

28-09-2025

## Agenda - ML - III

→ Gradient, Gradient Descent, Practical, Linear Regression (EDA - Model prediction)

### Dataset:

$y_{\text{exp}} = \exp(x_1)$	Completed_projects ( $x_2$ )	Salary_y_k ( $y$ )
1 $\times m_1$	0 $\times m_2$	39.761
2	1	48.400
3	1	56.978
4	2	68.240
5	3	77.867
6	4	85.022

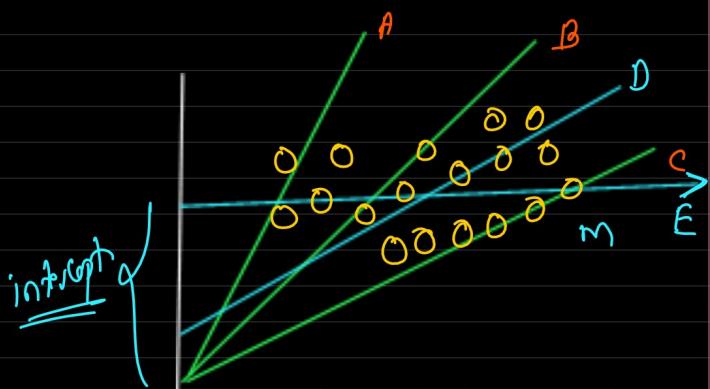
$$y = mx + c$$

$$y = mx + c \rightarrow \text{intercept } b \text{ (q)}$$

$$\begin{aligned} x &= 10 \\ m &= 1.5 \\ c &= 0 \\ y &= 20 \\ y &\rightarrow 15 \end{aligned}$$

$$\hat{y} = mx + c$$

$$\left. \begin{array}{l} \text{error} \rightarrow 5 \\ \text{1 data point} \\ \text{for prediction calculation} \end{array} \right\}$$



$x_1$	$y$	$\hat{y}$	$\hat{y} - y$	$m = 11$
1	11.5	11	11.5	
2	22.5	22	22.5	
3	33.5	33	33.5	
4	44.5	44	44.5	

vectorised method → array → 2D → matrix  
 → 3D → matrix  
 → ND → N-dim matrix

previous      new      multiple features  $[x_1, x_2]$

$$y = m \underset{m}{\underbrace{x}} + c$$

$$y = M \underset{\text{multiple } m}{\underbrace{x}} + c \quad \text{only one intercept}$$

$$y = \underset{M}{\underbrace{[m_1, f_1]}} + \underset{M}{\underbrace{[m_2]}} x_2 + c$$

## Matrix Multiplication:

Condition: Column of Matrix 1 = Row of Matrix 2

Rule of output shape:  $A(R) \times B(C)$

$$A = \begin{bmatrix} 1, 2, 1, 3, \\ 4, 5, 1, 6 \end{bmatrix} \quad \begin{bmatrix} 7, 8 \\ 9, 10 \\ 11, 12 \end{bmatrix} = B$$

output →  $2 \times 2$  matrix

Diagram showing dimensions:  
 Matrix A:  $2 \times 3$  (2 rows, 3 columns)  
 Matrix B:  $3 \times 2$  (3 rows, 2 columns)

$c(3) = A(3)$

$$\begin{bmatrix} 1, 2, 1, 3, \\ 4, 5, 1, 6 \end{bmatrix} \quad \begin{bmatrix} 7, 8 \\ 9, 10 \\ 11, 12 \end{bmatrix}$$

Calculated values:  
 $1 \times 7 + 2 \times 9 + 3 \times 11 \rightarrow 58$   
 $1 \times 8 + 2 \times 10 + 3 \times 12 \rightarrow 64$   
 $4 \times 7 + 5 \times 9 + 6 \times 11 \rightarrow 139$   
 $4 \times 8 + 5 \times 10 + 6 \times 12 \rightarrow 154$

$$\begin{bmatrix} x_0 & x_1 & x_2 \\ 1 & 1 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 1 \\ 1 & 4 & 2 \\ 1 & 5 & 3 \\ 1 & 6 & 4 \end{bmatrix}_{6 \times 3}$$

$$\begin{bmatrix} m_0 & m_1 & m_2 \end{bmatrix}_{1 \times 3}$$

$\downarrow$  Transpose  $\rightarrow 3 \times 1$

$$\begin{bmatrix} m_0 \\ m_1 \\ m_2 \end{bmatrix}_{3 \times 1}$$

$\rightarrow$  output shape:  $6 \times 1$

Python:

$$\begin{aligned}
 X @ B \\
 [6 \times 1] &\leftarrow np.dot(X, B) \\
 &\quad \xrightarrow{\text{interval}} \\
 &\quad (x_0 m_0) + x_1 m_1 \\
 &\quad + x_2 m_2
 \end{aligned}$$

Gradient Descent: —

Row of exp(x)	Completed projected(x)	Solar-Jk(y)
1	0	39.761
2	1	48.400
3	1	56.748
4	2	68.290
5	3	77.864
6	4	85.022

$$\begin{matrix} 0 \\ m_1 x_1 + m_2 x_2 + m_0 \end{matrix} \rightarrow \hat{y} \quad y$$

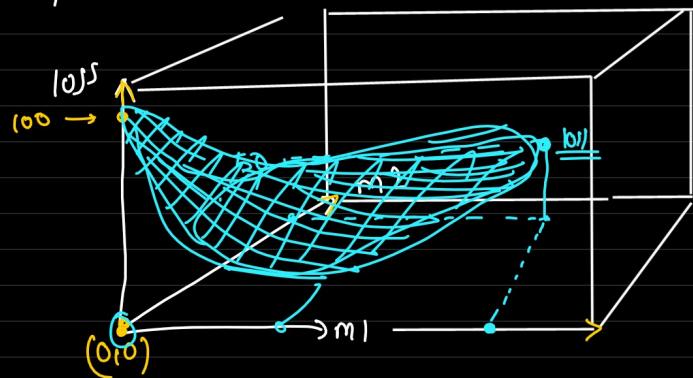
$$\sum_{i=1}^n \text{error} \rightarrow (y - \hat{y}) \rightarrow \text{Total loss} \\
 \uparrow \\
 \text{cost function}$$

$$m_1 = 0, m_2 = 0, m_0 = 0$$

$$\begin{array}{ccc}
 \hat{y} & ^\wedge & \text{error} \\
 \begin{array}{ccc}
 y & & \\
 39.761 & 0 & 39.761 \\
 48.400 & 0 & 48.400 \\
 56.748 & 0 & 56.748 \\
 68.290 & 0 & 68.290 \\
 77.864 & 0 & 77.864 \\
 85.022 & 0 & 85.022
 \end{array} & \left. \begin{array}{c} \uparrow \\ \text{add} \rightarrow \frac{1}{n} \sum (y - \hat{y})^2 \\ \rightarrow 200 \left( m_1 = 0, m_2 = 0, m_0 = 0 \right) \end{array} \right. & \rightarrow \text{new loss} \left( m_1 = 0.01, m_2 = 0.1, m_0 = 0.1 \right)
 \end{array}$$

for each value of parameter we will get different loss  
if we plot loss w.r.t  $m_1, m_2, m_0$

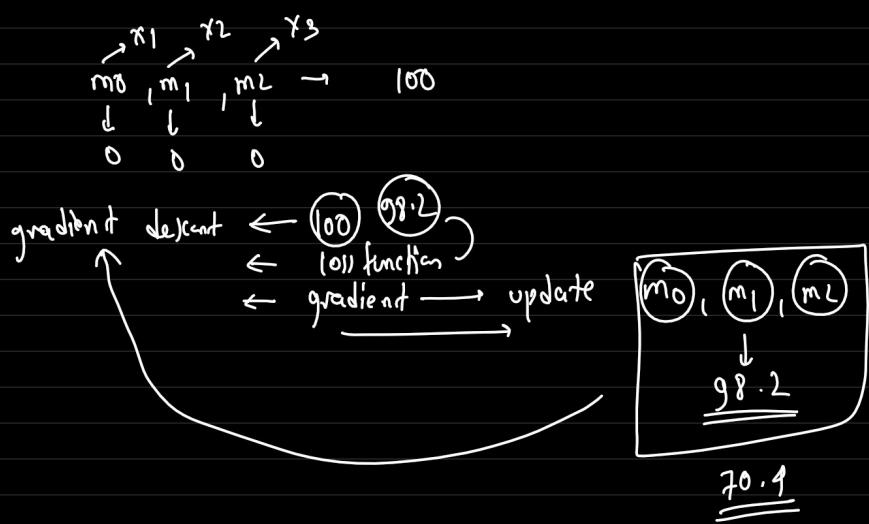
Set 1  $(m_0, m_1, m_2) \rightarrow 100$   
 Set 2  $(m_0, m_1, m_2) \rightarrow 20 \rightarrow$   
 Set 3  $(m_0, m_1, m_2) \rightarrow 400$   
 Set 4  $(m_0, m_1, m_2) \rightarrow 88$



Ideally we want to have lower loss  $\rightarrow m_0, m_1, m_2$



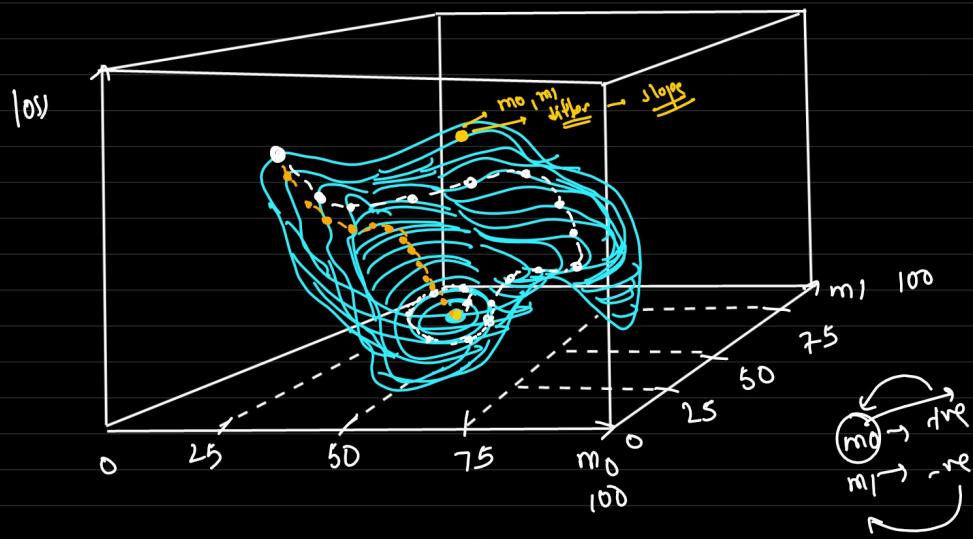
Gradient Descent is an optimization algo used in ML. Its a iterative method to adjust parameter of a model to find minimum loss of a function.



$$\begin{array}{ll}
 1 & 11.5 \\
 2 & 22.5
 \end{array}
 \quad
 \begin{array}{l}
 m=1 \xrightarrow{\text{up}} \text{true direction} \\
 m=2 \xrightarrow{\text{down}} \text{how much!} \\
 \vdots \\
 m=11 \rightarrow \approx 0 \\
 \hline
 m=100 \rightarrow 0.5
 \end{array}$$

$$\begin{array}{ccc}
 x_1 & \left[ \begin{matrix} m \\ 1 \end{matrix} \right] & y \\
 1 \times & 11.5 & 1 \\
 2 \times & 22.5 & 2
 \end{array}
 \quad
 \begin{array}{c}
 10.5 + 20.5 \rightarrow 33 \\
 | \\
 10.5 \rightarrow 28
 \end{array}$$

reduce loss, for lowest loss get parameters.

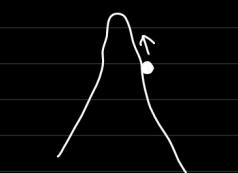


$$\text{loss} \approx 0 \rightarrow m_0 \rightarrow 25 \\ m_1 \rightarrow 75$$

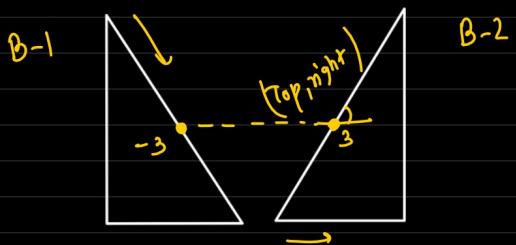
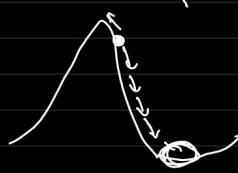
Gradient Descent: Goal is to find minimum loss

→ we use differentiation to calculate slope of the cost function at our current position.

- slope / gradient will tell us direction of the steepest ascent. (uphill)
- Take a step in the opposite direction. steepest descent. (downhill)



steepest gradient (uphill)  
direction of valley  $\approx$  opposite of  
steepest gradient



B-2 has higher slope

left to right upwards  $\rightarrow$  +ve slope

left to right downwards  $\rightarrow$  -ve slope

### Differentiation:

$$BM_1 = \frac{\omega}{h^2} \quad \text{parameter} \rightarrow \omega, h$$

height = 1.8

$\rightarrow$  change weight, keeping the height of constant.

$$\frac{\partial BM_1}{\partial \omega} = \frac{1}{h^2} = \frac{1}{(1.8)^2} \approx 0.309$$

$$\begin{aligned} \omega &= 100 \\ \omega &= 101 \end{aligned} \quad ] + 1 \text{ kg}$$

$$BM_1 = 1 \times 0.309 \rightarrow \underline{\underline{+0.309}}$$

$$\begin{aligned} \omega &= 100 \\ \omega &= 98 \end{aligned} \quad ] - 2 \text{ kg} \quad \approx -2 \times 0.309 \rightarrow -\underline{\underline{0.618}}$$

If I change my parameter by a small value, Derivative tells you how much the output changes per unit

Tells us: What will be the rate of change if input changes

start point  $\boxed{\quad}$  end point

Differentiation of the gradient:

$$J(\beta) = \frac{1}{2m} \sum (y - \hat{y})^2 \rightarrow \underline{\text{start point}}$$

$J \rightarrow \text{cost function}$   
 $\beta \rightarrow \text{parameters } (m_0, m_1, m_2)$

$$= \frac{1}{2m} \sum (y - X\beta)^2$$

$$= \frac{1}{2m} \sum (y - X\beta)^T (y - X\beta)$$

$y^2 \rightarrow \alpha \times y$   
 $A = [1, 2, 3]$

$A^2 \rightarrow A \times A^T$

←      →

keep this  
constant  
aside

$$S(\beta) = (y - X\beta)^T (y - X\beta)$$

$$\text{Rule 1: } (A - B)^T = A^T - B^T$$

$$\text{Rule 2: } (AB)^T = B^T \cdot A^T$$

$$S(\beta) = (y^T - (X\beta)^T) (y - X\beta) \rightarrow \text{Rule 1}$$

$$S(\beta) = (y^T - \beta^T X^T) (y - X\beta) \rightarrow \text{Rule 2}$$

$$(a \oplus b)(c \ominus d) \rightarrow \underline{ac - ad - bc + bd}$$

$$= y^T y - y^T X \beta - \underbrace{\beta^T X^T y}_{\downarrow y^T X \beta} + \underbrace{\beta^T X^T X \beta}_{}$$

simplif: since  $y^T x \beta$  is a scalar value, it is equal to its transpose

$$= \beta^T x^T y$$

$$S(\beta) = \underline{y^T y} - 2\underline{y^T x \beta} + \underline{\beta^T x^T x \beta}$$

Take the gradient:  $\frac{\partial S}{\partial \beta}$ , take partial derivative.

$$\frac{\partial}{\partial \beta} (y^T y) = 0$$

$$\frac{\partial}{\partial x} (\cdot) = 0$$

$$\frac{\partial}{\partial \beta} (-2y^T x \beta) = -2(x^T y)^T \beta = -2x^T y$$

$$\frac{\partial}{\partial x} (q^T x) = q$$

$$\frac{\partial}{\partial \beta} \left( \beta^T \underbrace{x^T x \beta}_{x^T \beta \downarrow x} \right) = \left( x^T x + (x^T x)^T \right) \beta$$

$$\frac{\partial}{\partial x} (x^T \beta x) = (x + x^T) \beta$$

$$\textcircled{x^T \cdot x}$$

$$\begin{aligned} &= (x^T x + x^T x) \beta \\ &= 2x^T x \beta \end{aligned}$$

combine:

$$\frac{\partial S}{\partial \beta} = 0 - 2x^T y + 2x^T x \beta$$

$$\begin{aligned}
 &= 2x^T x\beta - 2x^T y \\
 &= 2x^T (x\beta - y) \\
 &= 2x^T (\hat{y} - y)
 \end{aligned}$$

bringin) the constant :

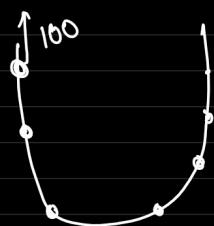
$$\begin{aligned}
 &= \frac{1}{m} \cdot 2x^T (\hat{y} - y) \\
 &= \cancel{x^T} (\hat{y} - y) / m
 \end{aligned}$$

final  
derivative  
of our cost  
function

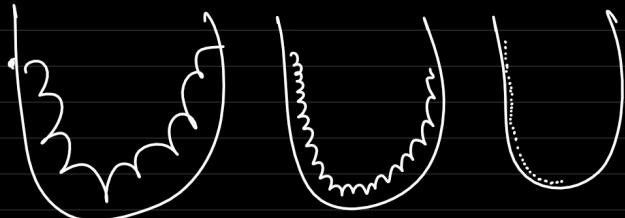
↓  
gradient

$$\text{gradient} = (X.T @ (y_{\text{pred}} - y)) / m$$

= direction & rate of change  $\times \underline{\underline{0.1}}$   
100 learning rate



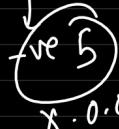
→ 100% depth → start, end, 0.1  
 → direction



$\beta \rightarrow$  matrix  $\rightarrow$  learnable parameter



+5



$$\beta_{\text{new}} = \beta_{\text{old}} - (\text{lr} \times \text{gradient})$$

| lr  $\rightarrow$  learning rate

x 0.01