

Review : Vanilla Gradient Descent  
"Batch GD" algorithm

→ update parameters  $\vec{\theta}$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \cdot \text{gradient of } J(\theta_{\text{old}})$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \nabla J(\theta_{\text{old}})$$

$$\frac{1}{2m} \sum_{j=1}^m ( \quad )$$

↑ summation over  
all data points (all  
examples in training  
dataset)

✓ → Variations of GD

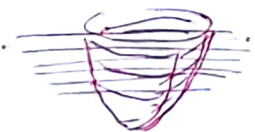
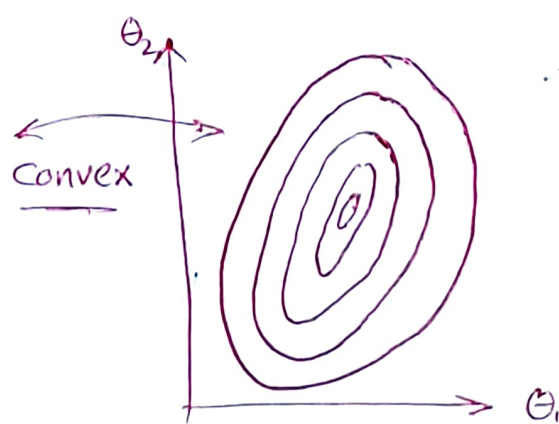
✓ → Problems associated with GD.

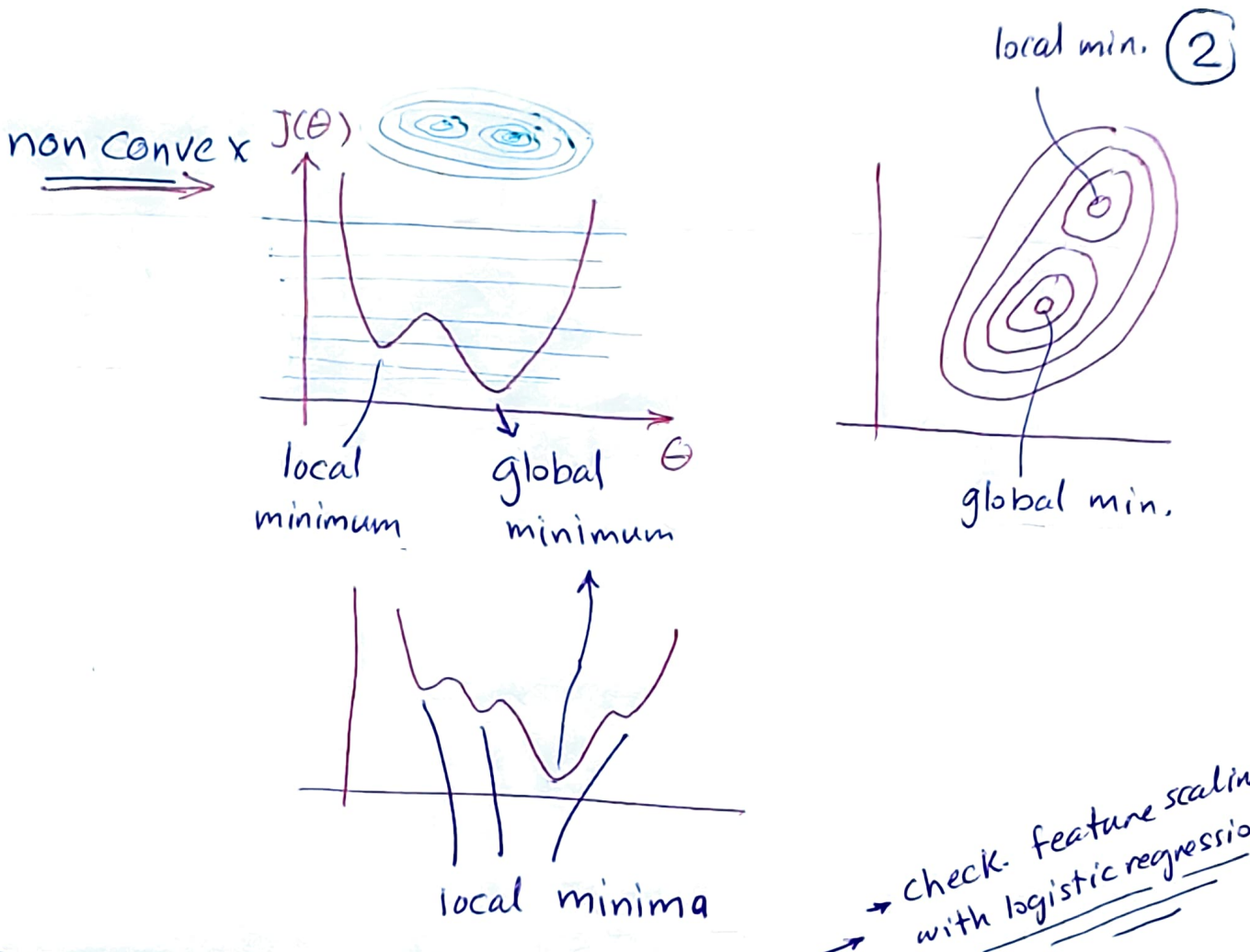
→ acceleration of GD ;

→ { momentum method  
Nesterov method }

• Adaptive GD  
• Adam "ADAM" } next session

→ feature scaling





## ① feature scaling

ex.  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (y - (\theta_0 + \theta_1 x_1 + \theta_2 x_2))^2$

1  
1.5  
2

in this example for simplicity

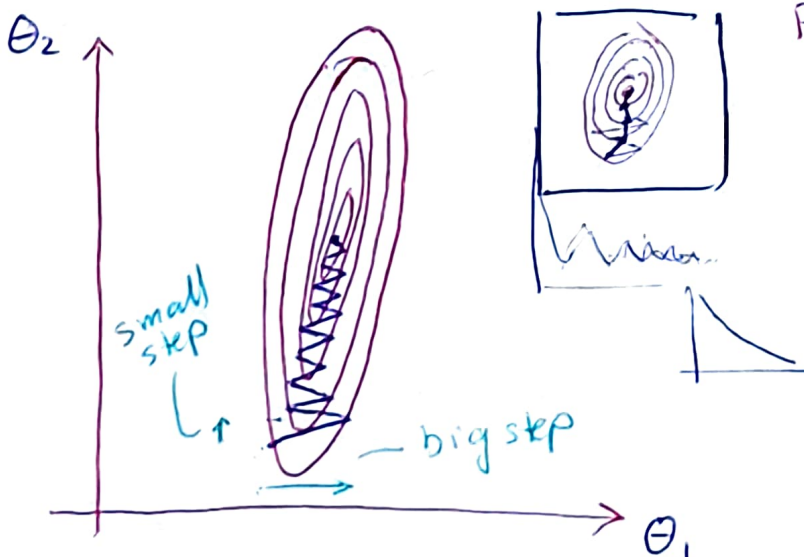
$\theta_0 x_0$

area	# rooms
100	2
200	5

$x_1 \in [100, 200]$

$x_2 \in [2, 5]$

remember; grad.  $\Rightarrow$  derivative  $\Rightarrow$  معدل تغير قيمة (J) إلى التغير في parameter ( $\theta$ )



Parameter ( $\theta$ )

$\hat{y} = 0 + \theta_1 x_1 + \theta_2 x_2$

1  
2  
7

area

$x_1$  has higher scale than  $x_2$

$\Rightarrow \theta_2 \gg \theta_1$

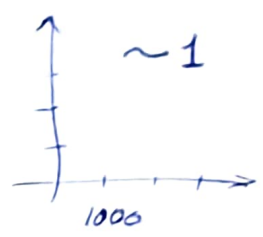
3

$$Y = \cancel{\theta_0} + \theta_1 X_1 + \theta_2 X_2$$

$$\hat{Y} = 0.001 X_1 + 10 X_2$$

$\theta_1$  (down arrow)     $\theta_2$  (up arrow)

$Y$	$X_1$	$X_2$
2	1000	0.1
3	1000	0.2
4	2000	0.2
	3000	



scaling : scaled features

$$X_1 \in [100, 200]$$

$$|X'| \sim 1 \sim \begin{matrix} 0 \rightarrow 1 \\ -1 \rightarrow 1 \end{matrix}$$

$$X_k \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \checkmark \begin{bmatrix} 0.1 \\ 1 \\ 0.8 \end{bmatrix}$$

scaled feature

$$X' \leftarrow X \xrightarrow{1000 \rightarrow 3000} \times \uparrow \uparrow \sim \begin{matrix} \theta_1: 0.0010 \rightarrow 0.0015 \\ \theta_2: 10 \rightarrow 15.01 \end{matrix}$$

$\theta \times$

scaled feature      original feature

i) min-max normalization

$$0 \leq x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \leq 1$$

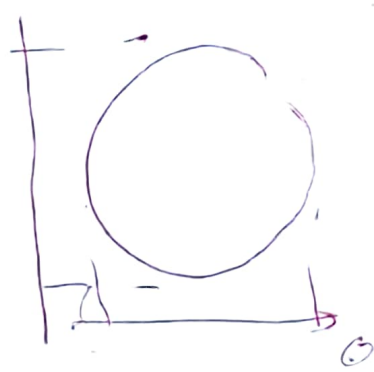
$$x' \in [0, 1]$$

e.g.

$$x_{\min} = 100$$

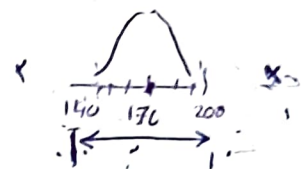
$$x_{\max} = 200$$

$$x_{\max} - x_{\min} = 200 - 100 = 100$$



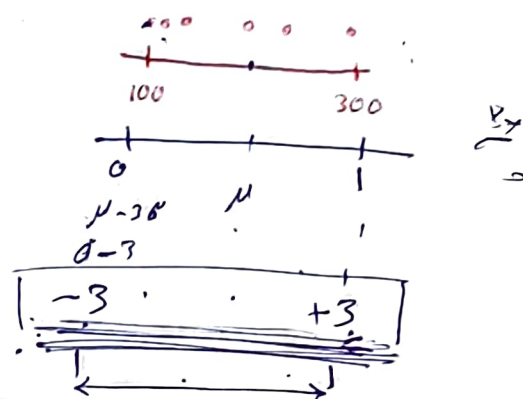
$X_1$	$x'_1$
$x_{\min} \leftarrow 100$	$\rightarrow 0$
100	$\rightarrow 0$
150	$\rightarrow 0.5$
$x_{\max} \leftarrow 200$	$\rightarrow 1$
180	$\rightarrow 0.8$
1000	$\rightarrow 1$

④



ii) mean-normalization (standardization)

$$x' = \frac{x - \bar{x}}{\sigma} \sim \text{mean}$$

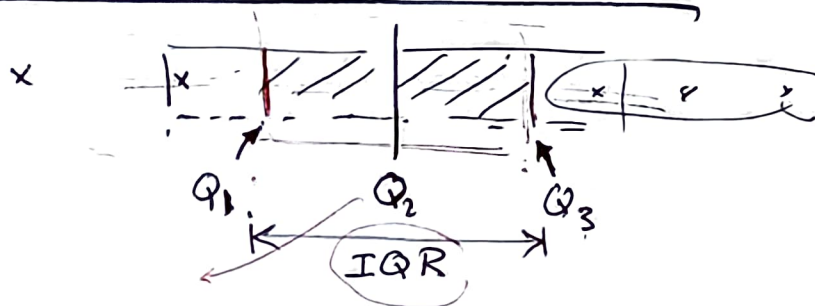


iii) Quantiles

$$IQR = Q_3 - Q_1$$

Robust scaling

$$x' = \frac{x - Q_2}{Q_3 - Q_1}$$



→ measures of central tendency

→ measures of spread (variability)

is feature scaling needed in all ML algorithms?

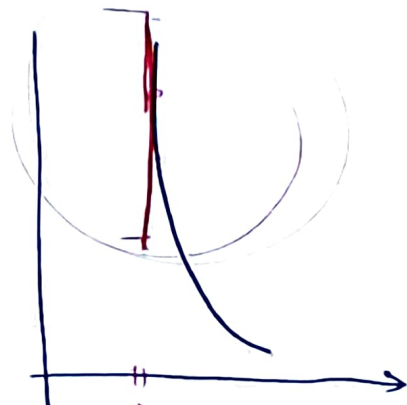
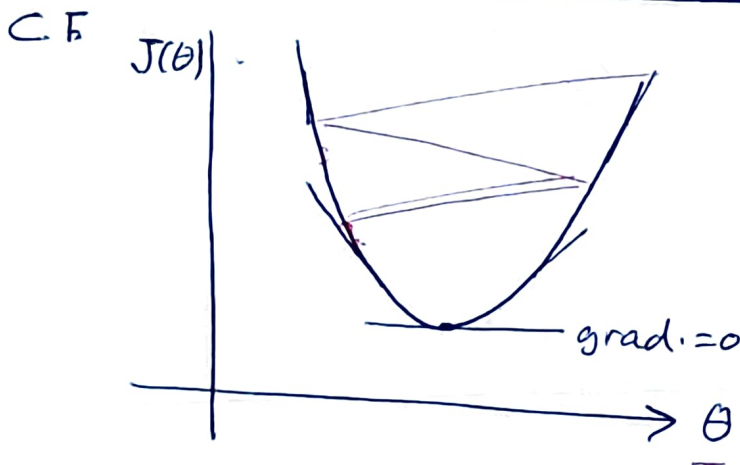
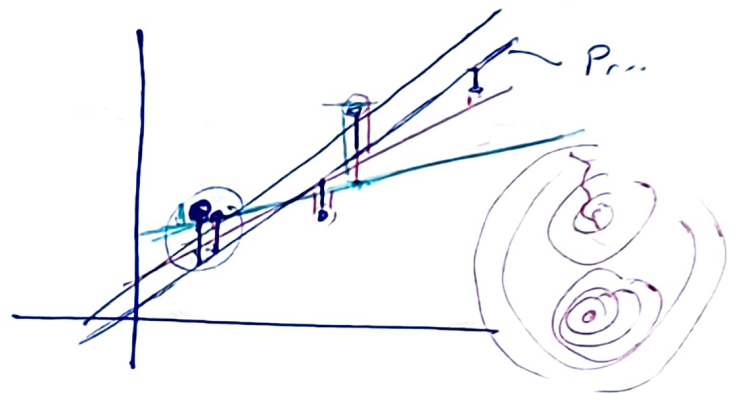
→ not in algorithms that are not distance-based.



## ② Variants of Batch GD

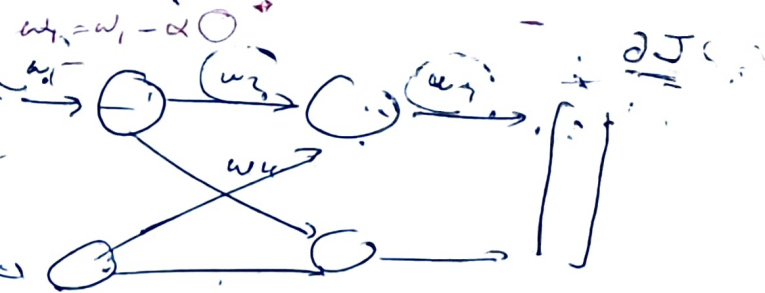
"Vanilla GD" ;

- stochastic gradient descent
- mini-batch gradient descent

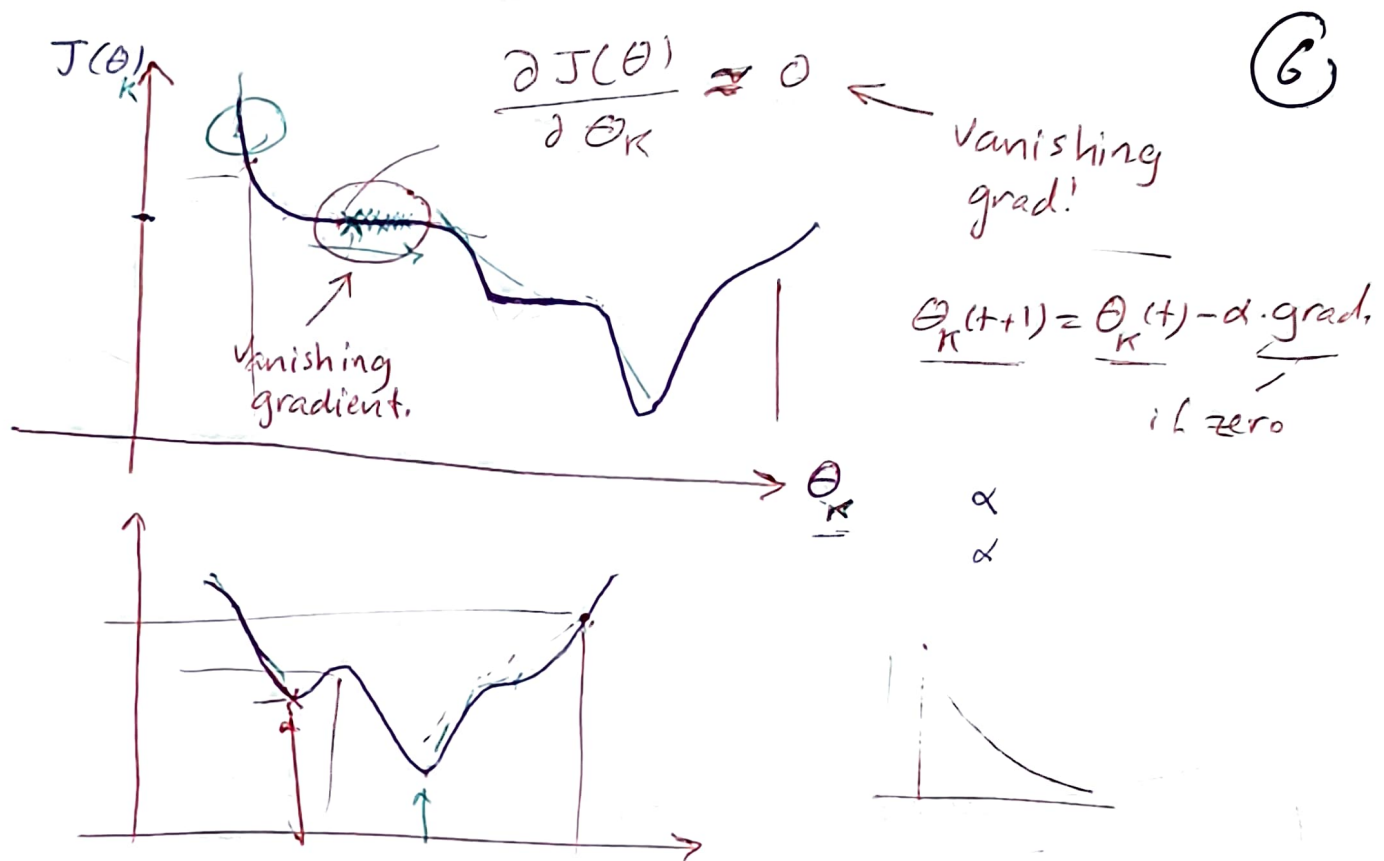


- exploding grad.
- vanishing grad.

$$\left( \frac{\partial C}{\partial \theta} \times \frac{\partial C}{\partial \theta} \times \frac{\partial C}{\partial \theta} \right)$$



← Chain rule



Gradient

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

cost function : scalar function

Jacobian

$\mathcal{F}$  (vector function)

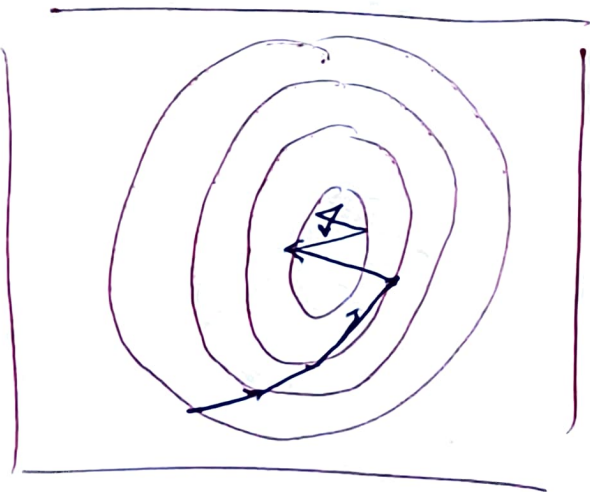
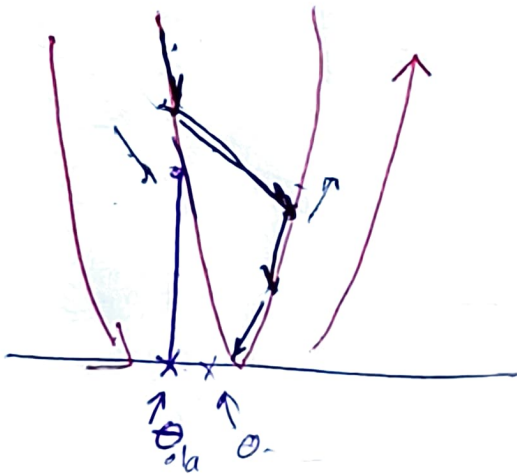
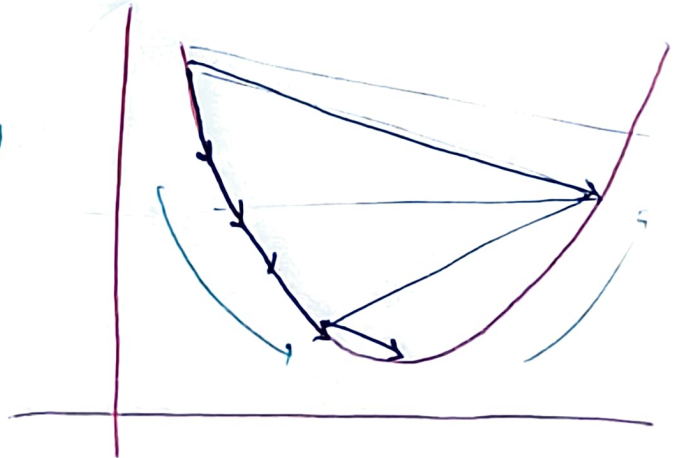
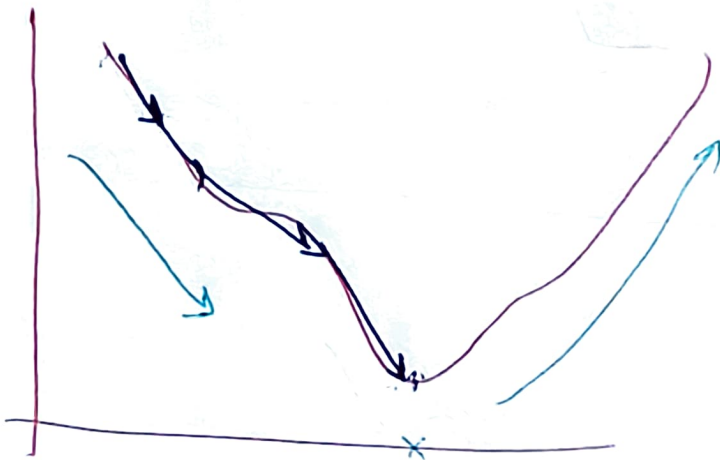
$$\mathcal{F}(\vec{F}(\vec{x}))$$

$$\vec{x}_1, \vec{x}_2 \quad \vec{x}_2 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\vec{F} = \begin{bmatrix} x_1 + x_2 \\ x_1 \times x_2 \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}$$

$$\mathcal{F}(\vec{F}) = \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} & \dots & \frac{\partial F_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} & \dots & \frac{\partial F_2}{\partial x_n} \end{bmatrix}$$

# acceleration methods of GD Algorithm



## Momentum GD

$$\Theta(t+1) = \Theta_{\text{new}} = \Theta_{\text{old}} - \alpha \text{gradient}$$

$\alpha \nabla J(\Theta_{\text{old}})$

$\Theta_{\text{old}} \Rightarrow \text{gradient}(\Theta_{\text{old}}) \Rightarrow \text{update parameters}$

$$\Theta_{\text{new}} = \Theta_{\text{old}} - \alpha \text{gradient}(\Theta_{\text{old}})$$

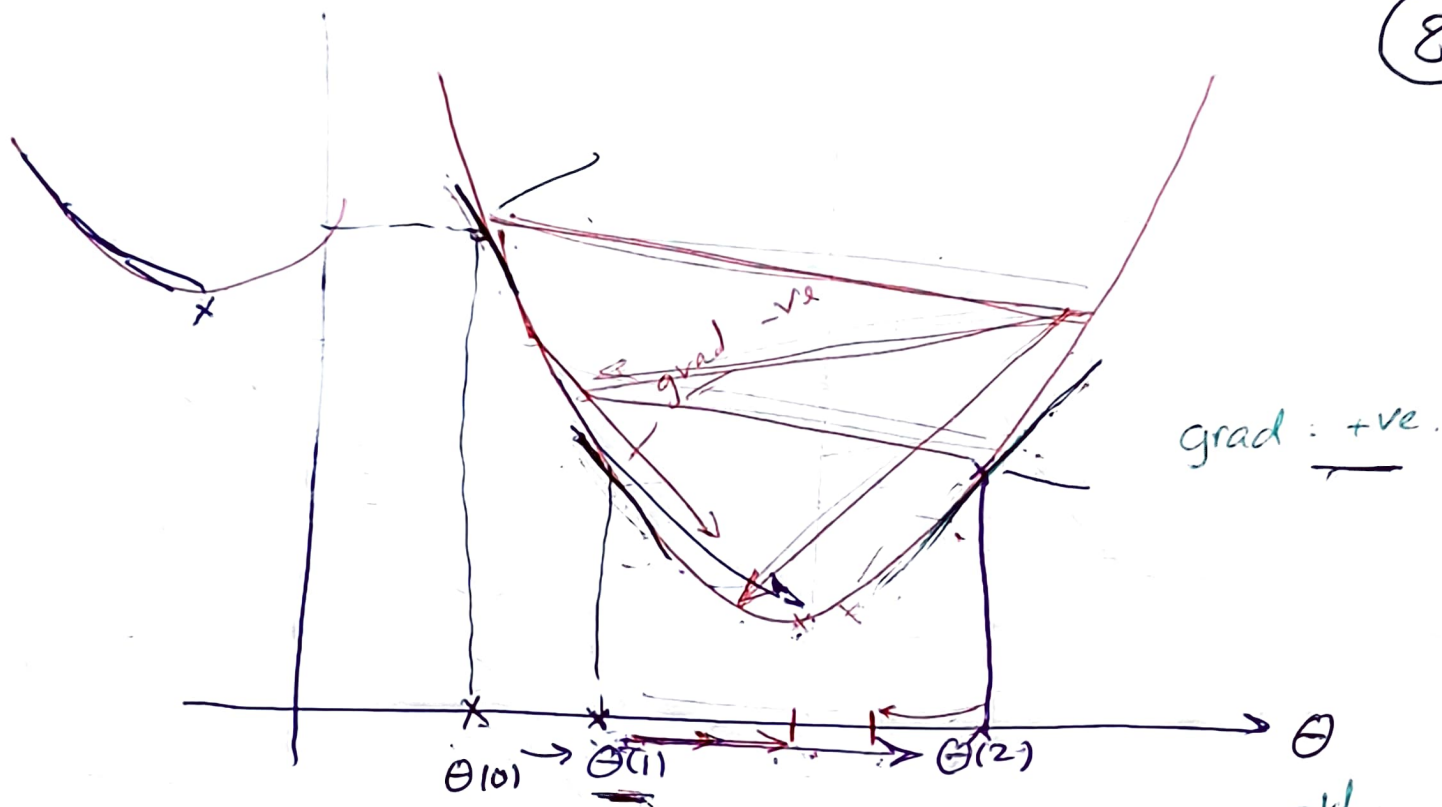
→ momentum

$$\Theta_{\text{new}} = \Theta_{\text{old}} - \text{new gradient}$$

*current gradient*

$$\Theta_{\text{old}} - (\gamma \times \text{old gradient} + \alpha \nabla J(\Theta_{\text{old}}))$$

⑧



$$\theta(1) = \theta(0) - \alpha \text{grad}(\theta_0)$$

momentum

$$\theta(2) = \theta(1) - \alpha \text{grad}(\theta_1)$$

$$\text{new gradient} = \gamma \cdot \text{old grad.} + \alpha \nabla J(\theta_1)$$

$$\theta(3) = \theta(2) - \left( \gamma \cdot \text{old grad.} + \alpha \nabla J(\theta_2) \right)$$

$\gamma < 1$

new grad. evaluated at  $\theta_2$

Nesterov

$$\theta_{\text{temp}} = \theta_{\text{old}} - \gamma \times \text{old grad}$$

$$\theta_{\text{new}} = \theta_{\text{temp}} - \alpha \nabla J(\theta_{\text{temp}})$$

$$\text{new gradient} = \gamma \times \text{old gradient} + \alpha \nabla J(\theta_{\text{temp}})$$