

Gradient Descent

Gradient descent is an Optimization algorithm used to reduce the cost (error) by repeatedly updating the model weights in the direction of steepest decrease of the cost function.

OLS VS Gradient Descent:-

Point	OLS	Gradient Descent
i) Computation Cost	very high \rightarrow matrix inversion $O(n^3)$	Much lower — scalable for large datasets
ii) Dataset Suitability	Works well for small/medium datasets	Best for large datasets/Big data.
iii) Feature Count	Struggles when features \gg samples (high-dimens)	works well in high dimens
iv) memory usage	Needs entire dataset in memory	Can use stochastic GD \rightarrow mini batches

GD

\rightarrow Generalize to Complex models:-

• OLS only works for linear regression, GD works for almost every ML algorithm.

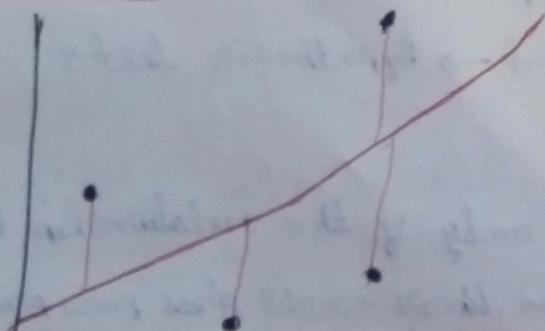
\rightarrow Avoids Expensive matrix inversion:-

No need to compute $(X^T X)^{-1}$

\rightarrow works even when $X^T X$ is non invertible
multicollinearity doesn't block GD.

\rightarrow Online learning possible :- can update made continuously (real-time learning).

Intuition:-



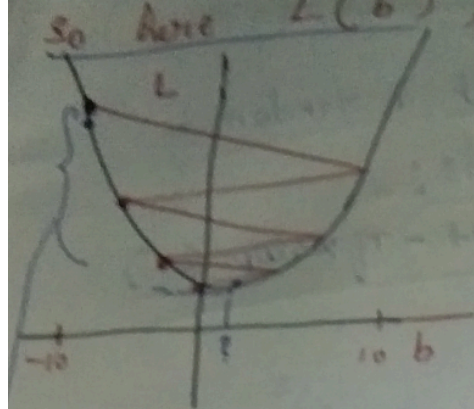
$$\frac{C_{ypa}}{C_{ypa}} \mid \frac{d_{pa}}{d_{pa}}$$

$$\hat{y}_i = mx_i + b$$

$$L = \sum_{i=1}^4 (y_i - \hat{y}_i)^2$$

$$L = \sum_{i=1}^4 (y_i - mx_i - b)^2 \quad \left(\begin{array}{l} \text{could write} \\ \text{consider} \\ \text{if } m = 78 \end{array} \right)$$

$$L = \sum_{i=1}^4 (y_i - 78.35 * x_i - b)^2$$



Step 1:- Select a random b_{min}
Such that L_{min} .

→ Slope = -ve \Rightarrow move forward

→ Slope = +ve \Rightarrow move backward

$$b_{new} = b_{old} - \text{Slope}$$

if $b = -10$, & consider Slope = 50

$$b_{new} = -10 - (-50) \Rightarrow 40 //$$

if $b = 10$, Slope = +50

$$b_{new} = 10 - (50) = -40 //$$

→ There is a drastic changes here, Some time we may miss actual minima point. To reduce, we have to add (η)

$$b_{new} = b_{old} - \eta \text{ slope}$$

→ η = learning rate
default { 0.01 }

$$i) b_{new} = -10 + (0.01 \times 50) = -9.5 //$$

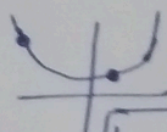
$$ii) b_{new} = -9.5 + (0.01 \times 50) = -9 //$$

here there is slight change

→ The idea is to take repeated steps in the opposite direction of gradient of the function at the current point. because this is the direction of steepest descent.

→ Stepping in the direction of the gradient will lead to local

when to stop:-



Should

$$i) \text{Diff } b_{old} \& b_{new} \Rightarrow b_{new} - b_{old} > 0.001$$

$$b_{new} - b_{old} = 0 \Rightarrow \text{should stop } \} \text{Convergence}$$

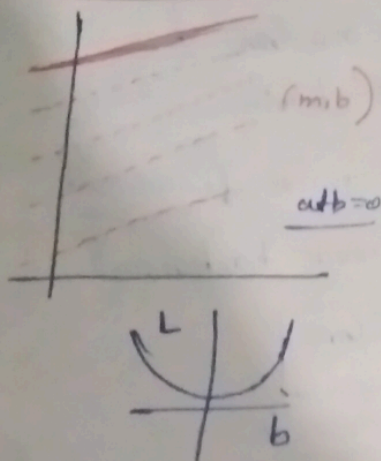
ii) $\text{max Iteration / epochs} \Rightarrow 1000, 500, 100, \text{etc}$

→ Convergence : Gradient becomes very small

If the changes in parameters is almost zero

→ (Divergence) : if learning rate is too high \rightarrow model may overshoot.

Mathematical Formulation: (consider we know $m = 78.35$)



→ start with a random b

for i in epochs: 1000
 $b_{\text{new}} = b_{\text{old}} - \eta \times \text{slope}(b=0)$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\frac{dL}{db} = \frac{d}{db} \sum_{i=1}^n (y_i - mx_i - b)^2$$

$$= 2 \sum_{i=1}^n (y_i - mx_i - b)(-1)$$

$$= -2 \sum_{i=1}^n (y_i - 78.35 x_i - 0)$$

at slope = 0,
 $b_{\text{new}} = b_{\text{old}} - \eta \times \text{slope}_{b=0}$ → stepsize
 then increment $i = 1$, calculate b_{new} again, the will loop
 run until the last epochs.

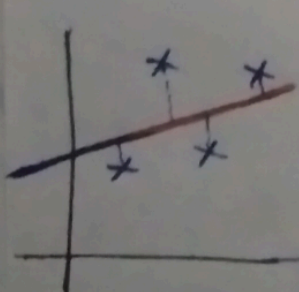
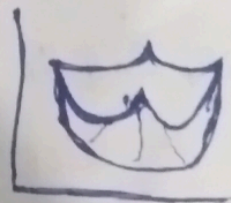
→ In previous mathematical induction, we considered max constant

Now Performing Gradient Descent by adding 'm':

Steps :- 1) initializing values for m and b

$$m = 1 \quad \text{and} \quad b = 0$$

2) epochs = 100, $\eta = 0.01$
 for i in epochs:
 $b = b - \eta \times \text{slope}_b$
 $m = m - \eta \times \text{slope}_m$



$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - mx_i - b)^2$$

→ here $L(m, b)$, L is a function (m & b)

$$b\text{-slope} = \frac{\partial L}{\partial b}$$

$$m\text{-slope} = \frac{\partial L}{\partial m}$$

$$\frac{\partial L}{\partial b} = -2 \sum_{i=1}^n (y_i - mx_i - b)(-1)$$

$$= -2 \sum_{i=1}^n (y_i - mx_i - b)$$

⇒ slope- b at $b = 0$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (y_i - mx_i - b)x_i$$

$$= -2 \sum_{i=1}^n (y_i - mx_i - b)x_i$$

⇒ slope- m at $m = 1$