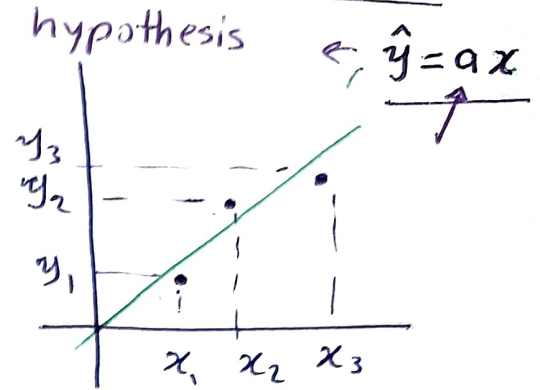
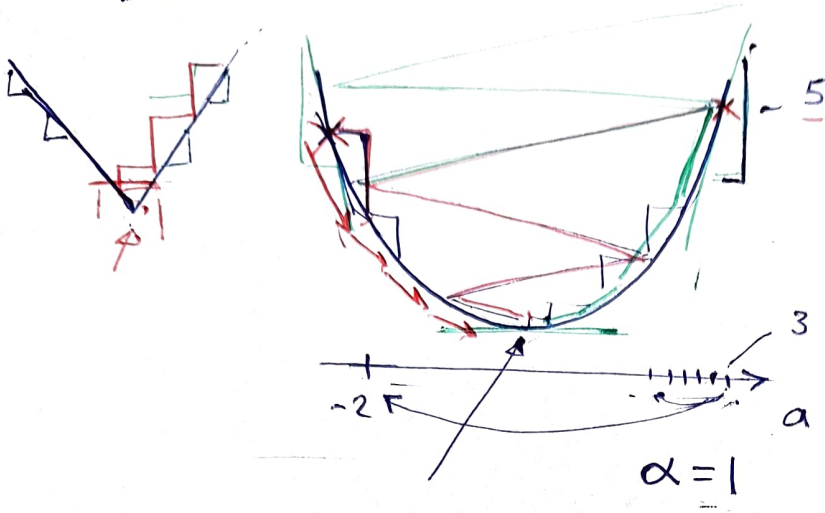


①

# Optimization "Mansoura" session 2

AI45

21/1/2025



$$a_{\text{new}} = a_{\text{old}} - \alpha \times \text{gradient}$$

learning rate

$$a_{\text{new}} = a_{\text{old}} - 1 \times 5$$

- cost function
- objective
- loss

hypothesis

$$h(x) = \hat{y} = ax + b$$

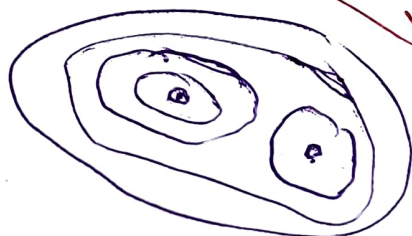
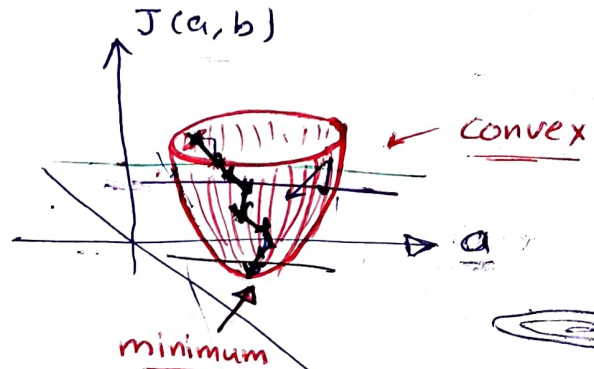
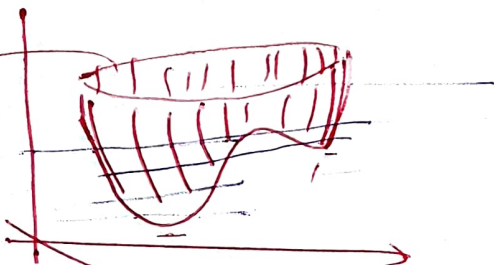
slope intercept

$$10,000x + 100,000$$

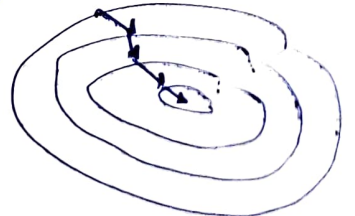
$$\text{MSE} \Rightarrow \frac{1}{m} \sum_i (y_i - \hat{y}_i)^2 = \frac{1}{m} \sum_{i=1}^m (y_i - (ax_i + b))^2 = J(a, b)$$

data points  $y_i$   $\hat{y}_i$   $h_i(a, b)$

non-convex



contour plot



# Gradient



$$\frac{1}{2m} \sum_{i=1}^m (y_i - (ax_i + b))^2$$

(2)

$$\frac{\partial}{\partial a} J(a, b) = \frac{1}{2m} \sum_{i=1}^m 2(y_i - (ax_i + b))(-x_i)$$

$$\frac{\partial}{\partial b} J(a, b) = \frac{1}{2m} \sum_{i=1}^m 2(y_i - (ax_i + b))(-1)$$

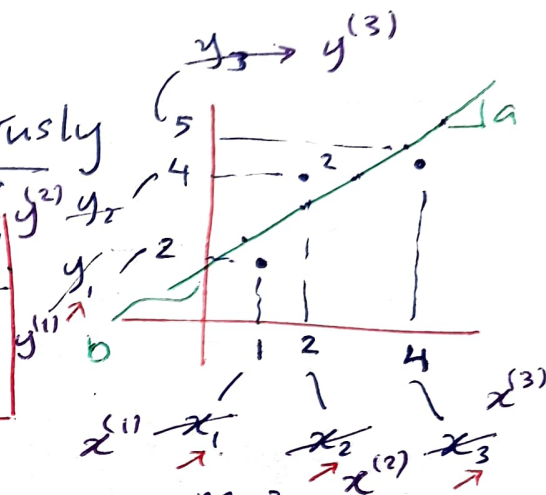
update

a, b

simultaneously

$$a_{\text{new}} = a_{\text{old}} - \alpha \text{grad}$$

$$b_{\text{new}} = b_{\text{old}} - \alpha \text{grad}$$



Repeat

until

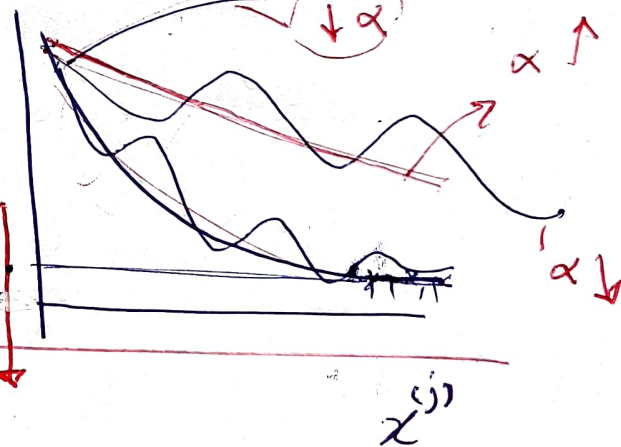
convergence

e.g., 10 → times  
100 →  
1000 →

convergence

and error

next, we consider fitting a model with n-parameters, dataset with m-point



## Multivariate linear regression

$$h(\theta_0) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

↑  
intercept / bias

$$h(\theta_0, \dots, \theta_n) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

↑  
area # of rooms  
↑  
floor #

ex  
area: 1st feature  
=  $\theta_1$   
area of apartment #2

→  $x_i$  = variable #  $i$  =  $i^{\text{th}}$  variable

→  $x_i^{(j)}$  ⇒  $i^{\text{th}}$  variable, data point #  $j$

ex. price

	$x_0$	$x_1$	$x_2$	...	$x_n$
→ 1	1	150	4	...	3
→ 2	0.8	120	3	...	4

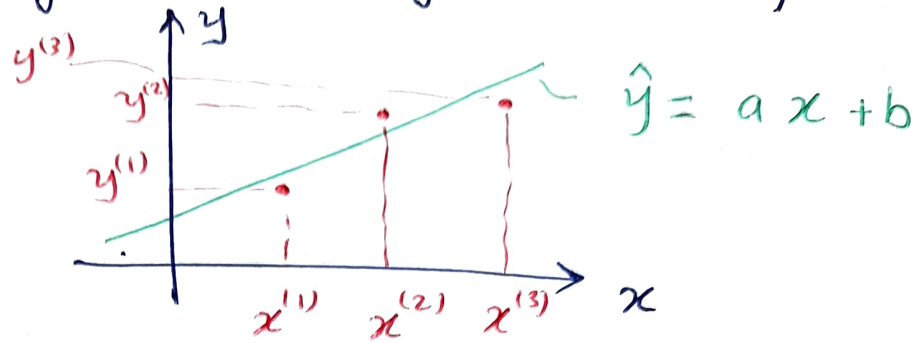
m data points

# linear regression (single variable)

$$y_1 = ax_1 + b$$

$$y_2 = ax_2 + b$$

$$y_3 = ax_3 + b$$



→ two approaches

① using Gradient descent algorithm (using mean squared error)

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

gradient of cost function

Repeat until convergence:

$$a := a - \alpha \cdot \text{gradient}$$

$$b := b - \alpha \cdot \text{gradient}$$

recursive solution  
"optimization problem"

② using Pseudo-inverse (Moore-Penrose inverse)  
(generalized inverse) (The normal equation)

⇒ least squares solution (LS):

$$y_1 = ax_1 + b$$

$$y_2 = ax_2 + b$$

$$y_3 = ax_3 + b$$

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$\underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix}}_X \quad \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{\vec{\theta}} \quad \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}}_{\vec{y}}$

$$X^T X \vec{\theta} = X^T \vec{y}$$

square matrix

$$X^T X \vec{\theta} = X^T \vec{y}$$

$$\mathbb{I} \leftarrow (X^T X)^{-1} (X^T X) \vec{\theta} = (X^T X)^{-1} X^T \vec{y} \Rightarrow \vec{\theta} = \begin{bmatrix} a \\ b \end{bmatrix}$$



hypothesis

Parameters

④

$$\hat{y} = h(\theta_0, \theta_1, \dots, \theta_n) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\vec{\theta} = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n]^T$$

linear regression  
"multivariate"

$$\hat{y} = h_{\theta}(\vec{x}) = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$x_0 = 1$

$\vec{x} \in \mathbb{R}^{n+1}$

$\vec{\theta} \in \mathbb{R}^{n+1}$

$$\vec{x} = [x_0 \ x_1 \ x_2 \ \dots \ x_n]^T$$

n features

n+1 vector

Gradient

del, nabla

divergence  $(\vec{x})$   
 $\vec{\nabla} \cdot \vec{x}$   
curl  $(\vec{x})$   
 $\vec{\nabla} \times \vec{x}$

$$\vec{\nabla} w = \frac{\partial w}{\partial x_1} \hat{x}_1 + \frac{\partial w}{\partial x_2} \hat{x}_2 + \dots = \begin{bmatrix} \frac{\partial w}{\partial x_1} \\ \frac{\partial w}{\partial x_2} \\ \vdots \\ \frac{\partial w}{\partial x_n} \end{bmatrix}$$

Cost / objective / loss function

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

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

MSE

m-data points

sum of squared errors

for all data points

"Vanilla" (Batch)

## Gradient descent algorithm

(5)

all model parameters  
from  $\theta_0$  to  $\theta_n$

( update (simultaneously) all  $\theta$  )

$$\overset{\text{new}}{\theta_0} := \overset{\text{old}}{\theta_0} - \alpha \frac{\partial J(\vec{\theta})}{\partial \theta_0}$$

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

$$\vdots$$
$$\theta_n := \theta_n - \alpha \frac{\partial J(\vec{\theta})}{\partial \theta_n}$$

$\ll$

update

"Repeat until convergence"

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\vec{\theta})$$

gradient of cost function

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

update

in vector form

cost function

vector of model  
parameters

learning rate

this can be  
also expressed  
as:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$

↑  
update