

3/1/2026

PAGE
DATE

①

NOFDS & ML session 3 Mansoura AI 46

① Feature Scaling

0 → 1

"measures of spread"

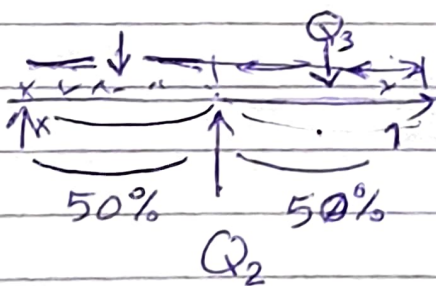
x_1	x_2
1	100
2	120
3	150

→ Range : $x_{\max} - x_{\min}$

→ Standard deviation $\sigma_x = \sqrt{s_x} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$

→ Interquartile range $IQR = Q_3 - Q_1$

Quartiles ;



Q_0 : minimum value

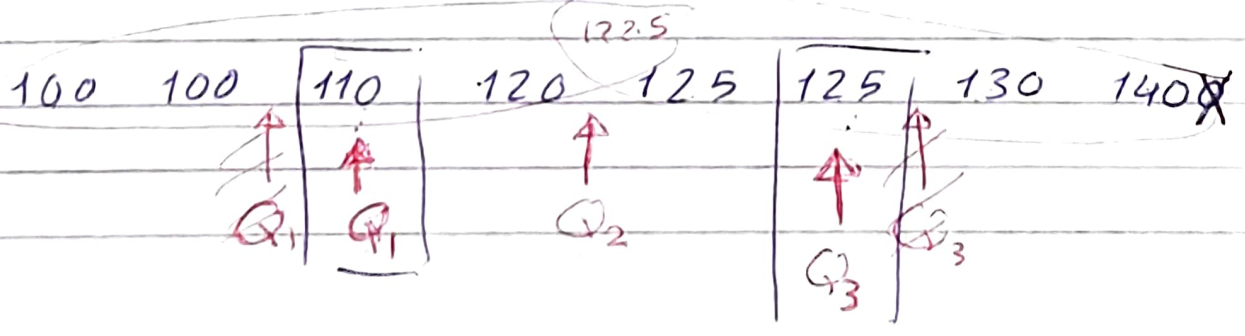
Q_1 : 25% : 75%

Q_2 : median

Q_3 : 75% : 25%

Q_4 : maximum value

$$Q_3 - Q_1 = \boxed{IQR} \rightarrow \text{Robust}$$



feature scaling $x_i \rightarrow x'_i$
 \rightarrow new scale (normalized)

$$\sim 0 \rightarrow 1$$

$$\sim -1 \rightarrow 1$$

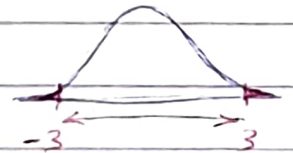
1) Range: min-max normalization

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

4, 2, 3, 3, 4, 5, ..., 5, 50

2) mean-normalization (standardization)

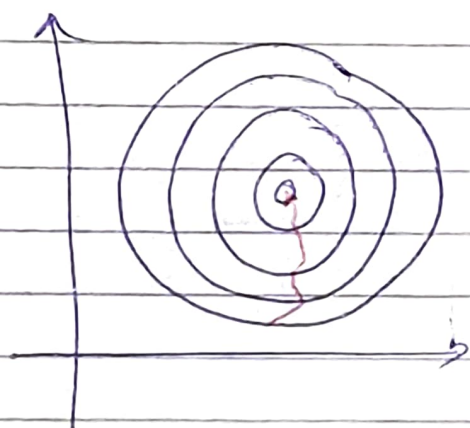
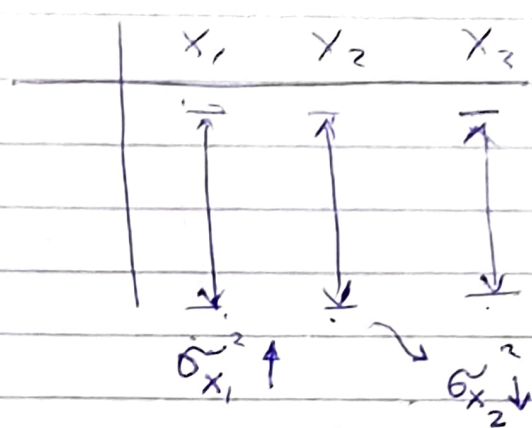
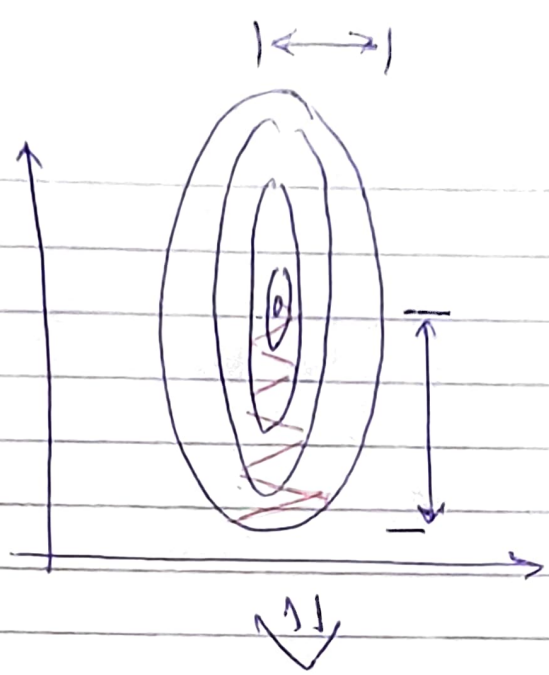
$$x'_i = \frac{x_i - \bar{x}}{\sigma_x}$$



zero-mean, standard dev. = 1

3) Robust feature scaling (Robust scaling)

$$x'_i = \frac{x_i - \text{median}}{IQR} = \frac{x_i - Q_2}{Q_3 - Q_1}$$



Kappa $\left[\kappa \approx \frac{\sigma_{x_1}^2}{\sigma_{x_2}^2} \right]$
 \downarrow
 κ : ratio
 \rightarrow condition number of optimization
 max. f. variance
 min. feature variance

Computations of Convergence $\sim O(\kappa)$

Gradient descent Algorithms

\rightarrow	<u>momentum</u>
\rightarrow	<u>Nesterov</u>
\rightarrow	$O(\sqrt{\kappa})$

Adaptive GD \rightarrow session 4

Review and further discussion

PAGE
DATE

4

→ Gradient (\vec{f}) = $\nabla J = \begin{bmatrix} \partial J / \partial \theta_1 \\ \partial J / \partial \theta_2 \\ \vdots \\ \partial J / \partial \theta_n \end{bmatrix}$

vector field

scalar-valued
function

↘

vector (after substituting for values of θ_i)

→ Jacobian of a vector valued function

$$\vec{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$$

scalar valued functions

$$J(\vec{f}) = \begin{bmatrix} \partial f_1 / \partial \theta_1 & \partial f_1 / \partial \theta_2 & \dots & \partial f_1 / \partial \theta_n \\ \vdots & \vdots & \ddots & \vdots \\ \partial f_m / \partial \theta_1 & \partial f_m / \partial \theta_2 & \dots & \partial f_m / \partial \theta_n \end{bmatrix}$$

$$J(\vec{f}) = \begin{bmatrix} \nabla^T(f_1) \\ \nabla^T(f_2) \\ \vdots \\ \nabla^T(f_m) \end{bmatrix}$$

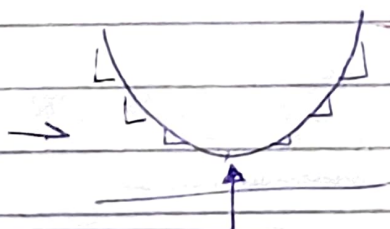
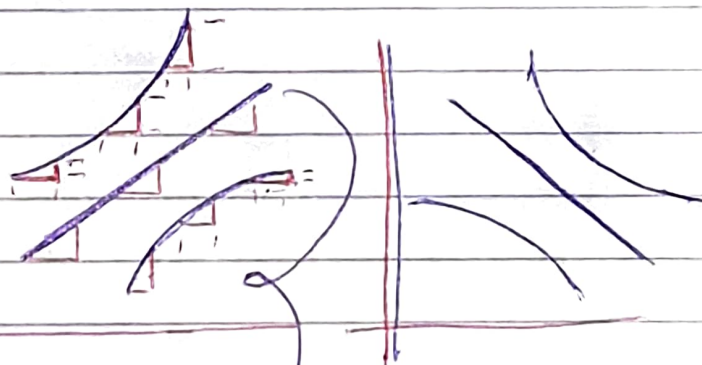
Hessian (J)

scalar-valued function

$$H(J) = \nabla^2(J) = \nabla \nabla(J)$$

$$H(J) = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_1^2} & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} & \dots & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 J}{\partial \theta_n \partial \theta_1} & \frac{\partial^2 J}{\partial \theta_n \partial \theta_2} & \dots & \frac{\partial^2 J}{\partial \theta_n^2} \end{bmatrix}$$

→ Convex
→ Concave
→ flat?
→ inflection?



1st derivative

2nd derivative

zero

$$\kappa = \frac{\lambda_{\max}(H)}{\lambda_{\min}(H)}$$

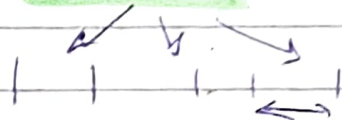
→ 1

$$O(\kappa)$$

GD Variants

Review

Vanilla GD / Batch GD



repeat
until
conver-
gence

$$\Theta^{\text{new}} = \Theta^{\text{old}}$$

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

for MSE:

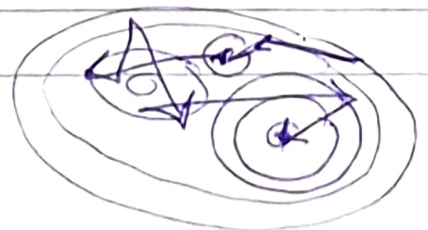
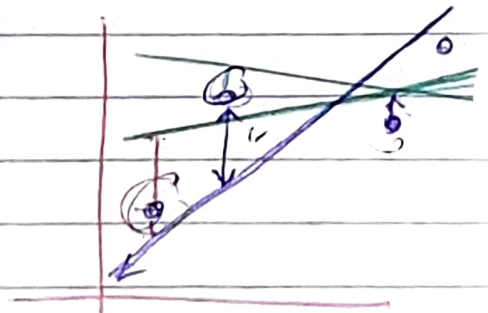
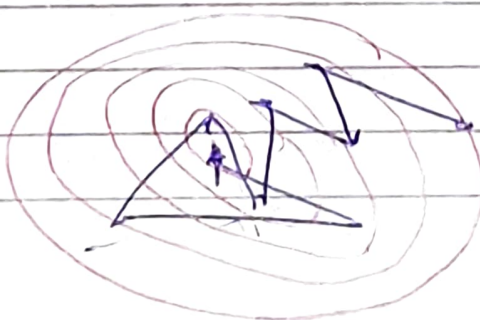
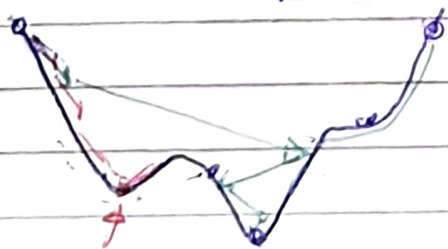
$$\nabla \left(\frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 \right)$$

e.g., in Linear
Regression

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

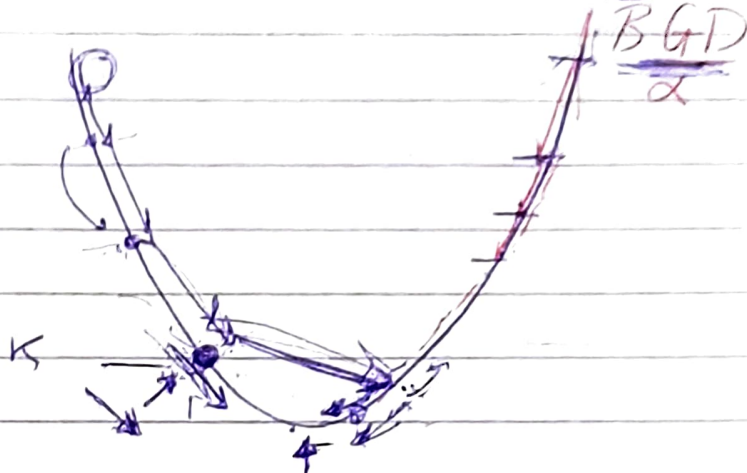
→ Stochastic GD "SGD"

→ mini-batch GD



Momentum

→ "heavy-ball"



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

$$\Theta^{(K+1)} = \Theta^{(K)} + \left[\begin{array}{c} \text{Current} \\ \text{Velocity} \end{array} - \begin{array}{c} \text{Current} \\ \text{gradient} \end{array} \right]$$

$$V^{(K+1)} = \mu V^{(K)} - \alpha \nabla J(\Theta^{(K)})$$

OR β OR γ

$$K=0 \quad V^{(0)} = 0$$

(μ=1)
(1-μ)
(1-β)

$$K=1 \quad V^{(1)} = \mu \times 0 - \alpha \text{ gradient}(\Theta^{(0)})$$

$$K=2 \quad V^{(2)} = \mu(V^{(1)}) - \alpha \text{ gradient}(\Theta^{(1)})$$

$$= \mu(-\alpha \text{ grad}(\Theta^{(0)})) - \alpha \text{ grad}(\Theta^{(1)})$$

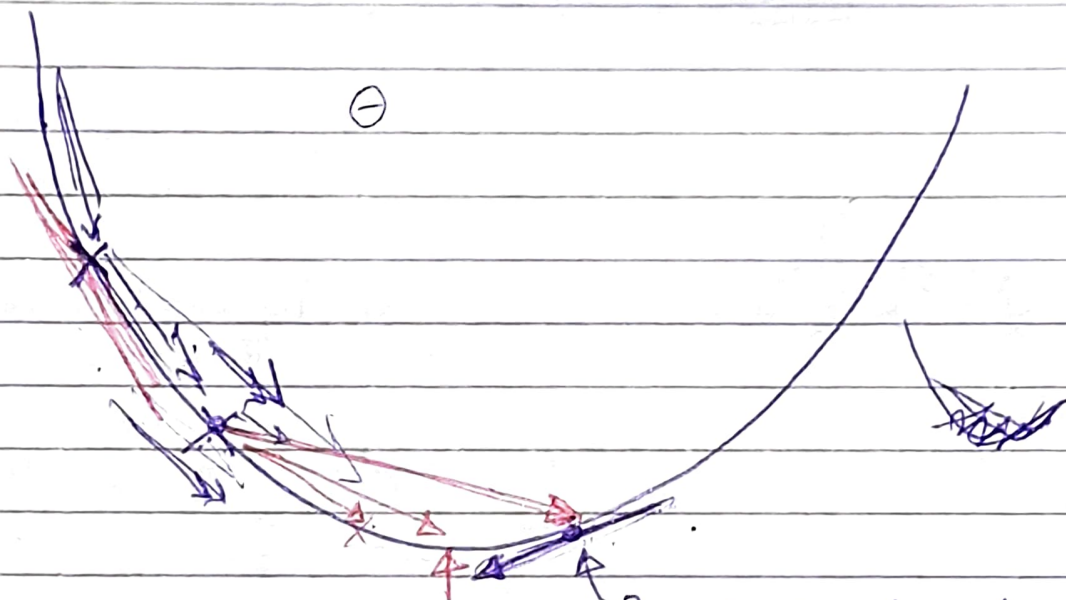
$$K=3 \quad V^{(3)} = \mu(V^{(2)}) - \alpha \text{ grad}(\Theta^{(2)})$$

$$= \mu(\mu(-\alpha \text{ grad}(\Theta^{(0)})) - \alpha \text{ grad}(\Theta^{(1)})) - \alpha \text{ grad}(\Theta^{(2)})$$

↑
EWA
EWMA

Nesterov method

\Rightarrow look ahead!



$\psi(k)$

? look ahead
gradient at Θ
look ahead.

β optimally ~ 0.9

$$\underline{0.8} < \beta < 0.999$$

Nesterov

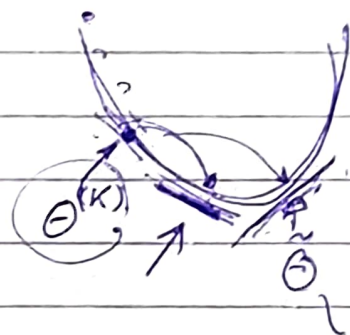
→ look-ahead parameter

→ temp. value of Θ ; $\tilde{\Theta}^{(k)}$

$$\tilde{\Theta}^{(k)} = \Theta^{(k)} + \mu \nabla J(\Theta^{(k)})$$

extrapolated value of Θ

~~$\nabla J(\Theta)$~~



$$\underline{\underline{v^{(k+1)} = v^{(k)} - \alpha \nabla J(\tilde{\Theta}^{(k)})}}$$

$\Theta_{temp.}$

$$\Theta^{(k+1)} = \Theta^{(k)} + v^{(k+1)}$$