 \\ -5 \end{bmatrix}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$$
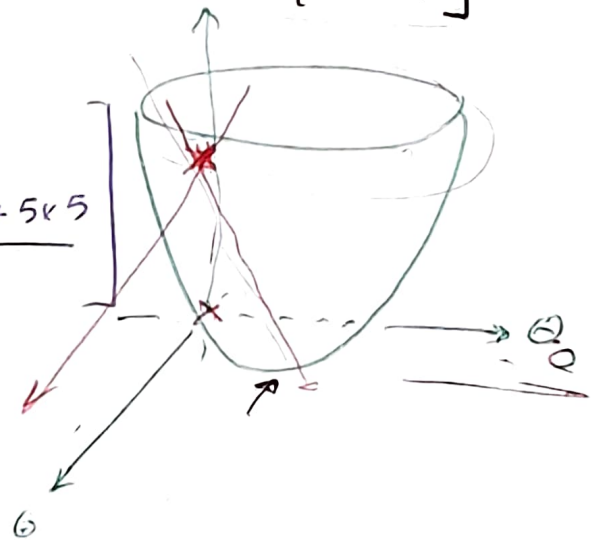
$$= \frac{1}{8}$$

$$\nabla J(\theta) = \begin{bmatrix} \partial J / \partial \theta_0 \\ \partial J / \partial \theta_1 \end{bmatrix}$$

for old $\theta$

$$\nabla J(\theta) \Bigg|_{\substack{\text{evaluated} \\ \text{at} \\ \text{old } \theta}} = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^{m} (\overset{0}{\theta_0} + \overset{0}{\theta_1} x_i - y_i) \\ \frac{1}{m} \sum_{i=1}^{m} (\overset{0}{\theta_0} + \overset{c}{\theta_1} x_i - y_i) x_i \end{bmatrix}$$

$$= \begin{bmatrix} -\frac{14}{4} \\ -\frac{1 \times 2 + 2 \times 3 + 4 \times 4 + 5 \times 5}{4} \end{bmatrix}$$

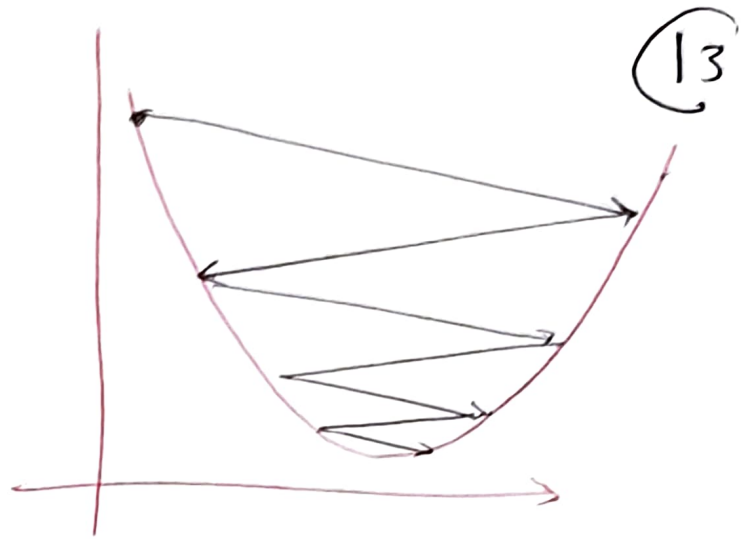$$= \begin{bmatrix} -14/4 \\ -49/4 \end{bmatrix}$$



$$\frac{\partial J}{\partial \theta_0} = -3.5 \qquad \frac{\partial J}{\partial \theta_1} = -12.25$$

$\Rightarrow$ update $\vec{\theta}$ ; update simultaneously $\theta_0, \theta_1$

✓ $\theta_0^{(new)} = \theta_0^{(old)} - \alpha(-3.5) = 0.35$ │ let $\alpha = 0.1$

✓ $\theta_1^{(new)} = \theta_1^{(old)} - \alpha(-12.25) = 1.225$

# learning rate $\alpha$ ↑



# leaning rate too big

## (diverge)



# learning rate too small

error



time / step / iteration

## dataset $\qquad$ m × n matrix

m rows: # of samples

n columns: # of features

area

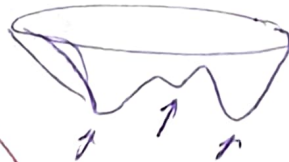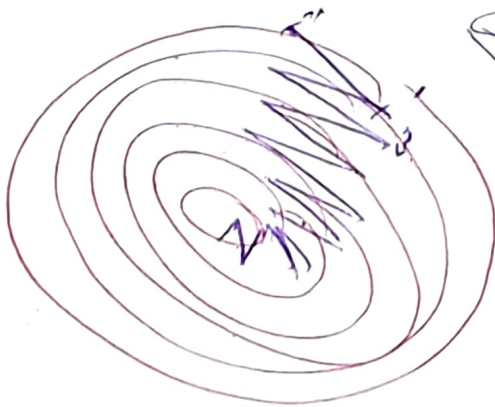| | $x_0$ | $x_1$ | $x_2$ | $x_3 \cdots x_n$ | y |
|---|---|---|---|---|---|
| i=1 | 1 | 100 | 1 | ✓ | --- |
| i=2 | 1 | 90 | 2 | | |
| : | : | 150 | 3 | | |
| : | : | 120 | 2 | | |
| : | : | | | | |
| i=m | 1 | | | | |

## linear regression:

intercept

$$\hat{y} = h_\Theta(x_1, \cdots, x_n) = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 + \cdots + \Theta_n x_n$$

n+1 parameters

$$\hat{y}_i = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \vdots \\ \Theta_n \end{bmatrix} \cdot \begin{bmatrix} x_{i0} \; 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{bmatrix}$$

add one extra feature

## feature scaling : ⟶ ☰

## contour plot