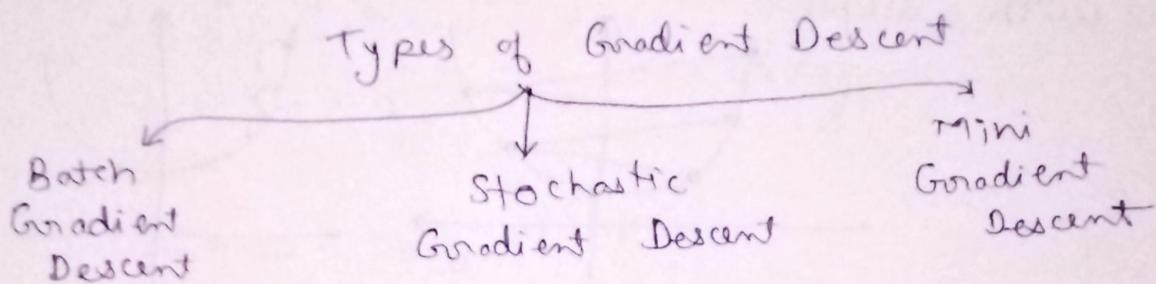


Date:  
09-02-24

# Gradient Descent

Gradient Descent:

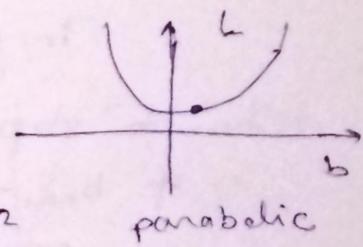
- Gradient Descent is a 1<sup>st</sup> order iterative optimization algorithm for finding a local minimum of a differentiable function.
- used in linear regression, logistic regression, t-sne



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

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

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



parabolic

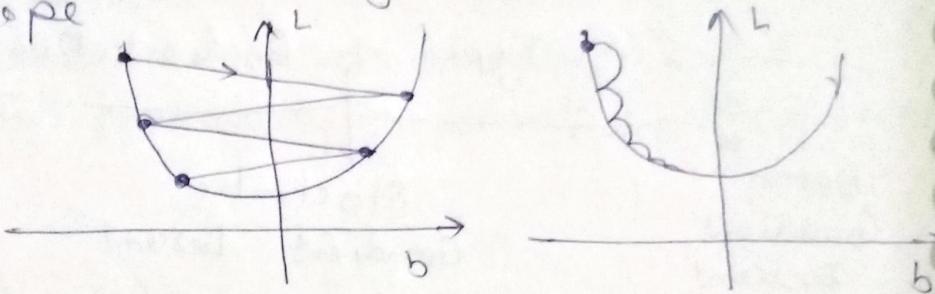
- + Here the relation b/w L and b is  $L \rightarrow b^2$
- Start w/ a value already known  $\frac{d}{db}$  above eqn  
in the given graph  $\frac{d}{db}$  at the start value  
निकालना  $\frac{d}{db}$  जैसे पर L is minimum

Approach of Gradient Descent:

step-1. Select the random value of b

step-2. find the slope. if slope = -ve more right (b increase) else more left (b decrement)

- In the Gradient Descent the idea is to take repeated steps in the opposite direction of the gradient (for approximate gradient) of the function at the current point, because this is the direction of steepest descent.
- for avoiding drastic change in the step size we have to multiply  $\eta$  (learning rate) with slope



$$\text{in general } \eta = 0.01$$

when to stop:

$$\text{if } b_{\text{new}} - b_{\text{old}} < 0.0001$$

or 1000 times iterations (epochs)

Mathematical Formulation:

let Here  $m$  is given

Step- start with a random  $b = b_0$   
Here  $\eta = 0.01$

for  $i$  in epochs:

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

let us assume  $b = 0$

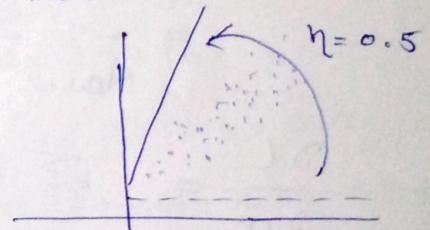
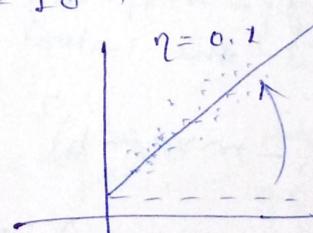
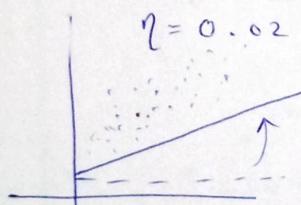
find slop at  $b = 0$

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

$$\begin{aligned}
 & \frac{d}{db} \sum_{i=1}^n (y_i - mx_i - b)^2 \\
 &= 2 \sum_{i=1}^n (y_i - mx_i - b) (-1) \\
 &= -2 \sum_{i=1}^n (y_i - mx_i - b) \\
 &\quad \text{above is equation of slope} \\
 &\quad \text{for } b=0 \\
 &= -2 \sum_{i=1}^n (y_i - mx_i) \\
 &= -2 \sum_{i=1}^n (y_i - 78.35x_i)
 \end{aligned}$$

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

Effect of Learning Rate:  
Here epoch = 10 for all three cases:



- gradient descent is work on all differentiable loss functions.

Calculating m and b both:

Steps:  
1. initial random value for m and b  
let  $m=1$  and  $b=0$

2. epochs = 100 and lr = 0.01

for i in epochs:

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

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

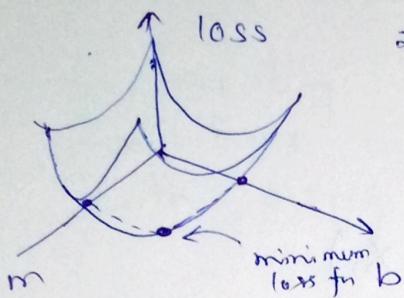
$$\text{cost/loss function } L = \sum (y_i - \hat{y}_i)^2$$

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

at the time of calculating  $b$   
Here  $m$  is constant but

Now  $L(m, b) = \text{cost/loss function}$

जब तक  $m$  और  $b$  की किसी विशेष मान नहीं लगती तब लॉस फंक्शन का मिनीमम वैल्यू होता है।



$b$  के विपरीत स्प्रेक्ट में स्लोप =  $\frac{\partial L}{\partial b}$

$m$  के विपरीत स्प्रेक्ट में स्लोप =  $\frac{\partial L}{\partial m}$

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

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

Now put the initial value of  $b$  ( $b=0$ )

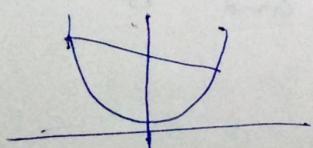
$$\frac{\partial L}{\partial m} = \sum (y_i - mx_i - b)^2$$

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

Now put initial value of  $m$  ( $m=1$ )

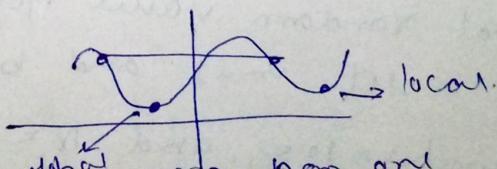
$$L = \sum (y_i - \hat{y}_i)^2 = \text{mean square error}$$

→ Convex function



one minima  
(global minima)

Convex loss function

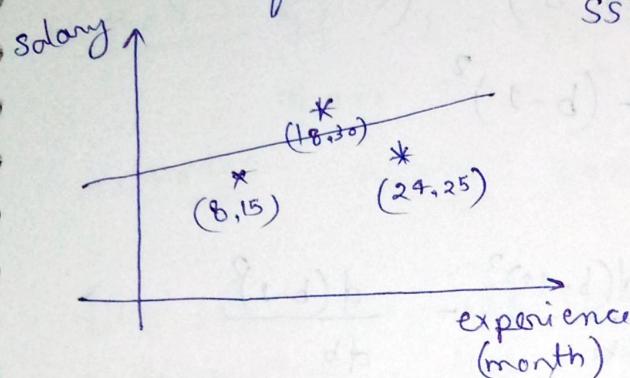


more than one minima

Non convex loss function

## Gradient Descent:

- This method is used when we have large number of dataset and large number of columns in dataset.
- This technique used in linear Regression, logistic regression, PCA and ~~Neuro~~ Neural Network
- it is optimization technique, parameter optimization technique.



$$SSR = (P_{y_1} - a_{y_1})^2 + (P_{y_2} - a_{y_2})^2 + (P_{y_3} - a_{y_3})^2$$

where  $P$  = predicted  
 $a$  = actual

$$\begin{cases} y = mx + b \\ \text{for } m=1 \quad b=0 \end{cases}$$

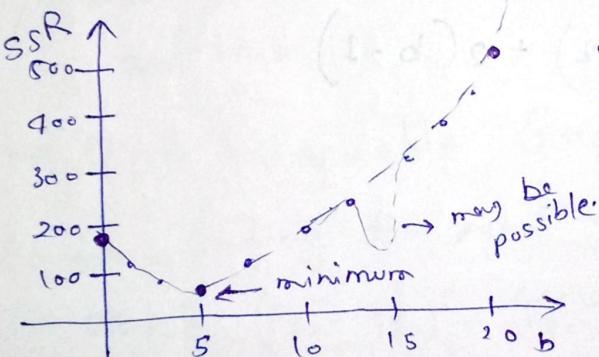
$$SSR = 49 + 44 + 1 = 194$$

$$\text{for } m=1 \quad b=5$$

$$SSR = 4 + 49 + 16 = 69$$

$$\text{for } m=1 \quad b=20$$

$$SSR = 169 + 64 + 361 = 594$$



if  $SSR$  is function of  $b$   
 $SSR = f(b) = b^2 - 4$

$$\frac{df}{db} = 2b \Rightarrow b=0$$

\*  $SSR$  is also called loss function

Step-1 find the equation of  $f(b)$

Step-2 differentiate

Step-3. put eqat to 0

then we get where curve goes minimum and if there is no any minima in curve then gradient descent calculate it

assume  $m=1$

$$SSR = (P_{y_1} - a_{y_1})^2 + (P_{y_2} - a_{y_2})^2 + (P_{y_3} - a_{y_3})^2$$

let us equation  $y = mx + b$

$$y = x + b$$

$$SSR = (P_{y_1} - 15)^2 + (P_{y_2} - 30)^2 + (P_{y_3} - 25)^2$$

$$SSR = ((8+b) - 25)^2 + ((18+b) - 30)^2 + ((24+b) - 25)^2$$

$$SSR = (b - 17)^2 + (b - 12)^2 + (b - 1)^2$$

$$SSR = f(b)$$

$$SSR' = \frac{df(b)}{db} = \frac{d}{db}(b-17)^2 + \frac{d}{db}(b-12)^2 + \frac{d}{db}(b-1)^2 \\ = 2(b-17) + 2(b-12) + 2(b-1)$$

$$SSR' = \frac{d f(b)}{db} = 6b - 60$$

$$\text{for } b = 0$$

$$SSR' = -60$$

learning Rate

$$\text{step size} = SSR' \times \text{learning rate}$$

$$= -60 \times 0.1$$

$$\text{step size} = 6 \cancel{-0.6}$$

$$\text{new intercept} = \text{old intercept} - \text{step size}$$

$$= 0 - (-6)$$

$$= 6$$

for  $b = 6$ ,  $SSR' = -24$

step size =  $-24 \times 0.1 = -2.4$

new intercept =  $6 - (-2.4) = 8.4$  (approx 9)

for  $b = 9$ ,  $SSR' = -6$

step size =  $-0.6$  new intercept = 9.6

\* How many times we do above steps:

method-1 if  $| > step | \leq 0.001$  (absolute)  
then stop the above step.

method-2 apply 1000 times above steps

fill Now in all in the above we take  
a assumption, the assumption is that the

~~m~~  $m=1$

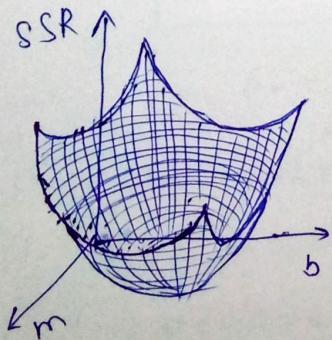
### Multiple Variable Gradient Descent:

for finding the equation of line  $y = mx + b$

$$SSR = (Py_1 - ay_1)^2 + (Py_2 - ay_2)^2 + (Py_3 - ay_3)^2$$

$$SSR = ((8m+b)-15)^2 + ((24m+b)-25)^2 + ((18m+b)-30)^2$$

$$SSR = f(m, b)$$



Step-1. find equation

Step-2. calculate derivative.

Step-3. putting value of m, b

Step-4. find minima

$$\frac{df(m, b)}{db} = 2((8m+b-15) \times \frac{d}{db}(8m+b-15)) + 2((24m+b-25)) + 2((18m+b)-30)$$

$$\frac{df(m, b)}{db} = 2((8m+b)-15) + 2((24m+b)-25) + 2((18m+b)-30)$$

Now,

$$\frac{df(m, b)}{dm} = 18((18m+b)-15) + 8((24m+b)-25) + 36((18m+b)-30)$$

let  $m = 1$  and  $b = 0$

$$(SSR)'_b =$$

$$\text{step size} = (SSR)'_b \times 0.1$$

$$NI = OI - \text{step size}$$

$$(SSR)'_m =$$

$$\text{step size} = (SSR)'_m \times 0.1$$

$$NI = OI - \text{step size}$$

- \* There are three variant of gradient descent, which differ in how much data we use to compute the gradient of the objective function. Depending on the amount of data, we make a trade-off between the accuracy of the parameter update and the time it takes to perform an update.

$$(d+n)7 = 922$$

storage bit = 8 qts

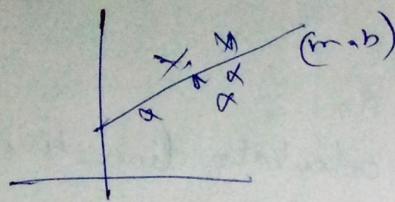
survived students = 8 qts

dead after cutting = 8 qts

min = 7 + 8 qts



# 1. Batch GD (ndim)



$$m=J \cdot b = 0$$

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

$$b_n = b_0 - \eta \times (\text{slope})_{b=0} \quad \rightarrow \frac{\partial L}{\partial b}$$

- \* Stochastic GD is suitable for large dataset
  - \* mini batch GD is intermediate b/w the Batch GD and stochastic GD
  - \* Batch GD used when we have the convex function. work only on small dataset
  - \* also called vanilla gradient descent.
- Here we work on more than two dimension data (ndim) example.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$\Rightarrow \{\beta_0, \beta_1, \beta_2, \beta_3\}$  ← coefficients.

→ Let us assume we have three column data set  
that is

cgpa $x_1$	iq $x_2$	LPA $y$
8.1	92	3.2
7.5	95	3.5

$x_{11} \quad x_{12} \quad y_1$   
 $x_{21} \quad x_{22} \quad y_2$

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

Steps:  
(1.) Start with Random values

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

(2.) epoch = 100, learning rate ( $\eta$ ) = 0.1

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

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

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

4-b graph.

(3)  $L(\beta_0, \beta_1, \beta_2)$  = loss function

$$\text{calculate } \frac{\partial L}{\partial \beta_0}, \frac{\partial L}{\partial \beta_1}, \frac{\partial L}{\partial \beta_2}$$

for n-dim we calculate derivative

from  $(n+1) \beta_0$  to  $\beta_n$

we know Loss function =  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  = MSE

for row = 2, cols = 2+1

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

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

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

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

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

Now, calculate slope-1, slope-2, slope-3.

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

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

for n rows

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

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

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

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

$$(II) \frac{\partial L}{\partial \beta_1} = \frac{1}{2} \left[ 2(y_1 - \hat{y}_1)(-x_{11}) + 2(y_2 - \hat{y}_2)(-x_{21}) + \dots + 2(y_n - \hat{y}_n)(-x_{n1}) \right]$$

for n-dimension. data

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

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

→ 1st column of data

$$(III) \frac{\partial L}{\partial \beta_2} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)(x_{i2}) = \text{slope}^{-3}$$

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

A problem with Batch GD:

1. let, n (number of rows) =  $10^5$

columns =  $10^2$

epoch =  $10^3$

~~total~~ total derivative calculate =  $10^{10}$

Hence high computation.  
the program will be slow

2. Hardware problem.

Code:

```
for i in range(nb_epochs):
    params_grad = evaluate_gradient(loss_function, data, params)
```

params = params - learning\_rate \* params\_grad

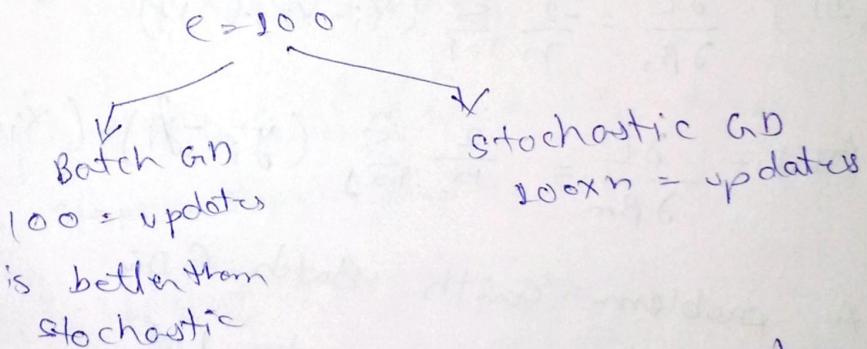
\* Here no. of epoch is equal to number of updates

## ② Stochastic Gradient Descent

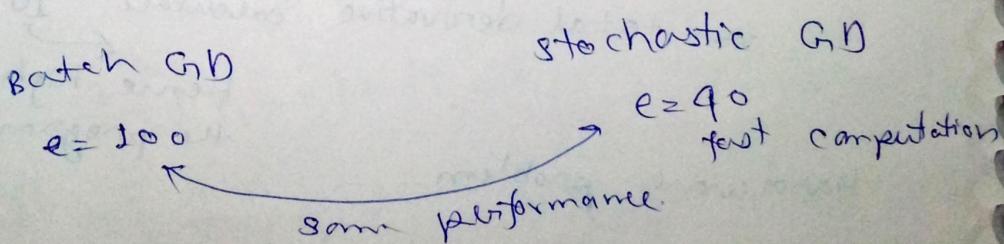
- it is better than the Batch GD because it is require less number of epochs.  
Here the update is done by row (randomly)
- it gives us faster convergence (less no. of epochs require)
- easy to handle.

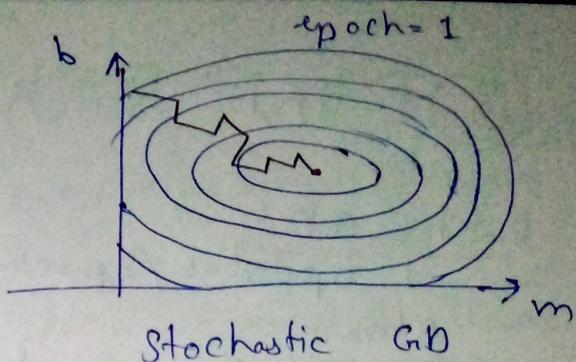
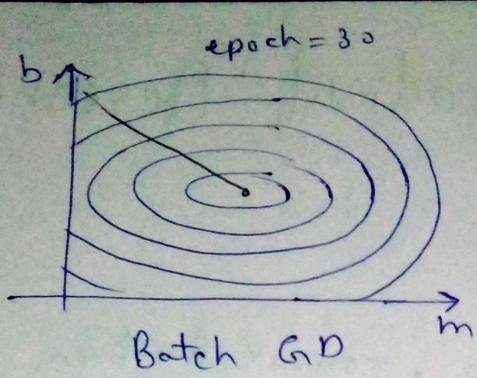
stochastic: having a random probability distribution or pattern that may be analysed statistically but may not be predicted precisely.

- it is not give steady solution.
- for the same number of epochs.



- but stochastic not require large number of epochs. (e.g. 100 here) it provide similar performance in less number of epochs.





When to use stochastic GD:

- when data is very big
- when we have non-convex functions

problem with stochastic GD:

- in the graph  $m$  and  $b$  value are very fluctuating
- To resolve this we use learning schedules

```

 $t_0, t_1 = 5, 50$ 
def learning_rate(t):
    return  $t_0 / (t + t_1)$ 
for i in range(epochs):
    for j in range(x.shape[0]):
        lr = learning_rate(i * x.shape[0] + j)
    
```

Here  $i$  = epoch.

e.g.  $n = 10^0$   
 $\text{epoch} = 1$        $l = 0.1$

### ③ Mini Batch Gradient Descent

Batch GD:

- 1 update/epoch
- slow
- useful for small data, convex function.

Stochastic GD:

- n(same) updates/epoch
- it is fast execution
- use full for big data, non convex function.
- not give us a optimal solution (randomness)

Mini Batch GD:

- Batch means group of same ( $n$ ) = 1000  
in one batch = 100  
hence number of batch = 10
- Batches update/epoch.
- It is intermediate b/w the Batch GD and Stochastic GD.
- It is the methods to reduce the randomness in stochastic GD.
- highly used in DL.