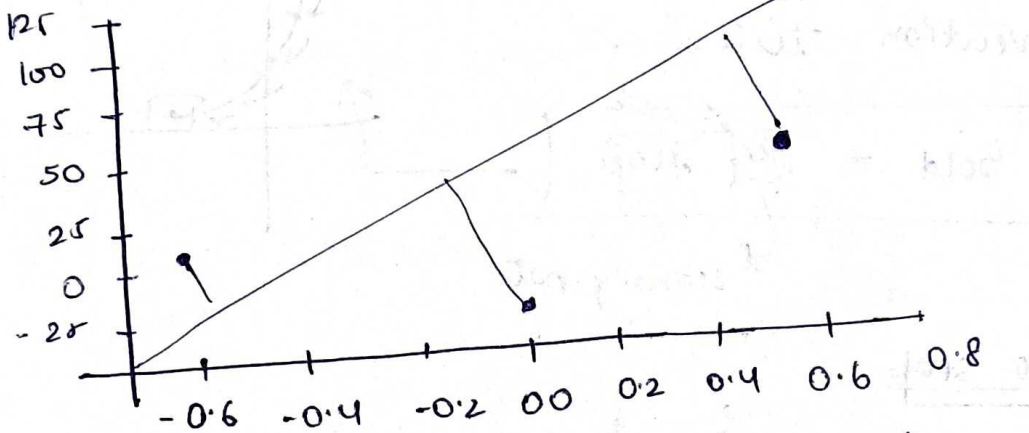


# Gradient Descent

## Intuition

Gradient Descent



4 rows

Cgpa	Lpa
—	—
—	—
—	—
—	—

4 rows, 2 columns

OLS

$$\hat{y}_i = m x_i + b$$

sk learn

klj coll-

$$m = 78.35$$

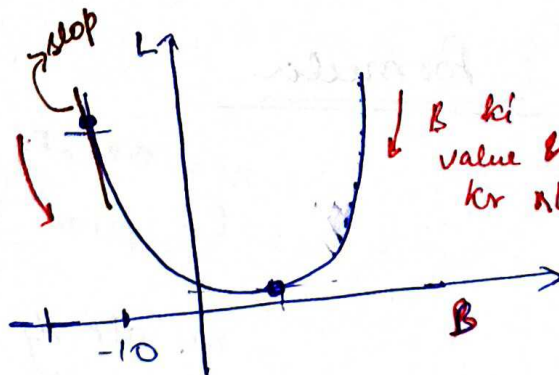
$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

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

$$L \rightarrow b^2$$

mai add  
value kr  
rhe

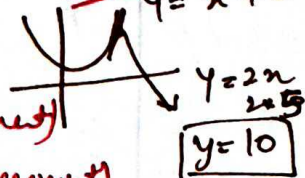


① Step:- select  
a random b

let random

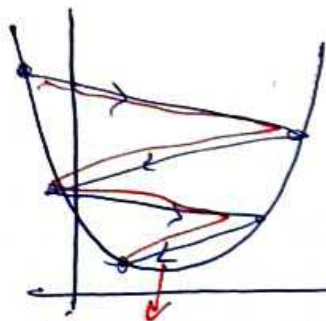
$$b = -10$$

derivative  $x=5$   
 $y = x^2 + 2$



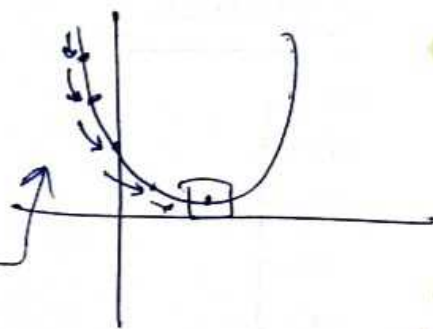
slope = -ve  $\rightarrow$  (B  $\rightarrow$  value increased)  
slope = +ve  $\leftarrow$  (B  $\rightarrow$  value decreased)

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



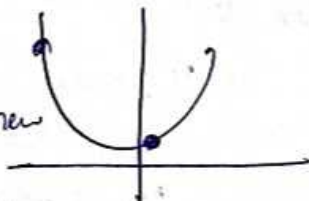
for prevention this

$$b_{\text{new}} = b_{\text{old}} - \underbrace{(\eta)}_{\text{learning rate}} \underbrace{\text{slope}}$$



### When to stop

- ① difference bet<sup>n</sup>  $b_{\text{old}}$   $b_{\text{new}}$   
 $= 0.001$   
 when multiple time difference  
 came same result then  
 need to stop.

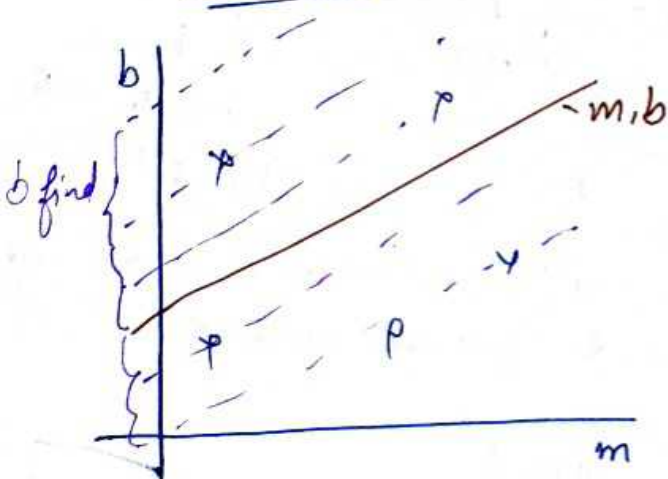


$$b_{\text{new}} - b_{\text{old}} = 0.001$$

$b_{\text{new}} - b_{\text{old}} \geq 0$   
 is close to zero.

- ② Iteration  $\Rightarrow$  limit eg 1000 times, 100 times  
 $\uparrow$   
 epoch  $\rightarrow$  set 1000 or 10000

### Mathematical Formula



$$m = 78.3\%$$

$\hookrightarrow$  from sklearn

m already know we just  
 find b

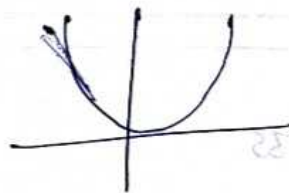
Step ① → start with a random  $b$

for  $i$  in epochs:

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

learning rate let  $\eta = 0.01$   
 $b_{\text{new}} = b_{\text{old}} - \eta \times \text{slope}$   
finding slope at  $b = b$

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



starting  $b = 0$

$$\frac{dL}{db} = \frac{d}{db} \left( \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

$$\frac{dL}{db} = \sum_{i=1}^n (y_i - m x_i - b) = 2 \sum_{i=1}^n (y_i - m x_i - b) (0 - 0 - 1)$$

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

let  $b = 0$

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

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

slope = after solve

at  $b = 0$

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

Step size



from sklearn.datasets import make\_regression  
import numpy as np.

X, y = make\_regression (n-samples=4, n-features=1,  
n-informative=1, n-targets=1, noise=80,  
random-state=13)

b = -100

m = 78.35

lr = 0.1

epochs = 15

for i in range (epochs):

loss-slops = -2 \* np.sum (y - m \* X.ravel() - b)

b = b - (lr \* loss-slops)

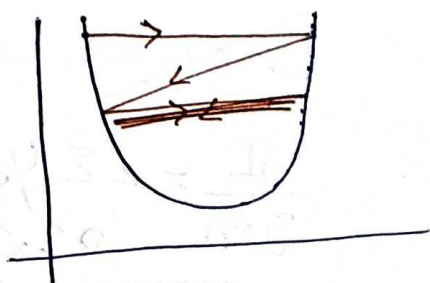
y-pred = m \* X + b

plt.plot (X, y-pred)

plt.scatter (X, y)

convert  
[2D or 3D to 1D]

## Visualization



if condition is hit then,  
learning rate  $\downarrow$   
epochs  $\uparrow\uparrow$

## # Adding m into the mix

### Steps

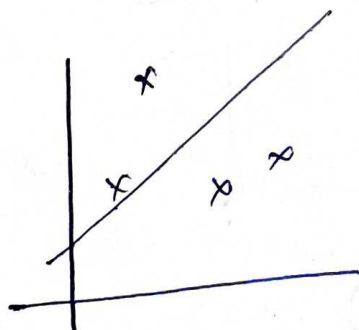
① init random values of  $m$  and  $b$   
 $m=1$  and  $b=0$

② epochs = 100,  $lr = 0.01$

for  $i$  in epochs:

$$b = b - n \text{ slope}$$

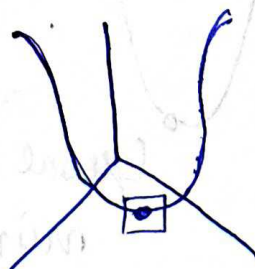
$$m = m - n \text{ slope}$$



$$b\text{-slope} = \frac{dL}{db}$$

$$L = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - mx_i - b)^2$$

$L(m, b)$  like  $f(x, y)$



$$m\text{-slope} = \frac{dL}{dm}$$

$$\sum (y_i - mx_i - b)^2$$

for calculating  $b$

$$\frac{\partial L}{\partial b} = -2 \sum (y_i - mx_i - b)$$

$$= -2 \sum (y_i - mx_i - b)$$

slope  $b$  at  $b=0$

for calculating  $m$

$$\frac{\partial L}{\partial m} = 2 \sum (y_i - mx_i - b) x_i$$

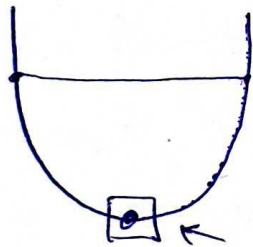
$$= -2 \sum (y_i - mx_i - b) x_i$$

Slope at  $m=1$

## # Effect of loss function

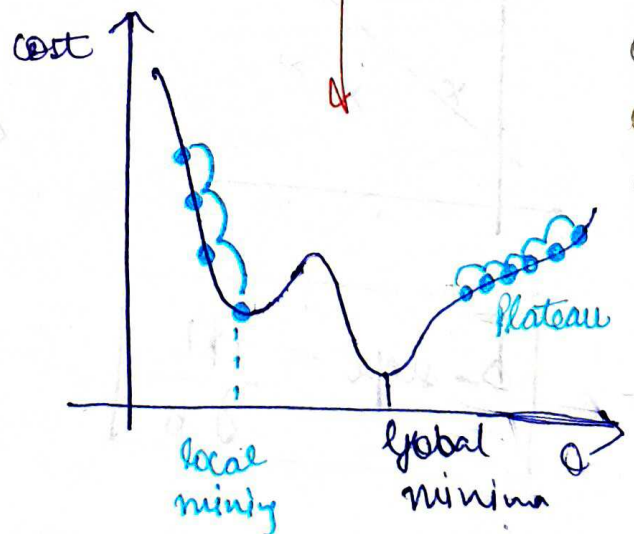
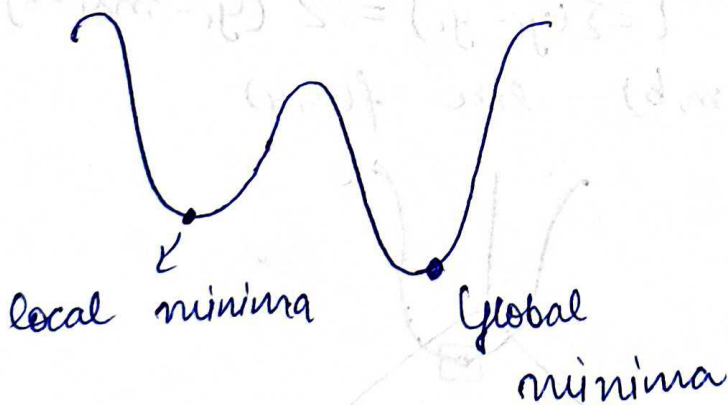
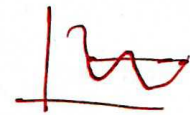
$$L = \sum (y_i - \hat{y}_i)^2$$

Only 1 minima  
Convex function



Global minima

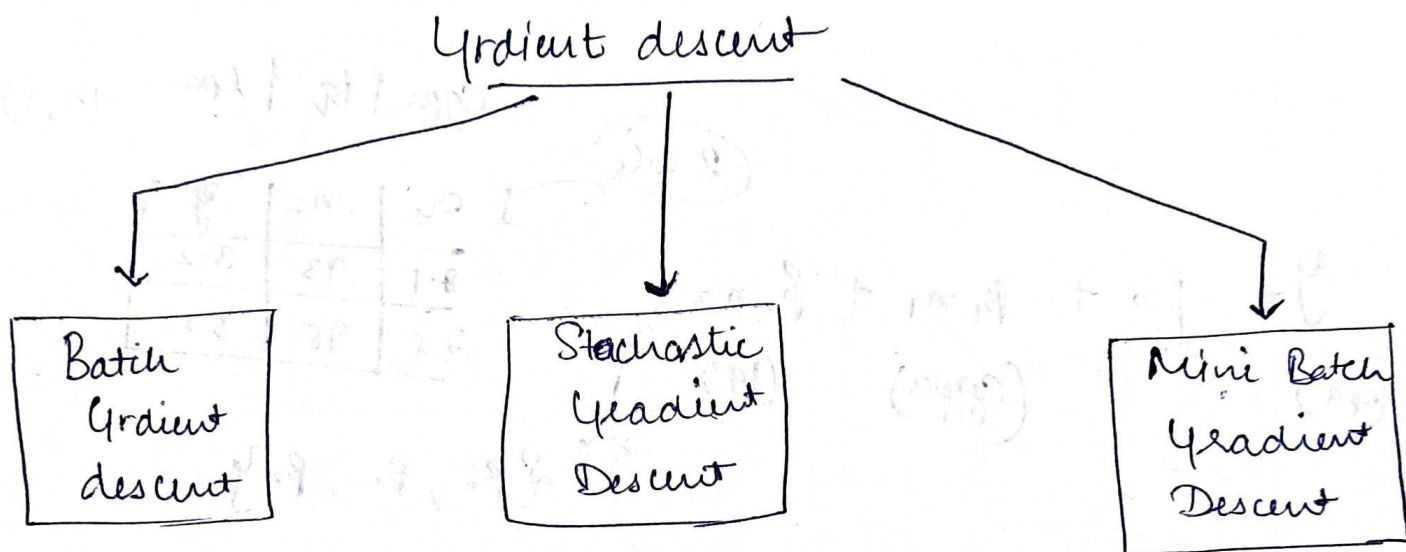
Non-convex



\* On this, we can not reach at global minima. we stuck at local minima.



# # Types of Gradient descent



$m=1, b=0$   
Iterate

Batch  $\rightarrow$  slow  
gd

$$\frac{dL}{dm_0} \rightarrow \frac{\partial L}{\partial m}$$

$$\rightarrow m_n = m_0 - \eta \times (\text{slope}) \quad m=0$$

$$b_n = b_0 - \eta \times (\text{slope}) \quad b=0$$

$$\frac{dL}{db} \rightarrow \frac{\partial L}{\partial b}$$

300

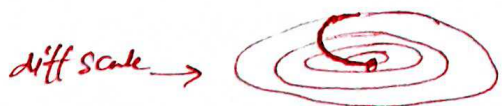
300 Rows ko dekhne ke bad update krte h

Stochastic GD  
Fast  
low computation  
1 row  
for larger data

mini batch GD  
both batch GD and Stochastic

Same Scale  $\rightarrow$  Data

bcz when data have different Scale so, it take much time to reach at ~~local~~ global minima.



# # Mathematical Formula

n-dimension dataset

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

(lpa) (cgpa) (iq)

cgpa | iq | lpa (2,3)

3 cols

$x_1$	$x_2$	$y$
8.1	93	3.2
7.5	95	3.8

$\{\beta_0, \beta_1, \beta_2\}$

① Step :- Random value

$$\beta_0 = 0, \beta_1; \beta_2 = 1$$

② Step epoch = 100,  $\eta = 0.1$

$$\beta_0 = \beta_0 - \eta \text{ slope}$$

$$\beta_1 = \beta_1 - \eta \text{ slope}$$

$$\beta_2 = \beta_2 - \eta \text{ slope}$$

$$L(\beta_0, \beta_1, \beta_2)$$

$$\frac{\partial L}{\partial \beta_0} = \frac{\partial L}{\partial \beta_1} \frac{\partial L}{\partial \beta_2}$$

if n-dim

(n+1)  $\beta_0 - \beta_n$

mean square error

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



row = 2, col = 2 + 1 = 3  
 input → output

$$= \frac{1}{2} \left[ (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 \right]$$

$$L = \frac{1}{2} \left[ (y_1 - [\beta_0 - \beta_1 x_{11} - \beta_2 x_{12}])^2 + (y_2 - [\beta_0 - \beta_1 x_{21} - \beta_2 x_{22}])^2 \right]$$

Annotations: Blue arrows point from 0 to the first term, from (-1) to the first term, and from 0 to the second term. A blue arrow points from  $\hat{y}_1$  to the first term, and from  $\hat{y}_2$  to the second term.

$x_1$	$x_2$	$y$
8.1 $x_{11}$	93 $x_{12}$	3.2
7.5 $x_{21}$	95 $x_{22}$	3.5

$$\frac{\partial L}{\partial \beta_0} = \frac{1}{2} \left[ 2(y_1 - \hat{y}_1) (-1) + 2(y_2 - \hat{y}_2) (-1) \right]$$

$$\begin{aligned} \hat{y} &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 \\ \hat{y}_1 &= \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} \\ \hat{y}_2 &= \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} \end{aligned}$$

$$\frac{dL}{d\beta_0} = -\frac{2}{2} [(y_1 - \hat{y}_1) + (y_2 - \hat{y}_1)]$$

$$= -\frac{2}{n} [(y_1 - \hat{y}_1) + (y_2 - \hat{y}_2) + (y_3 - \hat{y}_3) \dots (y_n - \hat{y}_n)]$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) = \frac{\partial L}{\partial \beta_0}$$

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

$$L = \frac{1}{2} [(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2]$$

$$L = \frac{1}{2} [(y_1 - [\beta_0 - \beta_1 x_{11} - \beta_2 x_{12}])^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2]$$

$$\frac{\partial L}{\partial \beta_1} = \frac{1}{2} [2(y_1 - \hat{y}_1)(-x_{11}) + 2(y_2 - \hat{y}_2)(-x_{21})]$$

$$\frac{\partial L}{\partial \beta_1} = \frac{-2}{n} [(y_1 - \hat{y}_1)(x_{11}) + (y_2 - \hat{y}_2)(x_{21}) + (y_3 - \hat{y}_3)(x_{31}) + \dots + (y_n - \hat{y}_n)(x_{n1})]$$

$$\frac{\partial L}{\partial \beta_1} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1}$$

$$x_{i1} \rightarrow \begin{matrix} x_{11} \\ x_{21} \\ x_{31} \\ x_{41} \\ \vdots \\ x_{n1} \end{matrix}$$

$$\frac{\partial L}{\partial \beta_2} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i2}$$

$$\frac{\partial L}{\partial \beta_m} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{im}$$

$\hat{y}$	$x_1$	$x_2$	$x_3$	$y$
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

$$\hat{y}_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13}$$

$$\hat{y}_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \beta_3 x_{23}$$

$$\beta_0 + \text{np.dot}(X_{\text{-train}}, \text{coef-})$$

$$(353, 10) \quad (10, 1)$$

$$(353, 1) \rightarrow \hat{y} = \text{np.dot}(\text{coef-}, X_{\text{-train}}) + \beta_0$$

$\uparrow$   
y-pred

$$\hat{y} = \beta_0 + [x_{11} \ x_{12} \ x_{13}] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

How to find coef-

$x_1$	$x_2$	$y$	$\hat{y}$
1	2	5	6
3	4	7	8

$$y - \hat{y} = \begin{bmatrix} 5 & 7 \end{bmatrix} - \begin{bmatrix} 6 & 8 \end{bmatrix} = \begin{bmatrix} -1 & -1 \end{bmatrix}$$

$$\frac{\partial L}{\partial \beta_1} = \frac{-2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1}$$

$$\frac{\partial L}{\partial \beta_1} = [y - \hat{y}] \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} -1 & -1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \end{bmatrix} \times -\frac{2}{4}$$

$$\frac{-2}{n} \begin{bmatrix} -1 & -1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$= 8$$



If you want this overall condition or find all derivatives.

$$[(y - \hat{y}) \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} x^{-\frac{2}{4}}$$

$$\frac{dL}{d\beta} \dots \frac{\partial L}{\partial \beta_{10}} = [(y_i - \hat{y}_i) x_{\text{train}}] x^{-\frac{2}{4}}$$

$y_{\text{train}}$

$$y_i = 353$$

transpose

2

$$(353, 1)$$

↓

$$(1, 353)$$

$$(353, 10)$$

$$(353, 10)$$

$$(353, 1)$$

$$(2, 1, 1)$$

$$(1, 10) \times \frac{-2}{4} = (1, 10)$$

↑

coef-der

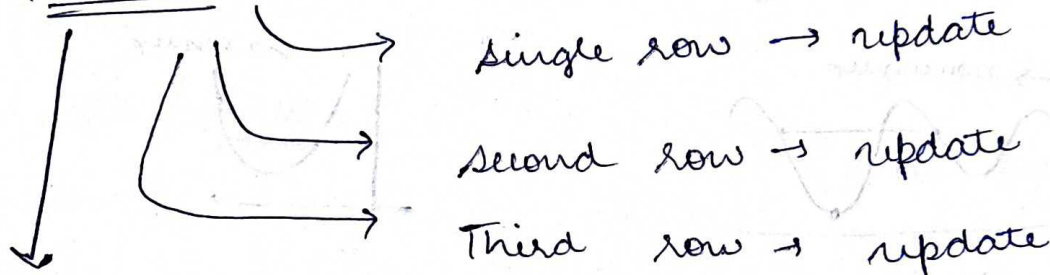
Disadvantage

- 1:- Very slow
- 2:- large calculation
- 3:- Hardware (less RAM and large data)

#

Stochastic Gradient Descent

\* Sklearn

Stochastic

\* Faster convergence

- less no. of epochs

n rows

1 epochs

n update

random select row

steady 50/100 %

correction

Time Comparisonno. of epochs  $\rightarrow$  fixed $e = 100$ 

batch

 $\rightarrow 100$ 

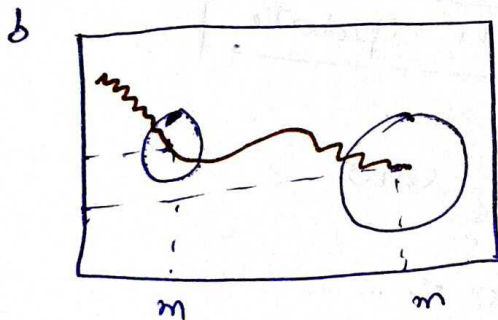
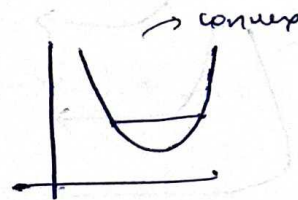
Stochastic

 $\rightarrow 100 \times n \rightarrow$  no. of rows $\rightarrow$  decrease no. of rowsno need that much epochs. ~~to~~

# # When to use stochastic GD

1) Big Data  $\rightarrow$  SGD

2) Non convex function



Problem :- <sup>still</sup> Fluctuate when near the destination

Solution :- learning schedule  $\rightarrow$  learning rate vary

$t_0, t_1 = 5, 50$

def learning-rate (t):

return  $t_0 / (t + t_1)$

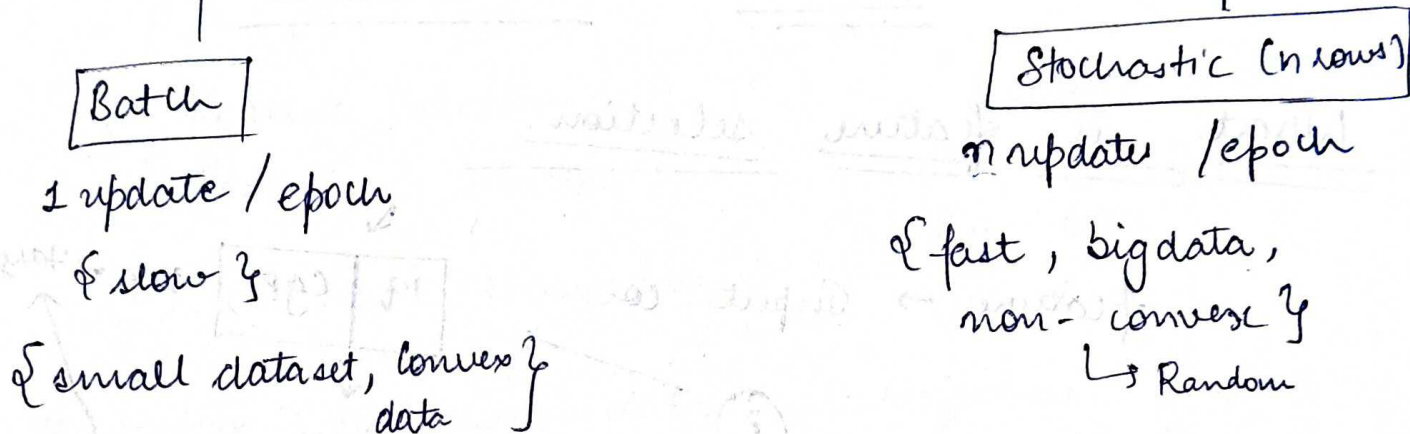
for i in range(epochs):

for in range(X.shape[0]):

$l_s = \text{learning-rate}(i * X.\text{shape}[0] + j)$



# Mini-Batch Gradient Descent



Mini batch

↳ group of rows

let  $n = 1000$

↳ batches  $\rightarrow 100$

10 batches

batches update / epoch

let  $n = 1000$

$n = 100 \rightarrow 10$  batches of 1000 rows  $\Rightarrow 10$  batches update / epoch