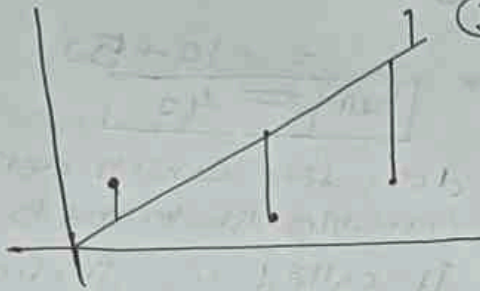# ✱ Gradient descent ✱

①Gradient descent is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function

OPtimization mean best Performance achive.

✳ intution : we want best fit line



✱we try to make minimum loss function

Value of

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

$$\hat{y_i} = mx_i + b$$

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

$\rightarrow L(m, b)$

Data

| cgpa | Package |
|------|---------|
| - | - |
| - | - |
| - | - |
| - | - |

THIS formula is depends on m & b

mean L is function depend on (m, b)
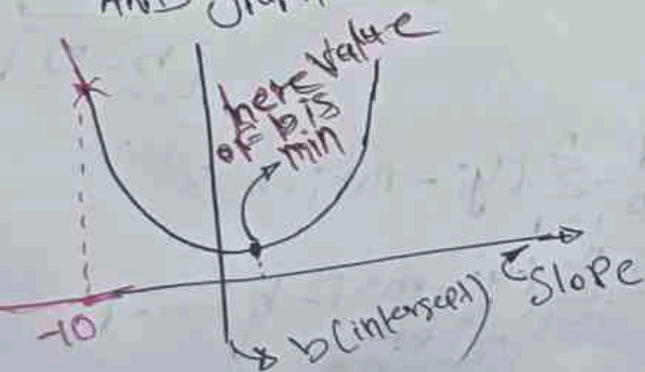
②suppose out data m value is known to us whichis

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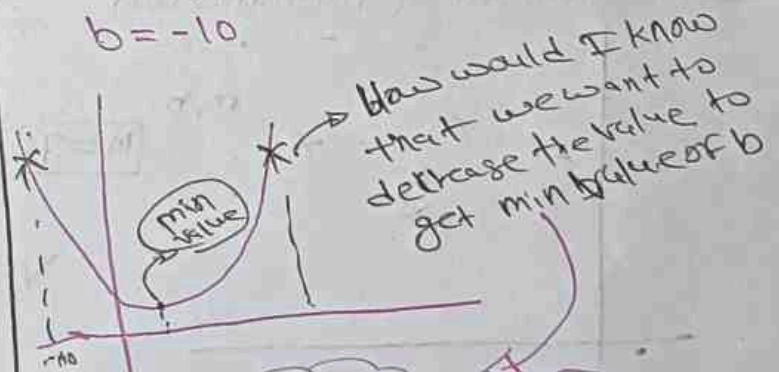
$$L = \sum_{i=1}^{4} (y_i - 78.35 x_i - b)^2$$

L(b)

$L \rightarrow b^2$

✱L and b relationship is depends on square
AND graph will be



here value of b is min

b(intersept) ← slope

✱step 1: Select random value of $b_{min}$

$$b = -10.$$



min value

→ How would I know that we want to decrease the value to get min value of b

ANS IS SLOP Finding

③understanding of slop
✱slop. It is how line is tilted
✱differentiation is a tool to find out slop when graph of function is curv

Ex: $y=x^2$

By differentiating

$$\frac{dy}{dx} = 2x$$

Slope(x) = $2x$

$x=1$ slope $2$
$x=2$ slope $4$
$x=-1$ slope $-2$

⊗ Very drastic changes in bnew values
To avoid this we multiply by learning rate $(\eta)$
Our equation becomes

$$bnew = bold - \eta \, slope$$

Ⓠ When to stop?

⊛ $bnew - bold \rightarrow$ is becomes very small $(0.0001)$
Then we can stop

⊛ Iteration $\rightarrow 1000, 100, -$
↳ epochs

⊛ Mathematical formulation



Here we find value of slope

$$L = \sum_{i=1}^{3} (y_i - \hat{y}_i)^2$$

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

It means
If we change b then How much change occurs in L

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

⊛ Slope value will tell us where is minima value of b

⊛ If slope is -ve then go forward
⊛ If slope is +ve then go backward

****

$$bnew = bold - slope$$

$bold = -10$ → assume slope
$bnew = -10 - (-50)$
$= -10 + 50$
$$bnew = 40$$

derivative which depends on two variables like m and b both then it called as _gradient_ or _approximate gradient_

⊛ gradient descent In terms of b

Step1: Start with a random value suppose $b = b'$
for i in epochs:   $\eta = 0.01$

$$bnew = bold - \eta \times slope$$

$$b = b$$

$m = 78.35$

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

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

no b term
it becomes 0    $-1$

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

↳ This is equation
of Slope
for b=0    (assume)
            m = 78.35

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

Slope (b=0)

Code Time :-

→This is use to make
dataset

From Sklearn.datasets import make_regression
import numpy as np

```
                          rows          shap of x        output feature
X,y = make_regression (n-samples=4, n_features=1, n_informative=1
                       n_targets =1, noise=80, random_state=13)
```
↓ how much
output

Randomness
NOT perfect
linear data

import matplotlib.pyplot as
    Plt.

Plt.scatter(x,y)



From Sklearn.linear_model import linearRegression

reg = LinearRegression()

reg.fit(x,y)

reg.coef_  ⟶ min slope (m)
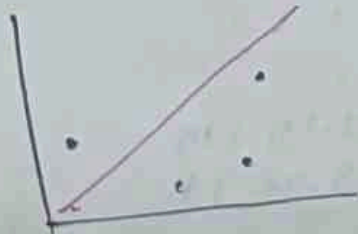
output⟶ 78.35

reg.intercept_

⟶ 26.1596  ⟶ min Intercept
              value (b)

Plt.scatter(x,y)
Plt.plot(x,reg.predict(x), color='red')

#Lets apply Gradient Descent by assuming slope is constant
m=78.35
# And let's assume the starting value of intercept b=0
   yPred = ((78.35 * x)+0).reshape(4)

m=78.35
b=0
loss_slope = -2 * np.sum(y- m*x.ravel()) - b)
loss-slope

      # xslope = loss_slope
         Now we want to do  bNew = bold - h_xslope
      h= 0.01
      b=0

$tepsize = 🔲 * loss-slope

   b=b-Stepsize
   b

→ 20.92

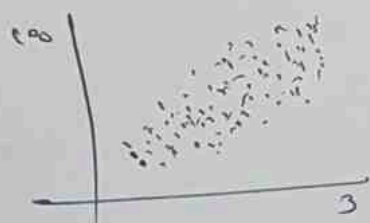⊛ Now we are making class to calculate only b for Now ⊛
   from sklearn.datasets import make-regression
   import matplotlib.pyplot as plt #AA
   import numpy as np.
   X,y = make-Regression (n-samples =100, n-Feature =1, n-informative
                                 n-targets =1, noise =20)      =1,

   plt.scatter (x,y)



#Actual result for comparing our result after making own class to
calculate b.
   from sklearn.linear-model import LinearRegression
   lr = LinearRegression()
   lr.fit(X,y)
   print(lr.coef_)
   print(lr.intercept_)

output:-
[29.19] m
[-3.35] b

m = 29.19

~~class~~
~~n~~ ~~def GD~~

```
class GDRegressor:
    def __init__(self, learning_rate, epochs):
        self.m = 29.19
        self.b = -120
        self.lr = learning_rate
        self.epochs = epochs
                    self.sqh
    def fit(x, y):
        # calculate the b using GD
        for i in range(self.epochs):
            loss_slope = -2 * np.sum(y - self.m*x.ravel() - self.b)
            self.b = self.b - (self.lr * loss_slope)
        print(self.b)

gd. GDRegressor(0.01, 100)
```

## ✴ Effect of learning rate ✴

$$\eta = 0.02 \qquad \eta = 0.1 \qquad \eta = 0.5$$

if very low
then very slow
- and more epochs

✴ The universality of Gradient Descent

✴ Now How to calculate m & b
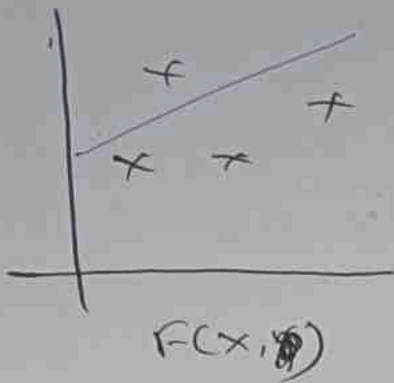   Step1:- initiate random vals for m and b
            m = 1 and b = 0
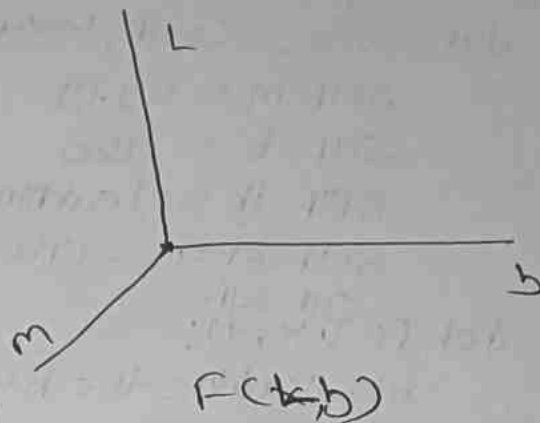
   ② epochs = 100   lr = 0.01
      for i in epochs:
            b = b - η.slope  ⎫ both slopes are different
            m = m - η.slope  ⎭

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

$$L(m,b)$$

$$F(x, y)$$

$$F(t, b)$$

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

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

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

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

$$\boxed{\text{Slope}_b = -2 \sum (y_i - mx_i - b)}$$

$$\boxed{= \text{Slope-b at } b = 0}$$

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

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

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

$$\boxed{\text{Slope-m at } m = 1}$$

⊛ Improvement In code ⊛

Now data is little bit change

X, Y = make-regression (n-samples = 100, n-features=1, n-informative=1, n-targets=1, noise=20, random-state=13)

```python
From Sklearn.linear_model import LinearRegression
lr = LinearRegression()
lr.fit(x,y)
Print (lr.coef_) #value of m
Print (lr.intercept_) #value of b
```
→ [27.8280] #m
   -2.2947   #b

```python
#This is calculating to Tally the ans with our class
from sklearn.model_selection import cross_val_score
np.mean (cross_val_score (lr,x,y, scoring ='r2', cv=10 ))
```

```python
class GDRegressor:
    def __init__ (self,learning_rate,epochs):
        self.m =100
        self.b = -120
        self.lr = learning_rate
        self.epochs =epochs

    def fit (self,x,y):
        for i in range (self.epochs):
            loss-slope-b = -2 * np.sum (y - self.m * x.ravel ()
                                        - self.b)

            loss_slope-m =-2 * np.sum((y - self.m *x.ravel()
                                       -self.b)* x.ravel())

            self.b = self.b - (self.lr * loss-slope-b )
            self.m = sdf.m -(self.lr * loss-slope-m)

        Print (self.m, self.b)

gd = GDRegressor (0.001,100)
gd.fit(x,y)
```
→ 27.8280.  -2.254
            ↑b
      ↓m

# Also making function to Predict

```
def Predict(self, x)
    return self.m * x + self.b
```

From sklearn.model-selection import train_test_split.

```
X_train, X_test, Y_train, Y_test = train_test-split (X, Y, test-
                              size = 0.2, random-state = 2)
```

```
Y_Pred = lr.Predict(x_test)
```

From sklearn.metrics import r2-score

r2-score (Y_test, Y_Pred)

→ 0.63 → this is ~~it~~ before class

## Now

```
gd.Fit (x_train, Y_train)
Y_Pred = gd.Predict (x_test)
```

From sklearn.metrics import r2-score

r2-score (Y_test, Y_Pred)

→ 0.63 → our class result

※ Impact on gradient Descent of Hyperparameter, loss Func &amp; Data.

① learning rate
② loss function
③ Data

① learning rate:-
  ⊙ low LR: ⊙ IF LR is 0.002
              ⊙ very small steps are geting
              ⊙ more epoch value
              ⊙ not much optimize

② Moderate LR : ⊙ Give optimize result
   (0.2)      ⊙ Not want more number
               of epochs
               ⊙ Get faster result


⊙ High LR : ⊙ It get bigger and
   0.80      bigger step
                         may
             ⊙ And some-times it ∧ not converge
             → to minimum value of loss.


② Effect of loss function :

⊙ $$L = \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$ → this loss function
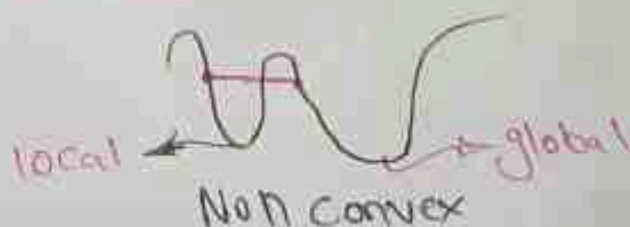                                             is convex function

⊛ Now what is convex function
   ⊙ If we draw a line between two points
   that line never cross the function that function is
Called  convex function.
         Ex:



Non convex Ex :-



local ←        → global
        Non convex

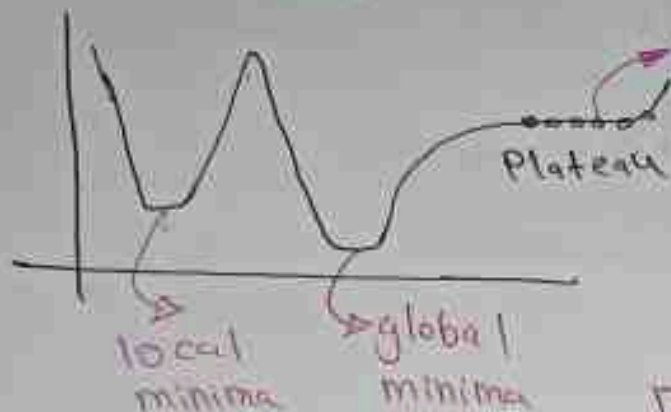   ⊙ In convex function there is only one minima
and that is global minima

   ⊙ In Non convex function there is more than one
minima we called them global minima and local minima

**※ In** Non convex Function there are two Problems

① ⊙ local minima is here :→ our algo. converge here and we should not go anyware.
   ⊙ more than 2 minima
   ⟹ Solution may be converge on local minima.

② Plateau :-



Plateau

local minima    global minima

To get out from Plateau there are we want to increase our epochs
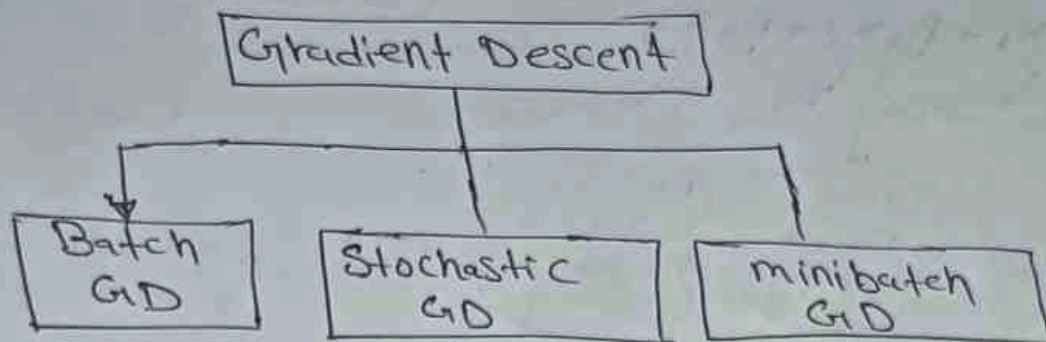
※ And more time taken to train

In 3D Plateau become Saddle Point

only Name change idea is Same

**✲ Effect of Data :-**
   ⊙ IF data is well scaled then our algorithm will fast converge towards the minima.

✳Types of gradient descent✳

$$\boxed{\text{Gradient Descent}}$$

$$\boxed{\begin{array}{c}\text{Batch}\\\text{GD}\end{array}} \quad \boxed{\begin{array}{c}\text{Stochastic}\\\text{GD}\end{array}} \quad \boxed{\begin{array}{c}\text{minibatch}\\\text{GD}\end{array}}$$

$$m_n = m_0 - h \times (slope)_{m=0}$$
$$b_n = b_0 - h \times (slope)_{b=0}$$
} We update the value of m & b by Viewing all data This is batch GD.

✳Stochastic GD✳
⊙ In this we update the values of m & b by Viewing Single row
⊙ This is Fast
⊙ Suitable for big/large data

✳minibatch GD✳
✳ fixed Batch size

✳ IF data is of 300 rows and batch size=30 our values of m & b will change after 30 rows not after full data.

✳ we already learn batch gradient descent to understand the ⊠ gradient descent for two Variables
we learn for: CgPa || IPa

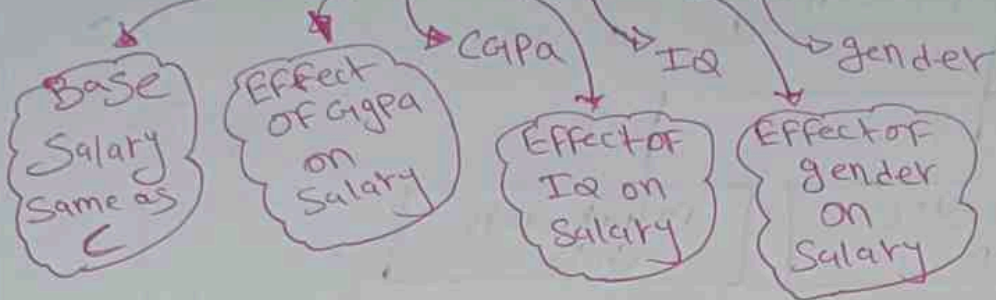Now we are learning Batch gradient descent for more than two Variable
CgPa | IQ | gender | cPa

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

# ✳ understand the equation ✳

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

- Base Salary Same as C
- Effect of cgpa on Salary ← cgpa
- Effect of IQ on Salary ← IQ
- Effect of gender on Salary ← gender

## ✳ N Dimensional Data ✳

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

Values required to find the Value of y Predict

## ✳ Mathematical Formulation ✳

| cgpa | iq | lpa |
|------|----|-----|
| $x_1$ | $x_2$ | y |
| 8.1 | 93 | 3.2 |
| 7.5 | 95 | 3.5 |

$$\underset{(lpa)}{y} = \beta_0 + \beta_1 \underset{(cgpa)}{x_1} + \beta_2 \underset{(lpa)}{x_2}$$

Step 1: Start with random value

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

Step 2: epoch = 100, lr = 0.1

$$\beta_0 = \beta_0 - \eta \, slope \longrightarrow \frac{\partial L}{\partial \beta_0}$$

$$\beta_1 = \beta_1 - \eta \, slope \longrightarrow \frac{\partial L}{\partial \beta_1}$$

$$\beta_2 = \beta_2 - \eta \, slope \longrightarrow \frac{\partial L}{\partial \beta_2}$$

Intercept

$$\boxed{L(\beta_0, \beta_1, \beta_2)}$$

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

$\downarrow$
(MSE)

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

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

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

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

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 8.1 | 93 | 3.2 |
| 1.5 | 95 | 3.5 |

$x_{12}$ , $x_{11}$ , $x_{21}$ , $x_{22}$ , $y_1$ , $y_2$

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

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

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

for N rows

$$\frac{\partial L}{\partial \beta_0} = \frac{-2}{n} \left[ (y_1 - \hat{y}_1) + (y_2 - \hat{y}_2) + (y_3 - \hat{y}_3) + \dots + y_n - \hat{y}_n \right]$$

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

$\rightarrow$ This is give the
value of slope

$$\beta_0 = \beta_0 - n(slope)$$

this
value
is calculated
by $\frac{\partial L}{\partial \beta_0}$

✱ Now calculating of $\beta_1$ slope ✱

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

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

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

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

$$\boxed{\frac{\partial L}{\partial \beta} - \beta_1 x_{11} = -x_{11}}$$

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

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

$x_{i1}$ is represent all values
of column 1

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 8.1 | 93 | 3.2 |
| 75 | 95 | 3.5 |

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

✱ for m cols ✱

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

```python
Class GDRegressor:
    def __init__(self, learning_rate=0.01, epochs=100):
        Self.coef_ = (None)
        Self.intercept_ = None
        Self.lr = learning_rate
        Self.epochs = epochs.

    def Fit(self, x_train, y_train):
        Self.intercept_ = 0
        Self.coef_ = np.ones(x_train.shape[1])
        for i in range(Self.epochs):
            # Updating all the Coef & Intercept
```

First understand the intercept - formula.

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

$(y_i - \hat{y}_i) \rightarrow$ this is nothing but our y-train

$(y_i - \hat{y}_i) \rightarrow x$

$$= -\frac{2}{n} \sum_{i=1}^{n} x$$

→ This is nothing but formula for calculating <u>mean</u>.

understand calculation to finding y hat

| x1 | x2 | x3 | 4 |
|----|----|----|---|
| -  | -  | -  | - |

→ for first row

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

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

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

we can calculate for all values in x_train By using my vectorization

$$Y = \beta_0 + np.dot(Coef_, xttrain)$$

$$\beta_0 + np.dot(X_train, coef_)$$

Formula becomes → (353,10) (10,1) x_train shape coef shape (353,1)

****

$y$-hat $= np.dot(x$-train, self.coef_) + self.intercept_

intercept_der $= -2 * np.mean(y$-train $- y$hat)

self.intercept_ $=$ self.intercept_ $- ($self.lr $*$ intercept_der)

# Now time to calculate coef_

for single: $\frac{\partial L}{\partial \beta_1}, \frac{\partial L}{\partial \beta_2}, \frac{\partial L}{\partial \beta_n}$

~~single1~~

※ How to find out this at game time?

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

$=$

| x 1 | x2 | y | $\hat{y}$ |
|---|---|---|---|
| 1 | 2 | 5 | 6 |
| 3 | 4 | 7 | 8 |

$$\left[ [y-\hat{y}] \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \right] \times \frac{-2}{n}$$

$[5 \ 7] \quad [6 \ 8]$

$y - \hat{y} = [-1, -1]$

$[-1, -1] \begin{bmatrix} 1 \\ 3 \end{bmatrix}$

$[-1 + -4]$

$\boxed{-8}$

$\frac{\partial L}{\partial \beta_1} = \cdots \frac{\partial L}{\partial \beta_{10}}$

$= [(y_i - \hat{y}_i) \times \text{train}] \times \frac{-2}{n}$

$(353,1)^T \qquad (353,10)$

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

$(1,10) \times \frac{-2}{n} = (1,10) \rightarrow$ coef_der

```python
        coef_der = -2 * np.dot((y_train - y_hat), x_train)/
                                                    x_train.shape[0]
        self.coef_ = self.coef_ - (self.lr * coef_der)


    def predict(self, x_test):
        return np.dot(x_test, self.coef_) + self.intercept_


gdr = GDRegressor(epochs=100, learning_rate=0.5)
gdr.fit(x_train, y_train)

152.0135   [14.3891   -173.72   491.54   323.91,
                    -39.32  -116.010,  -194.04  103.38
             451.63   97.57]


y_pred = gdr.predict(x_test)
r2_score(y_test, y_pred)
0.45
```

# ✳ Stochastic Gradient Descent ✳

## ✳ Problem with Batch GD ✳

$n = 1000$
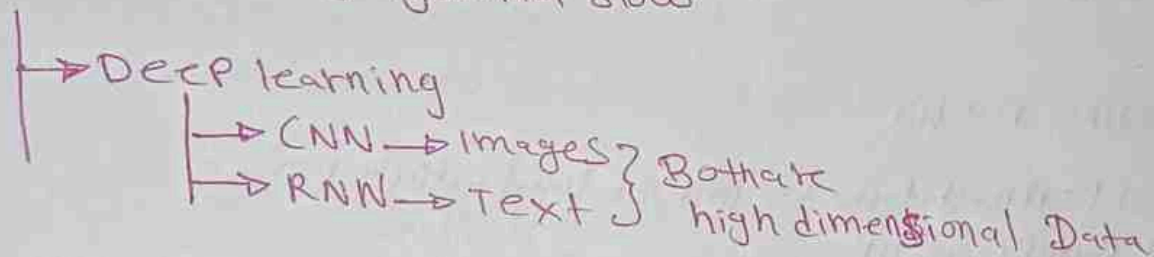
$cols = 5 \longrightarrow 6\ coefs$

$epoch = 50$

$50$

$1 \longrightarrow 1000$

$6 \longrightarrow 6000$ ⟶ derivatives want to calculate

$6000 \times 50\ \boxed{300000}$

⊙ It makes algorithm slow

⟶ Deep learning
⟶ CNN ⟶ Images ⎤ Both are
⟶ RNN ⟶ Text ⎦ high dimensional Data

✳ And In Deep Learning most of the time we are not use Batch GD

In coding Part we did vectorization to calculate y-hat

$y\text{-hat} = np.dot(x\text{-train}, self.coef\_) + self.intercept\_$

In this we used numpy and avoid loop

✳ When we are calculating y-hat we load full x-train data load on RAM.

✳ And if data is very large then it could give (error) Not exactly error but system not support (Hardware Problem)

✳ To resolve this all Problem we shifted toward Stochastic GD

✱ In Batch gradient descent Value of coef. is change after checking full data. only one time changes

✱ In Stochastic we change value of coef. by Seeing/Checking a single row.

———— ↳ row update

> ✱ In 1 single epoch n updates will do

⎰ ↳ faster convergom
⎱ ↳ row is select randonly
  ↳ Not give stady (ans) soln

✱ Now, coding Part ✱

From Skleann·datasets import load-diabetes.
import numpy as np
From Skleann.linear_model import LinearRegression
From Skleann.metrics import $r_2$-Score.
From sklearn.model-Selection import train-test-split.

X,Y = load-diabetes (return_x - y = True)

Print (x.Shape)
Print (y.Shape)

→ (442, 10)
→ (442, )

x (X, y = load-diabetes (return-x-y = True))

X-train, X-test, Y-train, Y-test = train-test-split (X, Y, test-siz = 0.2, random-state = 2)

reg = LinearRegression()
reg.fit (x-train, y-train) [-116 -205 46 516.69 340.62 -895 54
Print (Reg.coef_) ##m →  561.21 153·88 126·73 961·121 52·41]
Print (reg. intercept_) ## b → 151·88

Y_pred = reg.predict(x_test)
r2_score(y_test, y_pred)
→ 0.439938

⊛Now we make our class to predict coefficient⊛

Class SGDRegressor:

    def __init__ (self, learning_rate=0.01, epochs=100):
        self.coef_ = None
        self.intercept_= None
        self.lr = learning_rate
        self.epochs = epochs

    def fit (self, x_train, y_train):
        self.intercept_= 0
        sef.coef_= np.ones(x_train.shape[1])

        for i in range(self.epochs):
            for j in range(x_train.shape[0]):
                idx = np.random.randint(0, x.train.shape[0])
                y_hat = np.dot(x_train[idx], self.coef_) + self.intercept_

In BGD

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

1

collectively

Not required in SGD

*for SGD*    **

$$\boxed{\frac{\partial L}{\partial \beta_0} = -2(y_i - \hat{y}_i)}$$

intercept_der = -2 * (y_train[idx] - y_hat)
self.intercept_ = self.intercept_ - (self.lr * intercept_der)

Coef_der.

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

↓

In Batch GD

$$\frac{\partial L}{\partial \beta_i} = -2 \, (y_i - \hat{y}_i) \, x_{i1}$$

Coef_der = -2 * ((y_train[idx] - y_hat), X_train[idx])
Self.coef_ = Self.coef_ - (Self.lr * Coef_der)

X_train.shape
→ (353, 10)

⊛ Rule of matrix multiplication ⊛

A ⟶ (m×n) Shape

B ⟶ (n×p) Shape

C ⟶ (m×p) Shape

(353,10) * (1,1) = (10)

⊛ In this SGD Algorithm we converge to our minimum
loss value by giving less number of epochs.

Sgd = SGD Regressor (learning_rate=0.01, epochs=50)

Sgd. Fit(x_train, y_train)

→ 150.89,8 [54.28, -67.43  851.20  251.38  17.26
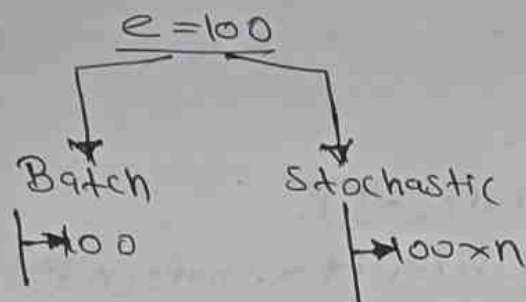-28.74 -165.85  123.27  314.84'  123.52]

y_Pred = Sgd. Predict (x_test)
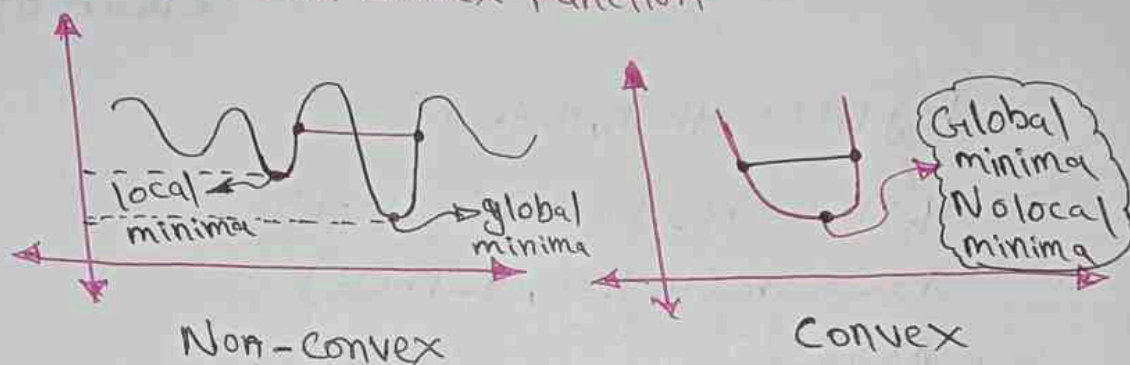
r2_score (y_test, y_pred)

→ 0.4325

⊛Time comparison⊛

no. of epoch → fixed

e = 100

Batch          Stochastic

├→100          ├→100×n

*When to use Stochastic GD

① Big data

② when we have non convex function



Non-Convex                    Convex

⊛ Because of Stochastic nature of SGD it could be ~~throw~~ get out from local minima (It have that velocity which is not stop it to local minima.)

⊛ Problem with SGD

⊙ In Stochastic GD, even after reaching the optimal solution, the algorithm continues to update parameters and does not Stop exactly at the minimum.

⊛ To solve this problem we use one concept called learning schedule.

⊛ learning schedule :- Learning rate schedule is a strategy in machine learning where the learning rate is changed (usually reduced) during training instead of keeping it constant, so that the model learn fast at the beginning and converges smoothly later.

**※Code:-**

```
t0,t1 = 5,50
def learning_rate(t):
    return t0(t#t1)
for i in range(epochs):
    for j in range(x.shape[0]):
        lr = learning_rate(i*x.shape[0]+j)
```

## ※ Using Sklearn SGD Regressor ※

From sklearn.linear_model import SGD Regressor

reg = SGDRegressor (max_iter=100, learning_rate = "constant", eta=0.01)
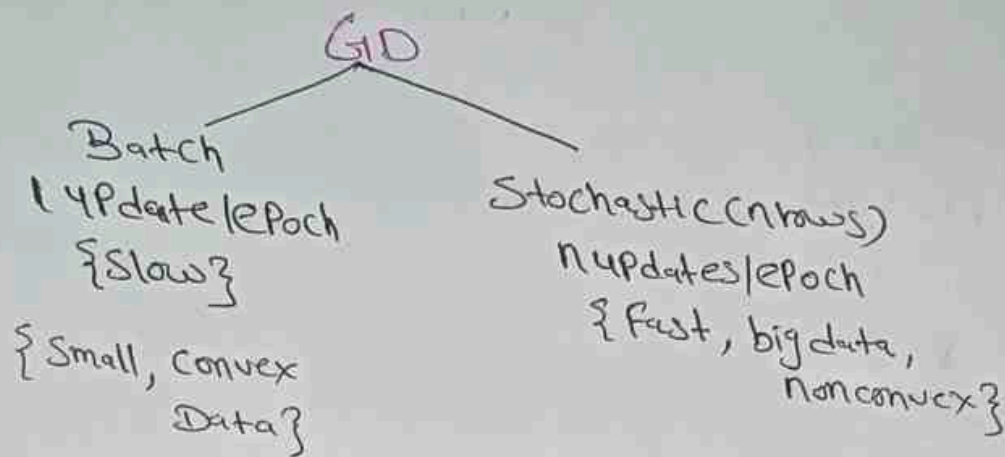
reg.fit(x_train, y_train)

y_pred = reg.Predict(x_test)

r2_score(y_test, y_pred)

→ 0.4305

# * Mini-Batch Gradient Descent *

GD

Batch
1 update/epoch
{Slow}

{Small, Convex
Data}

Stochastic (n rows)
n updates/epoch
{ Fast, big data,
nonconvex}

minibatch GD := Combination of Batch & Stochastic GD

Batch → Group of rows

n=1000
└→ Batches──(100) ──→ TB→10

Batches updates/epochs

1000
10
(Batches) ──→ 10 updates

## *Coding Part*

```
From skleath.datasets import load_diabetes
import numpy as np
From sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
From sklearn.model_selection import train_test_split
Import random
X,y = load_diabetes(return_X_y=True)
Print(x.shape)
Print(y.shape)
```

──→ 442, 10
(442,)

```
X_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
                                                     random_state=2)

reg = LinearRegression()
reg.fit(x_train, y_train)
Print(reg.coef_)
Print(reg.intercept_)
```

→ [-9.160 -205.46 516.68 340.62 -895.54 561.21
   153.88 126.73 86.12 52.41 151.88]

151.88

```
y_pred = reg.Predict(x_test)
r2score = (y_test, y_pred)
```

→ 0.4399

```
class MBGDRegressor:
    def __init__(self, batch_size, learning_rate=0.01,
                 epochs=100):

        self.coef_ = None
        self.intercept_ = None
        self.lr = learning_rate
        self.epochs = epochs
        self.batch_size = batch_size

    def fit(self, x_train, y_train):
        self.intercept_ = 0
        self.coef_ = np.ones(x_train.shape[1])

        for i in range(self.epochs):
            for j in range(int(x_train.shape[0]/
                              self.batch_size)):
                idx = random.sample(range(x_train.shape
                                          [0]),
```

```python
y-hat = np.dot(x_train[idx], self.coef_) +
                              self.intercept_
intercept_der = -2 * np.mean(y_train[idx] -
                              y-hat)
coef_der = -2 * np.dot((y_train[idx] - y_hat),
                              x_train[idx])
self.coef_ = self.coef_ - (self.lr * coef_der)

Print(self.intercept_, self.coef_)
def predict(self, x_test)
    return np.dot(x_test, self.coef_) + self.intercept_


mbr = MBGDRegressor(batch_size = int(x_train.shape[0]/10),
                    learning_rate = 0.01,
                    epochs = 100)


mbr.fit(x_train, y_train)
```

→154.8374   [38.40  -142.67  457.28  303.60  -17.99
              -85.81  -192.04  116.18  407.24  105.80]


```python
y_pred = mbr.predict(x_test)
r2_score(y_test, y_pred)
```

→ 0.45188

⊛ How to use MBGDR in sklearn ⊛

```python
from sklearn.linear_model import SGDRegressor
sgd = SGDRegressor(learning_rate = 'constant', eta0 = 0.2)
# partial_fit(x, y, sample_weight = None)
    ∟using this we can pass subset of x_train,
                                      y_train)
```

```
batch_size = 35
  for i in range(100)
      idx = random. Sample (range (x_train .shape [0]), batch_size)
      Sgd . Partial_ Fit (x_train [idx], y_train [idx])

      Sgd. Coef_
  → [49.19 , -67.84 , 338.57 , 247.97, 25.30 , -24.71,
                -155.45, 116.19, 312.91, 133.36]

      Sgd. intercept_
       → 148-61
  y_Pred = sgd.Predict (x_test)
  r2_score (y_test, y_Pred)
       → 0.4271
```