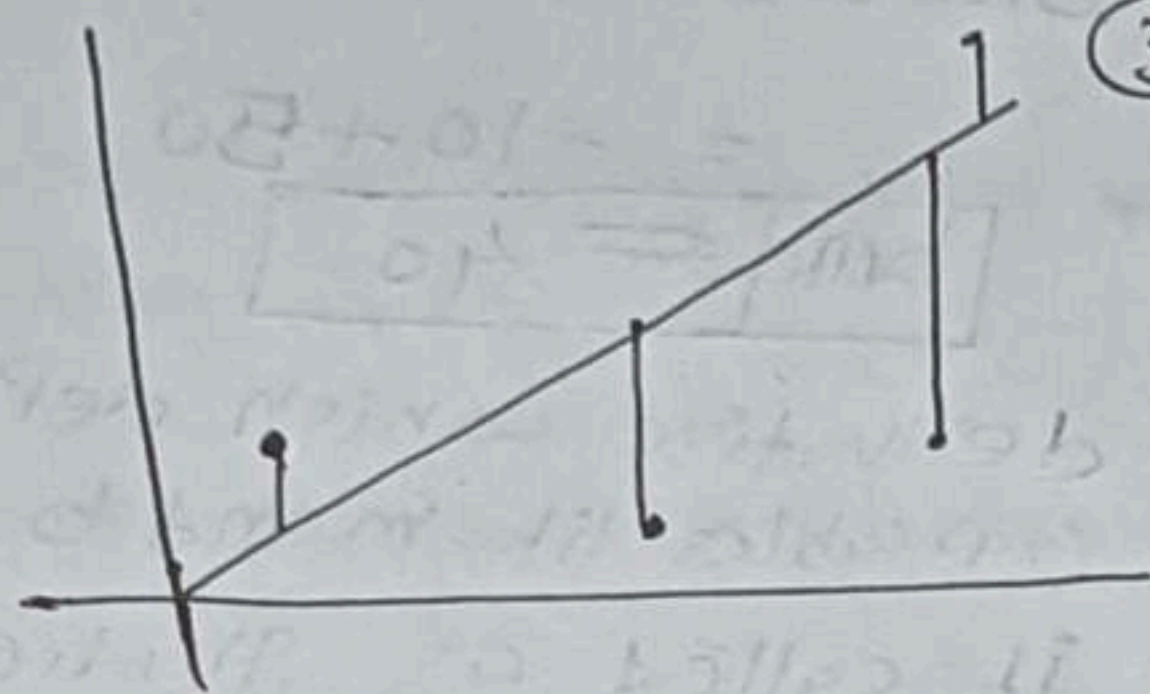


* Gradient descent *

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

Optimization mean best performance achieve.

* intuition: we want best fit line



* we try to make minimum loss function

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

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

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

Data

cgpa	Package
-	-
-	-
-	-
-	-

This formula is depends on m & b

$L(m, b)$

mean L is function depend on (m, b)

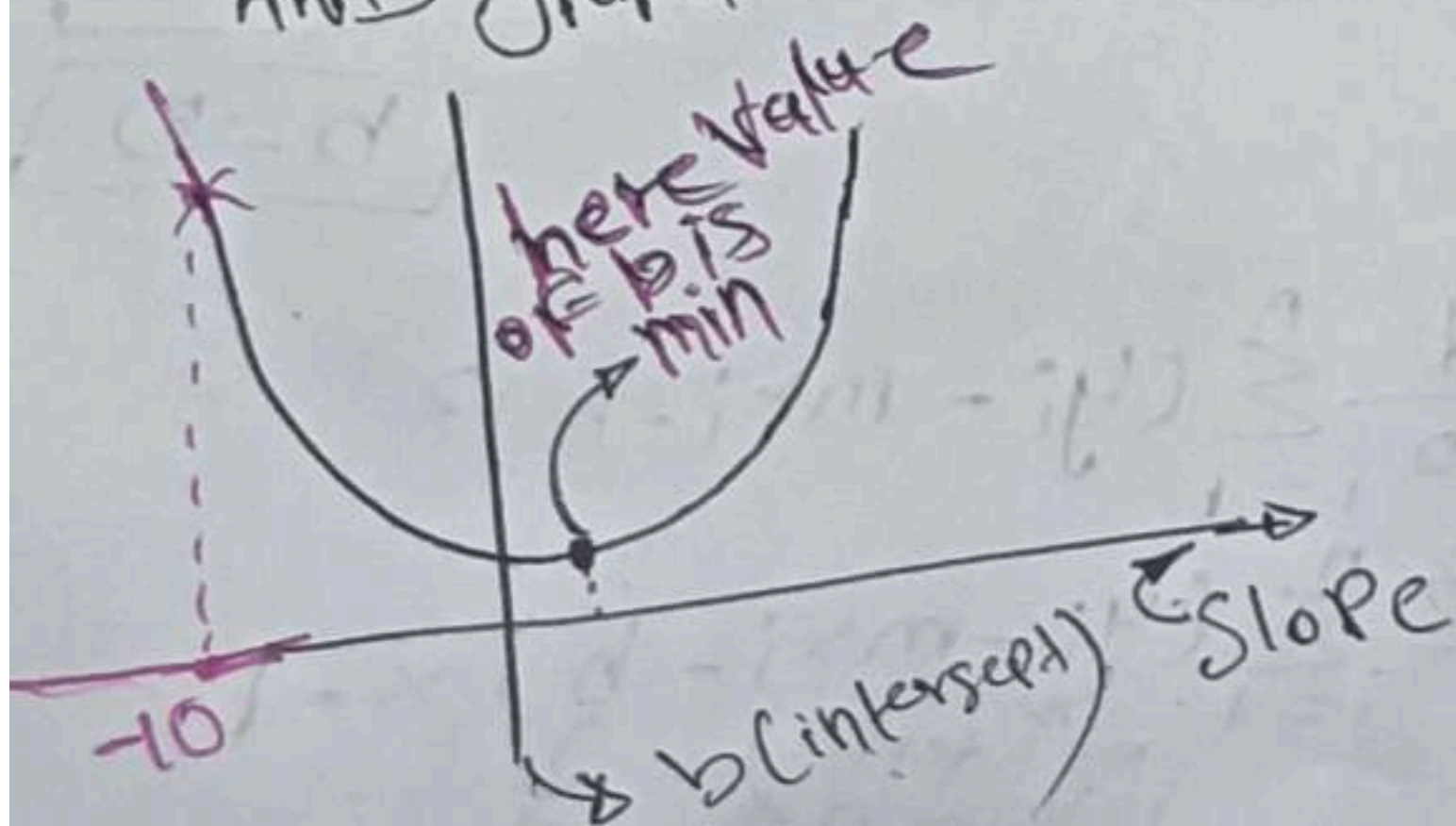
* Suppose our data m value is known to us which is

$$m = 78.35$$

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

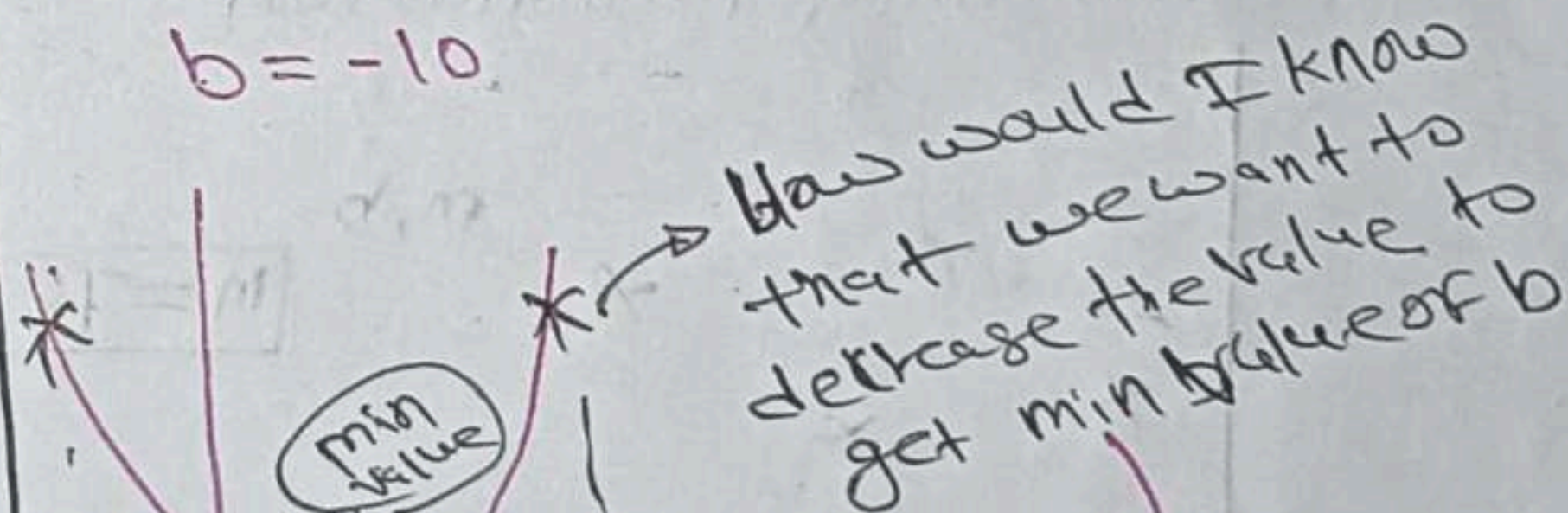
$$L(b) \quad L \rightarrow b^2$$

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



* Step 1: Select random value of b_{min}

$$b = -10$$

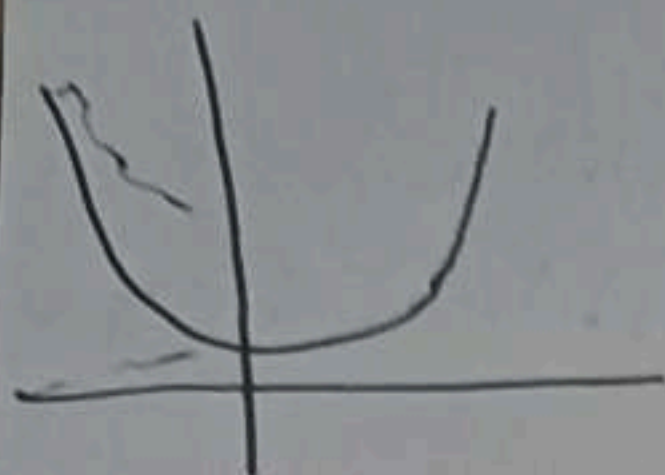


ANS IS Slope Finding

Understanding of Slope

* Slope: It is how line is tilted
Differentiation is a tool to find out Slope when graph of function is curve

Ex: $y = x^2$



By differentiating

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

$$\text{Slope}(x) = 2x$$

$$x = 1 \text{ Slope} = 2$$

$$x = 2 \text{ Slope} = 4$$

$$x = -1 \text{ Slope} = -2$$

Very drastic changes in b_{new} values

To avoid this we multiply by learning rate (η)

Our equation becomes

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

Slope value will tell us where is minimum value of b

IF Slope is -ve then go forward

IF Slope is +ve then go backward

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

$$b_{\text{old}} = -10$$

assume slope

$$b_{\text{new}} = -10 - (-50)$$

$$= -10 + 50$$

$$b_{\text{new}} = 40$$

derivative which depends on two variables like m and b both then

it called as gradient or

Approximate gradient

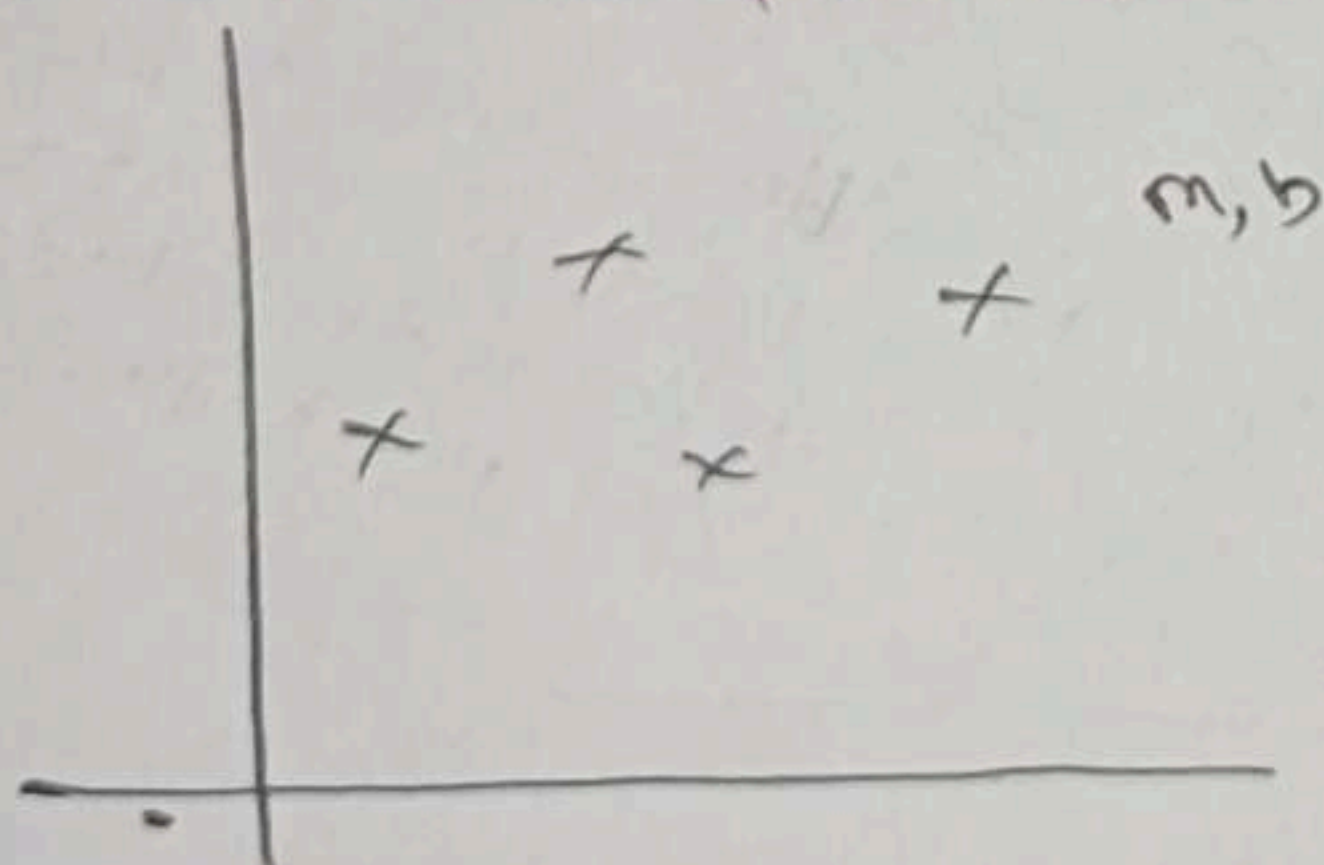
When to stop?

$b_{\text{new}} - b_{\text{old}}$ is ~~very~~ becomes very small (0.0000...1)

Then we will stop

Iteration \rightarrow 1000, 100, ...
Epochs

Mathematical Formulation



m, b

$$m = 78.35$$

gradient descent
In terms of b

Step 1: Start with a random value suppose $b = b'$

for i in epochs:

$$\eta = 0.01$$

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

$$b = 0$$

Here we find value of Slope

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

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

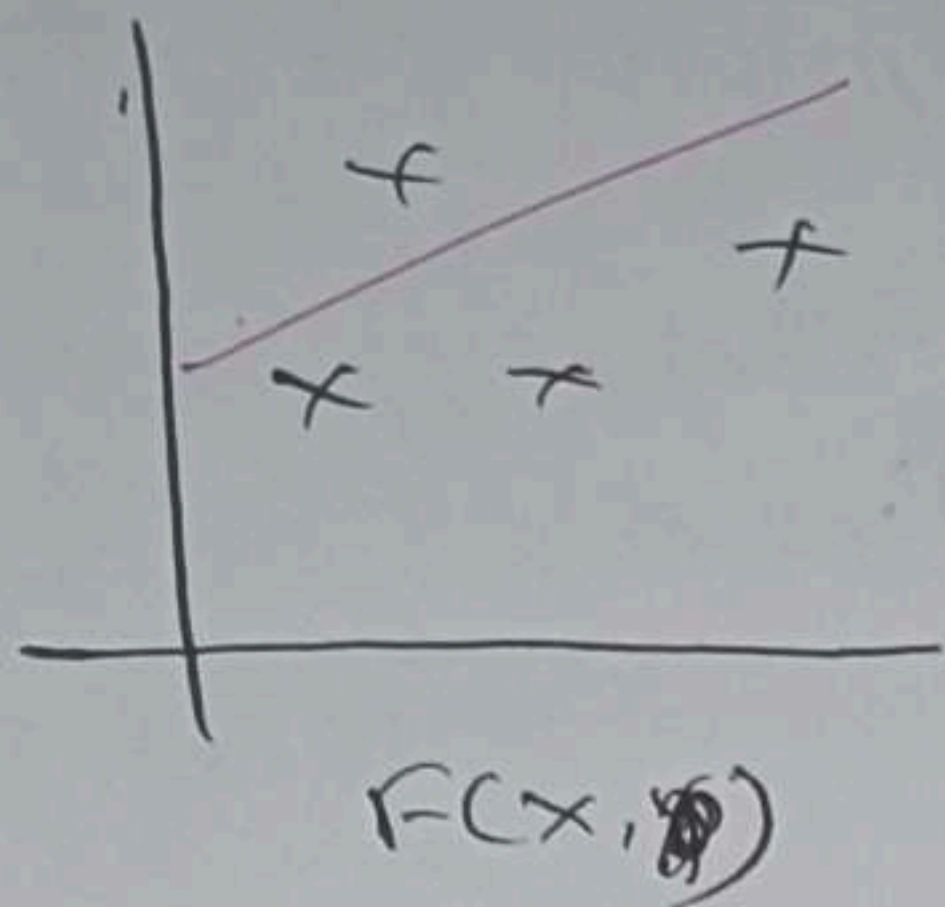
It means
If we change b then how much change occurs in L

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

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

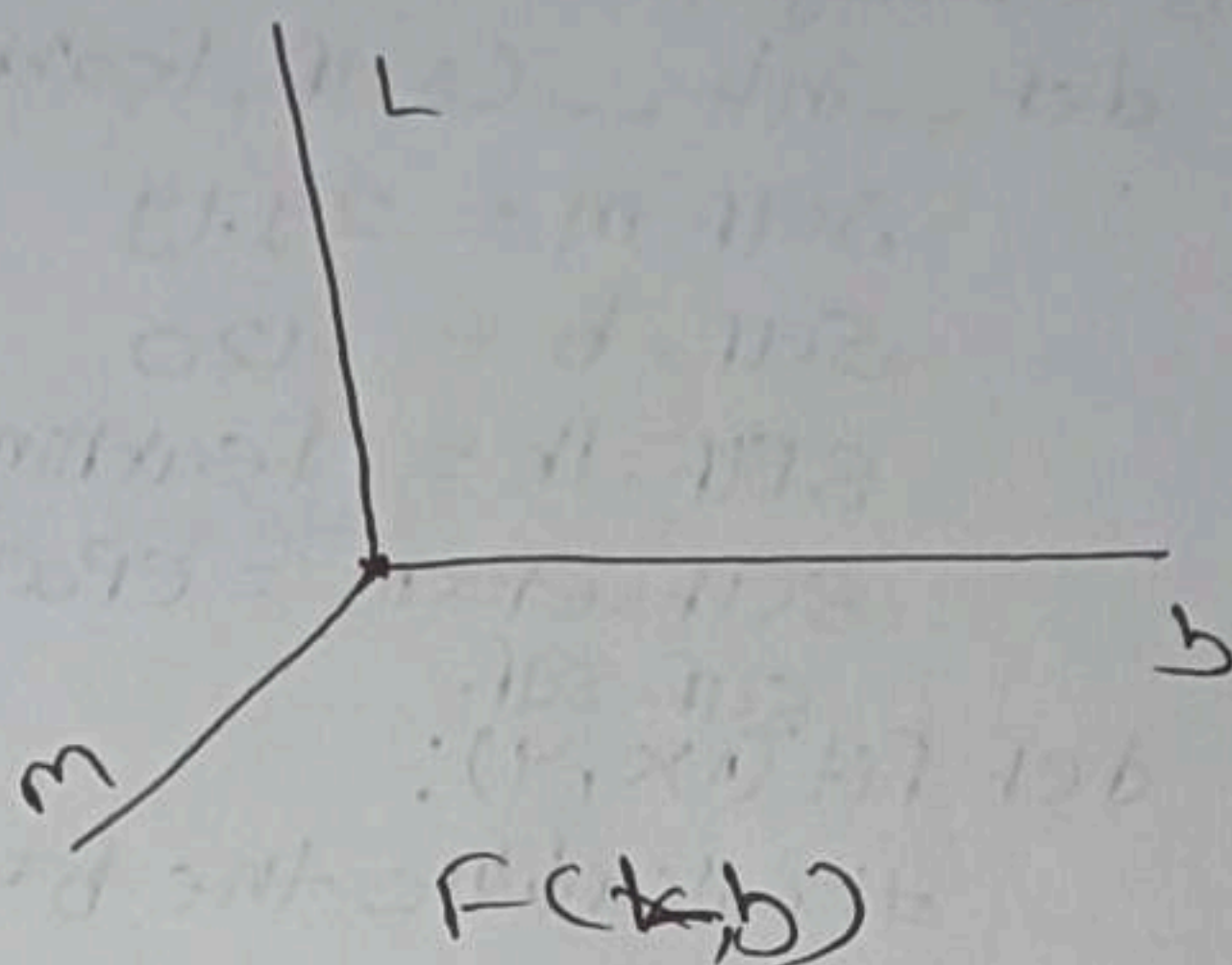
$$2 \sum_{i=1}^n (y_i - mx_i - b) \times -1$$

No b term it becomes 0 $\rightarrow -1$



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

$$L(m, b)$$



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

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

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

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

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

$$= \text{Slope}_b \text{ at } b=0$$

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

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

$$= -2 \sum x_i$$

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

$$\text{Slope}_m \text{ at } m=1$$

* Improvement in code *

Now data is little bit change

$X, y = \text{make_regression}(n\text{-Samples}=100, n\text{-Features}=1,$
 $n\text{-informative}=1, n\text{-target}=1,$
 $\text{noise}=20, \text{random_state}=13)$

$m = 29.19$
~~class~~
~~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.suf}
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 *

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

if very low
then very slow
and more epochs

* The universality of Gradient Descent

* Now How to calculate m & b :

Step 1: initiate random vals for m and b

$m = 1$ and $b = 0$

② epochs = 100 lr = 0.01

for i in epochs:

$b = b - \eta \text{slope}$ } both slopes are different
 $m = m - \eta \text{slope}$

Lets apply Gradient Descent by assuming slope is constant

$$m = 78.35$$

And let's assume the starting value of intercept $b = 0$

$$y_{pred} = (78.35 * x) + 0, \text{ reshape(1)}$$

$$m = 78.35$$

$$b = 0$$

$$\text{loss_slope} = -2 * \text{np.sum}(y - m * x.\text{reshape(1)} - b)$$

$$\text{loss_slope}$$

$$\# \text{XSlope} = \text{loss_slope}$$

$$\text{Now we want to do } b_{\text{new}} = b_{\text{old}} - \eta * \text{XSlope}$$

$$\eta = 0.01$$

$$b = 0$$

$$\text{StepSize} = \eta * \text{loss_slope}$$

$$b = b - \text{StepSize}$$

$$b$$

$$\rightarrow 20.92$$

* Now we are making class to calculate only b for now *

from sklearn.datasets import make_regression

import matplotlib.pyplot as plt

import numpy as np

$X, y = \text{make_regression}(n_samples=100, n_features=1, n_informative=1, n_targets=1, \text{noise}=20)$

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:-

$[29.19] m$

$[-3.35] b$

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

This is equation of slope

for $b=0$

assume

$m = 78.35$

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

Slope ($b=0$)

Code Time:-

From sklearn.datasets import make_regression
import numpy as np

This is use to make dataset

$X, y = \text{make_regression}(\text{rows} \rightarrow \text{n_samples}=4, \text{shape of } x \rightarrow \text{n_features}=1, \text{output feature} \rightarrow \text{n_informative}=1, \text{n_targets}=1, \text{Noise}=80, \text{random_state}=13)$

import matplotlib.pyplot as plt.

How much output

Randomness
Not perfect linear data

plt.scatter(X, y)



From sklearn.linear_model import LinearRegression

reg = LinearRegression()

reg.fit(X, y)

reg.coef_ \rightarrow min slope (m)

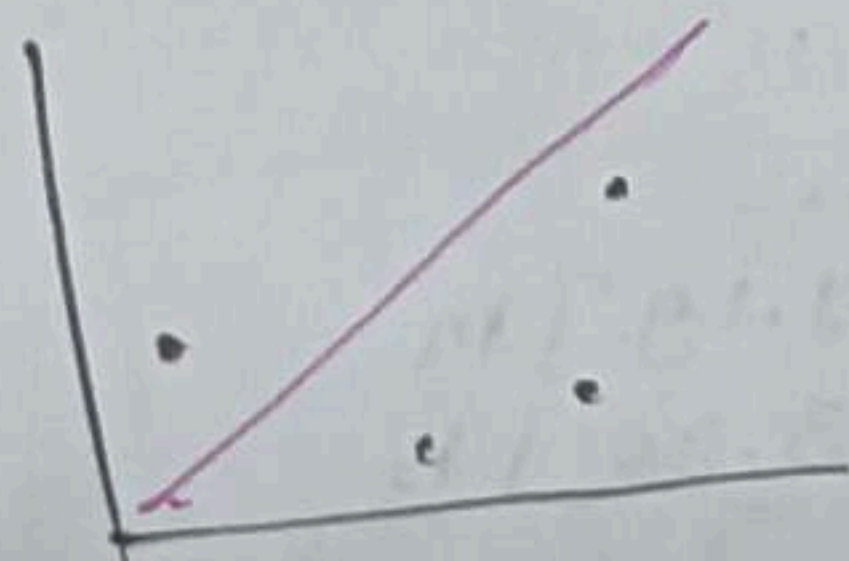
output $\rightarrow 78.35$

reg.intercept_

$\rightarrow 26.1596$ \rightarrow min intercept value (b)

plt.scatter(X, y)

plt.plot(X, reg.predict(X), color='red')




```
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
```

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))
```

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 = self.m - (self.lr * loss_slope_m)
```

```
    print(self.m, self.b)
```

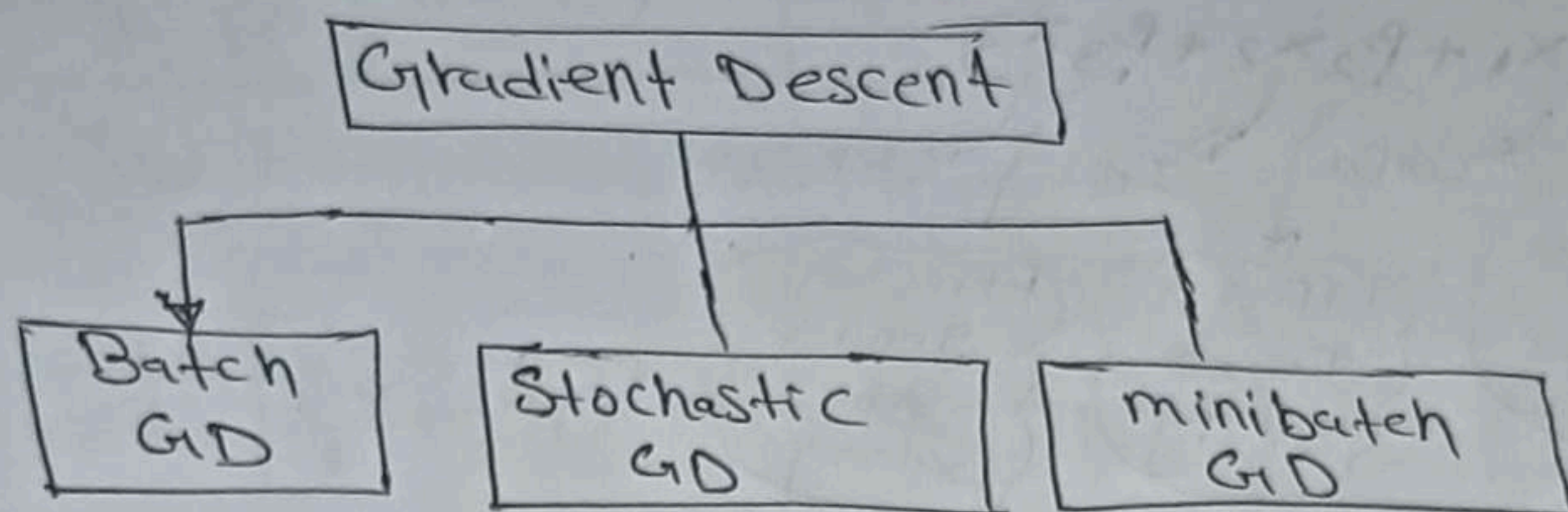
```
gd = GDRegressor(0.001, 100)
```

```
gd.fit(x, y)
```

```
→ 27.8280. -2.254
```

m b

Types of gradient descent



$$\left. \begin{aligned} m_n &= m_0 - \eta \times (\text{slope})_{m=0} \\ b_n &= b_0 - \eta \times (\text{slope})_{b=0} \end{aligned} \right\} \begin{array}{l} \text{We update the value of} \\ m \text{ \& } b \text{ by viewing all data} \\ \text{This is batch GD.} \end{array}$$

Stochastic GD

① In this we update the values of m & b by viewing single row

① This is Fast

① Suitable for big/large data

mini batch GD

*Fixed Batch Size

*IF data is of 300 rows and batch size = 30
Our values of m & b will change after 20 rows not after full data.

*We already learn batch gradient descent to understand ~~the~~ gradient descent for two variables

we learn for: GyPa / IPa

Now we are learning Batch gradient descent for more than two variable

GyPa / IQ / gender / CPA

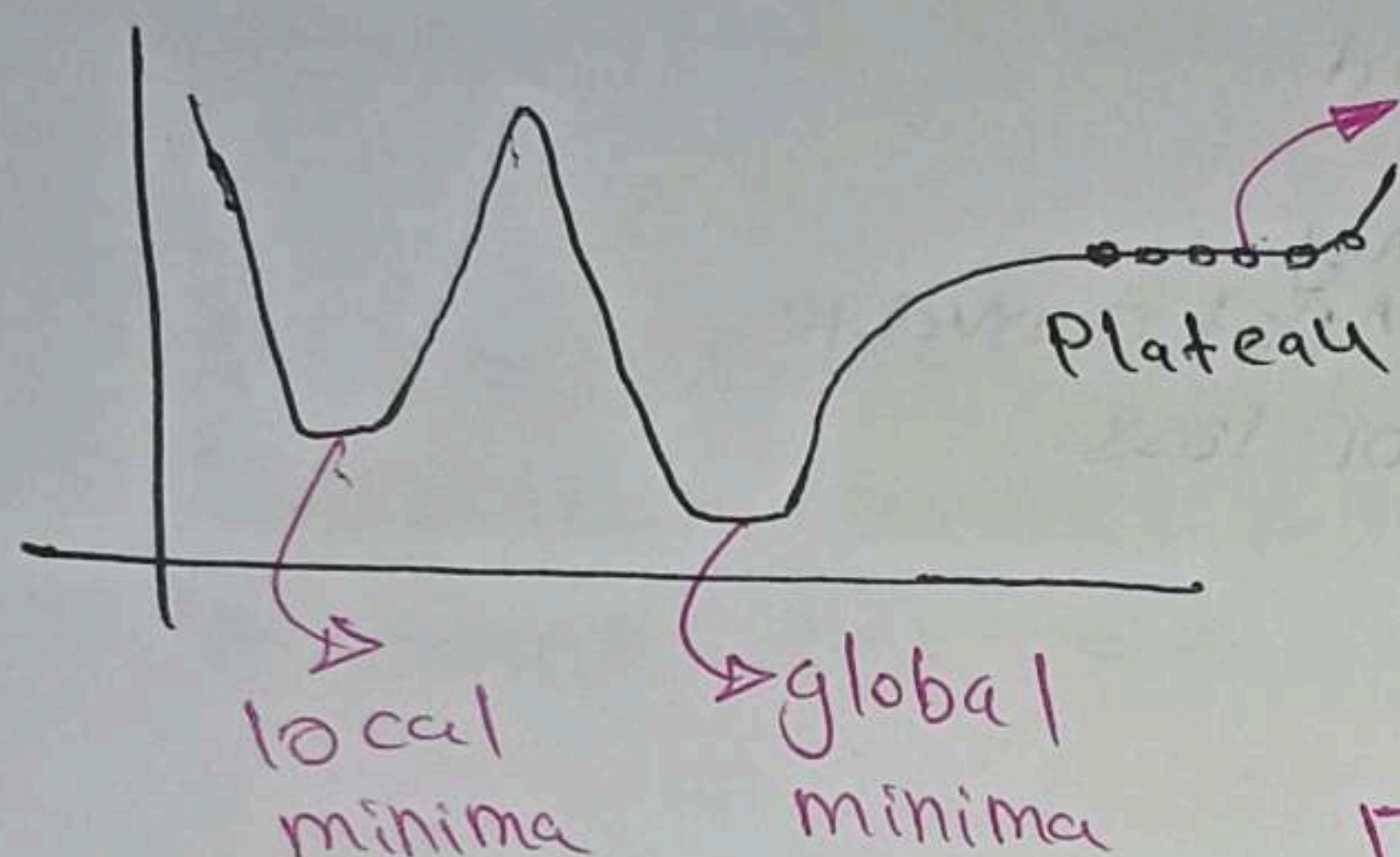
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

* In Non convex Function there are two Problems

- ① local minima is here : our algo. converge here and we should not go anyware.
 - ② more than 2 minima

4 solution may be converge on local minima.

② Plateau:-



To get out from Plateau there are we want to increase our epochs

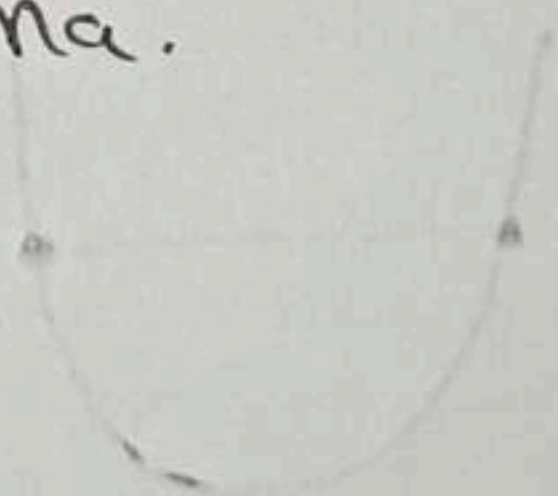
* And more time taken to train

In 3D Plateau become Saddle Point

only Name change
idea is same

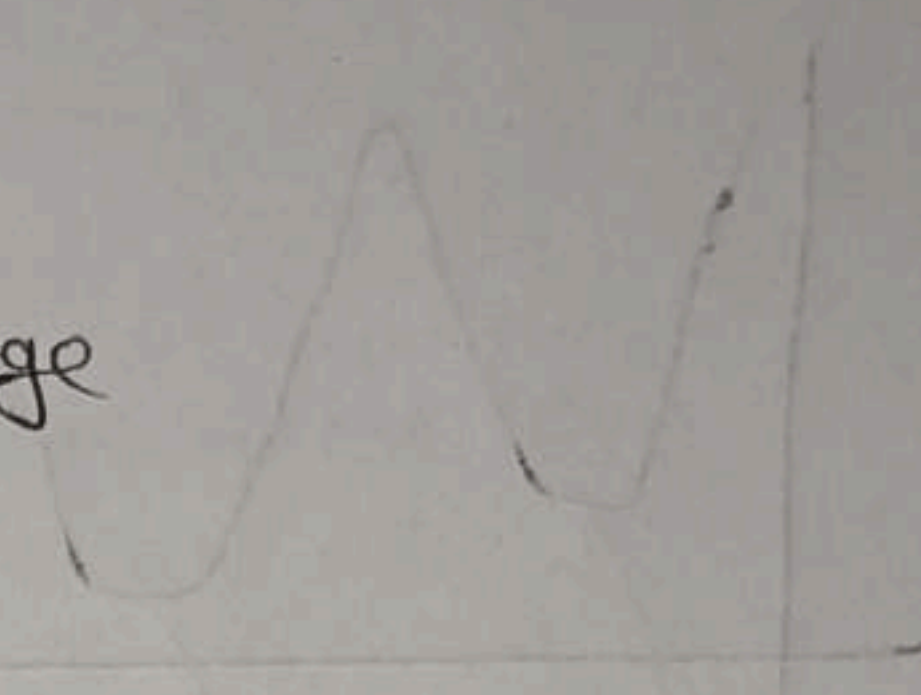
* Effect of Data:

① IF data is well scaled then our algorithm will fast converge towards the minima.



- ② Moderate LR:
 - ① Give optimize result (0.2)
 - ② Not want more number of epochs
 - ③ Get faster result

- ③ High LR:
 - ① It get bigger and bigger step 0.80
 - ② And sometimes it ^{may} not converge to minimum value of loss.



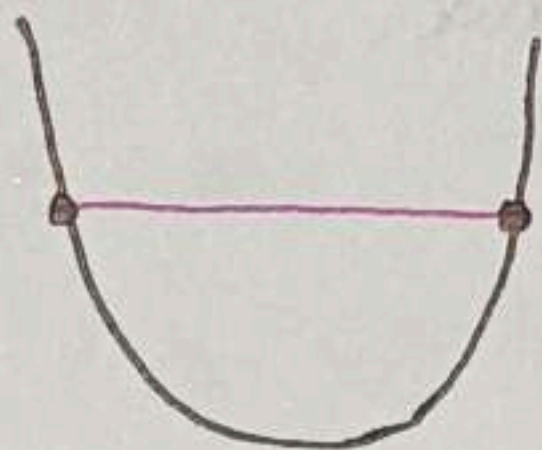
② Effect of loss function:

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

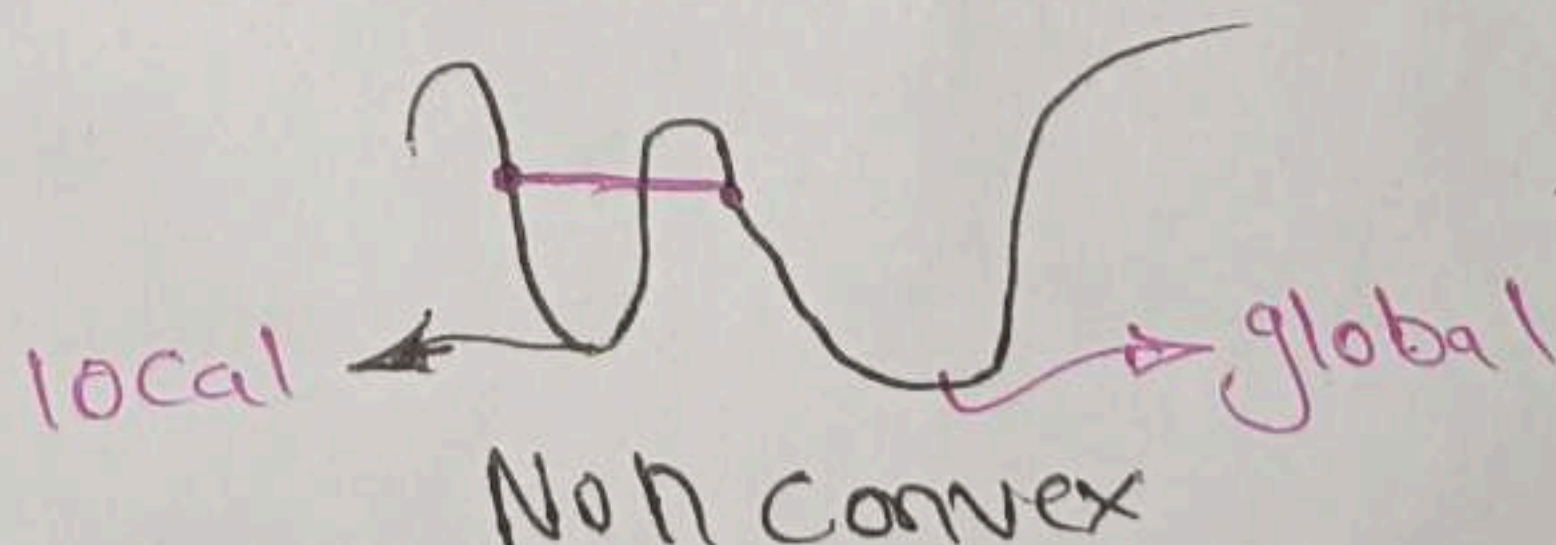
* Now what is **convex** function

① IF we draw a line between two points that line never cross the function that function is called **convex function**.

Ex:



Non convex Ex:-



① In convex function there is only one minima and that is **global minima**

① In non convex function there is more than one minima we called them **global minima** and **local minima**.

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 ~~the~~ ~~best~~ score 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 Function & Data.

① learning rate

② loss Function

③ Data.

① learning rate:

① low LR: ① IF LR is 0.002

① very small steps are getting

① more epoch value

① not much optimize

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.

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

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

→ this is nothing but our y_train

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

→ This is nothing but formula for calculating mean

understand calculation to finding y hat

x1	x2	x3	y
-1	-	-	-

→ for first row

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

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

$$\hat{y}_i = \beta_0 + [x_{i1} \ x_{i2} \ x_{i3}] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

We can calculate for all values in x_train by using vectorization

$$\hat{y} = \beta_0 + \text{np.dot}(\text{coef_}, \text{x_train})$$

Formula becomes

$$\beta_0 + \text{np.dot}(\text{x_train}, \text{coef_})$$

(353, 10) (10, 1)

x_train shape coef shape (353, 1)

* Now calculating of β_1 slope *

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

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

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

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

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

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

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

x_1	x_2	y
8.1	93	3.2
7.5	95	3.5

x_{i1} is represent all values of Column 1

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

* for m cols *

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

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

(MSE)

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

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

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

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

x_1	x_2	y
8.1	93	3.2
7.5	95	3.5

Annotations: x_{11} points to 8.1, x_{12} points to 93, y_1 points to 3.2, x_{21} points to 7.5, x_{22} points to 95, y_2 points to 3.5.

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

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

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

for N rows

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

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

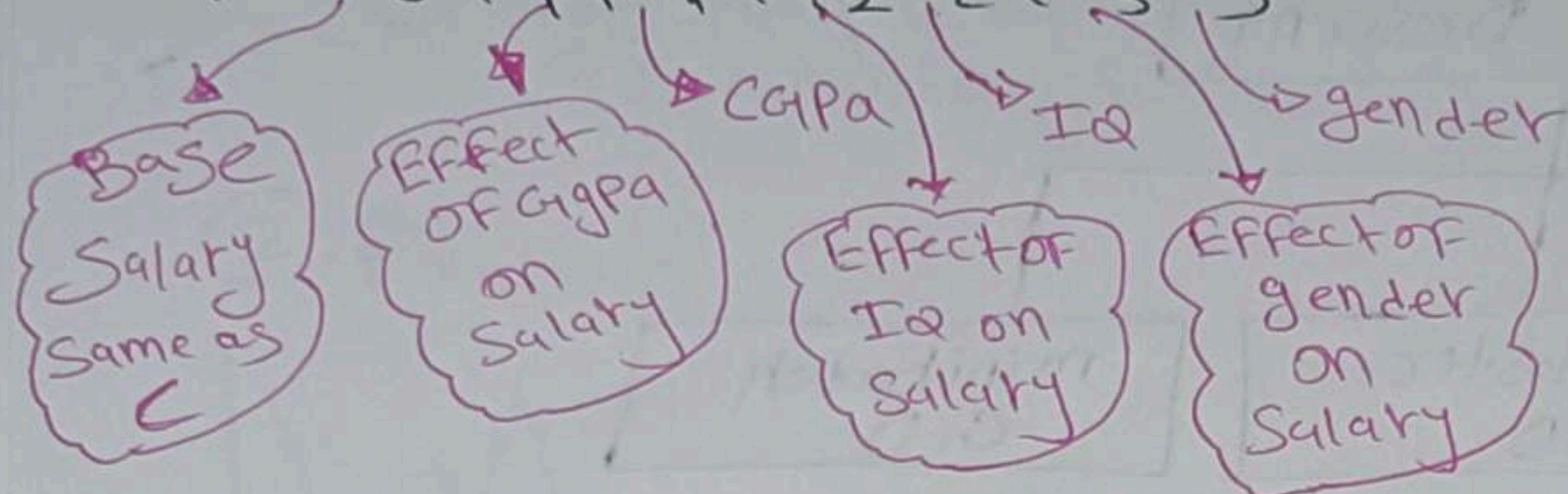
→ This is give the value of slope

$$\beta_0 = \beta_0 - n(\text{slope})$$

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

* Understand the equation *

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



* 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 | ipa

x_1	x_2	y
8.1	93	3.2
7.5	95	3.5

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

(IPA) (cgpa) (IPA)

Step 1: Start with random value

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

Step 2: Epoch = 100, lr = 0.1

$$\beta_0 = \beta_0 - \eta \text{slope} \rightarrow \frac{\partial L}{\partial \beta_0}$$

$$\beta_1 = \beta_1 - \eta \text{slope} \rightarrow \frac{\partial L}{\partial \beta_1}$$

$$\beta_2 = \beta_2 - \eta \text{slope} \rightarrow \frac{\partial L}{\partial \beta_2}$$

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

* 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 convergence
- row is select randomly
- Not give steady (ans) soln

* Now, coding Part *

```
from sklearn.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.
```

```
X, y = load_diabetes(return_X_y=True)
```

```
Print(X.shape)
```

```
Print(y.shape)
```

→ (442, 10)

→ (442,)

```
X, y = load_diabetes(return_X_y=True)
```

```
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_) # m → [-3.16 -205.46 516.68 840.62 -895.54
```

```
Print(reg.intercept_) # b → 157.88
```


* Stochastic Gradient Descent *

* Problem with Batch GD *

$$n = 1000$$

$$\text{cols} = 5 \rightarrow 6 \text{ Coefs}$$

$$\text{Epoch} = 50$$

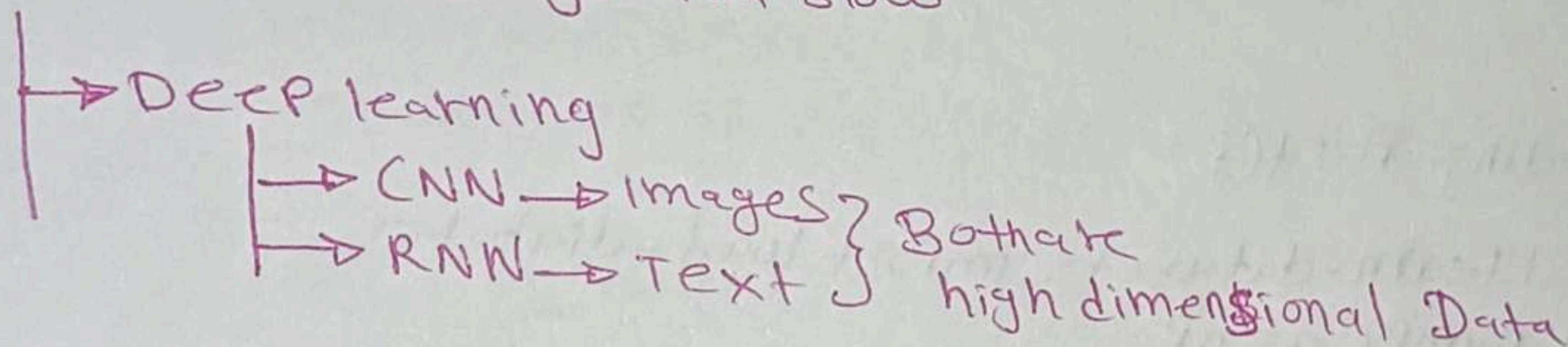
$$50 \rightarrow 1000$$

$$6 \rightarrow 6000$$

$$6000 \times 50 = 300000$$

derivatives
want to calculate

① It makes algorithm slow



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

In coding Part we did vectorization to calculate \hat{y}

$$\hat{y} = \text{np.dot}(x_{\text{train}}, \text{self.coef_}) + \text{self.intercept_}$$

In this we used numpy and avoid loop

* When we are calculating \hat{y} 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

$$\text{Coef_der} = -2 * \text{np.dot}((y_train - y_hat), x_train) / x_train.shape[0]$$

$$\text{self.coef_} = \text{self.coef_} - (\text{self.lr} * \text{coef_der})$$

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

$$\hat{y} = \text{np.dot}(X_{\text{train}}, \text{self.coef_}) + \text{self.intercept_}$$

$$\text{intercept_der} = -2 * \text{np.mean}(y_{\text{train}} - \hat{y})$$

$$\text{self.intercept_} = \text{self.intercept_} - (\text{self.lr} * \text{intercept_der})$$

Now time to calculate coef_

For single:

$$\frac{\partial L}{\partial \beta_1}, \frac{\partial L}{\partial \beta_2}, \frac{\partial L}{\partial \beta_n}$$

~~Simple~~
 (*) 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	x_2	y	\hat{y}
1	2	5	6
3	4	7	8

$$y = [5 \ 7] \quad \hat{y} = [6 \ 8]$$

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

$$\begin{bmatrix} -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$[-4 + -3]$$

$$\textcircled{-8}$$

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

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

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

$$\begin{matrix} \nearrow & \searrow \\ (353, 1)^T & (353, 10) \end{matrix}$$

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

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

$y_pred = \text{reg.predict}(x_test)$

$r2_score(y_test, y_pred)$

→ 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

self.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

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

1

collectively

Not required in SGD

for SGD

**

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

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

$$\frac{\partial L}{\partial \beta_1} = -2 (y_i - \hat{y}_i) x_{i1}$$

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 = SGDRegressor(learning_rate=0.01, epochs=50)

Sgd.fit(X_train, y_train)

→ 150.8918 [54.28 -67.43 851.20 251.58 17.26
-28.74 -165.85 123.27 314.84 123.52]

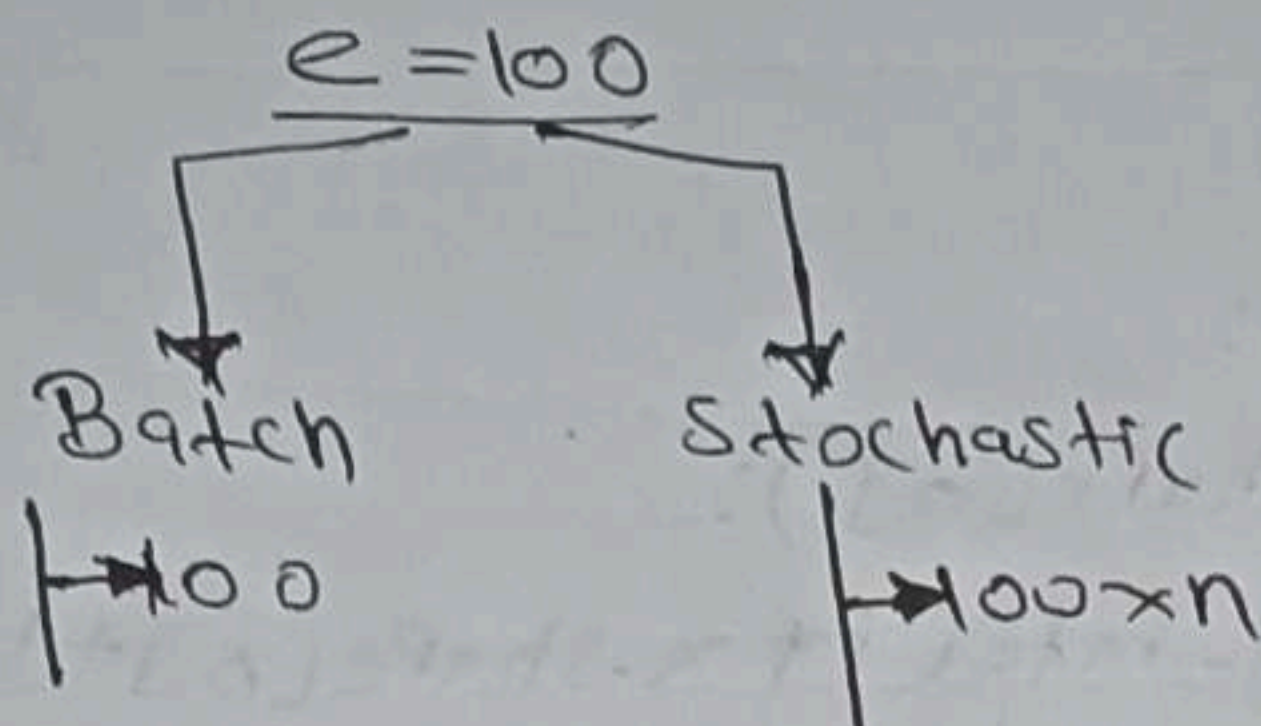
y_pred = Sgd.predict(X_test)

r2_score(y_test, y_pred)

→ 0.4325

* Time comparison *

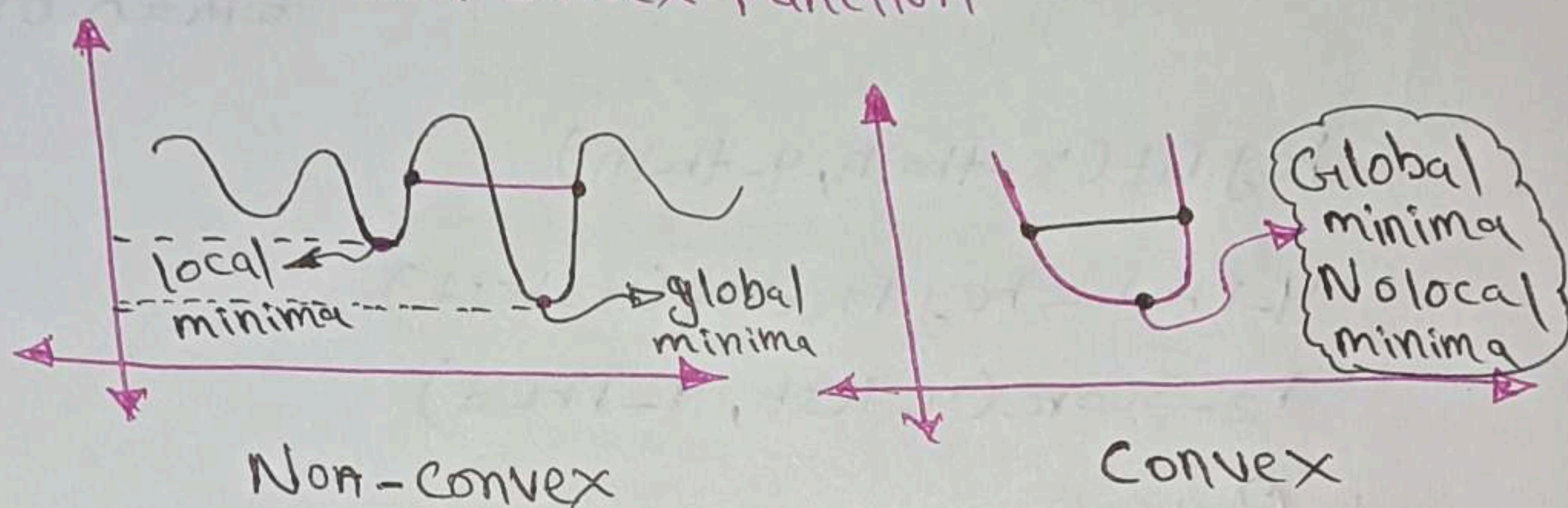
no. of epoch \rightarrow fixed



* When to use Stochastic GD

① Big data

② When we have non convex function



* Because of stochastic nature of SGD it could be ~~that~~ get out from local minima (It have ^{that} velocity which is not stop it to local minima.)

* Problem with SGD

① 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.

$y_hat = np.dot(x_train[idx], self.coef_) + self.intercept_$

$intercept_der = -2 * np.mean(y_train[idx] - y_hat)$

$coef_der = -2 * np.dot((y_train[idx] - y_hat), x_train[idx])$

$self.coef_ = self.coef_ - (self.lr * coef_der)$

$Print(self.intercept_, self.coef_)$

$def Predict(self, x_test)$

$return np.dot(x_test, self.coef_) + self.intercept_$

$mbt = MBGDRegressor(batch_size = int(x_train.shape[0]/10),$
 $learning_rate = 0.01,$
 $epochs = 500)$

$mbt.fit(x_train, y_train)$

$\rightarrow 154.8374 \quad [38.40 \quad -142.67 \quad 457.28 \quad 303.60 \quad -17.99$
 $\quad -85.81 \quad -192.04 \quad 116.18 \quad 407.24 \quad 105.80]$

$y_pred = mbt.predict(x_test)$

$r2_score(y_test, y_pred)$

$\rightarrow 0.45188$

(*) How to use MBGDR in sklearn (*)

$from sklearn.linear_model import SGDRegressor$

$sgd = SGDRegressor(learning_rate = 'constant', eta0 = 0.2)$

$\# Partial_Fit(x, y, sample_weight = None)$

\hookrightarrow Using this we can pass subset of $x_train,$
 y_train


```
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  
    153.88 126.73 861.2 52.41 151.88]  
    151.88.
```

```
y_Pred = reg.predict(X_test)
```

```
r2_score = (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
```

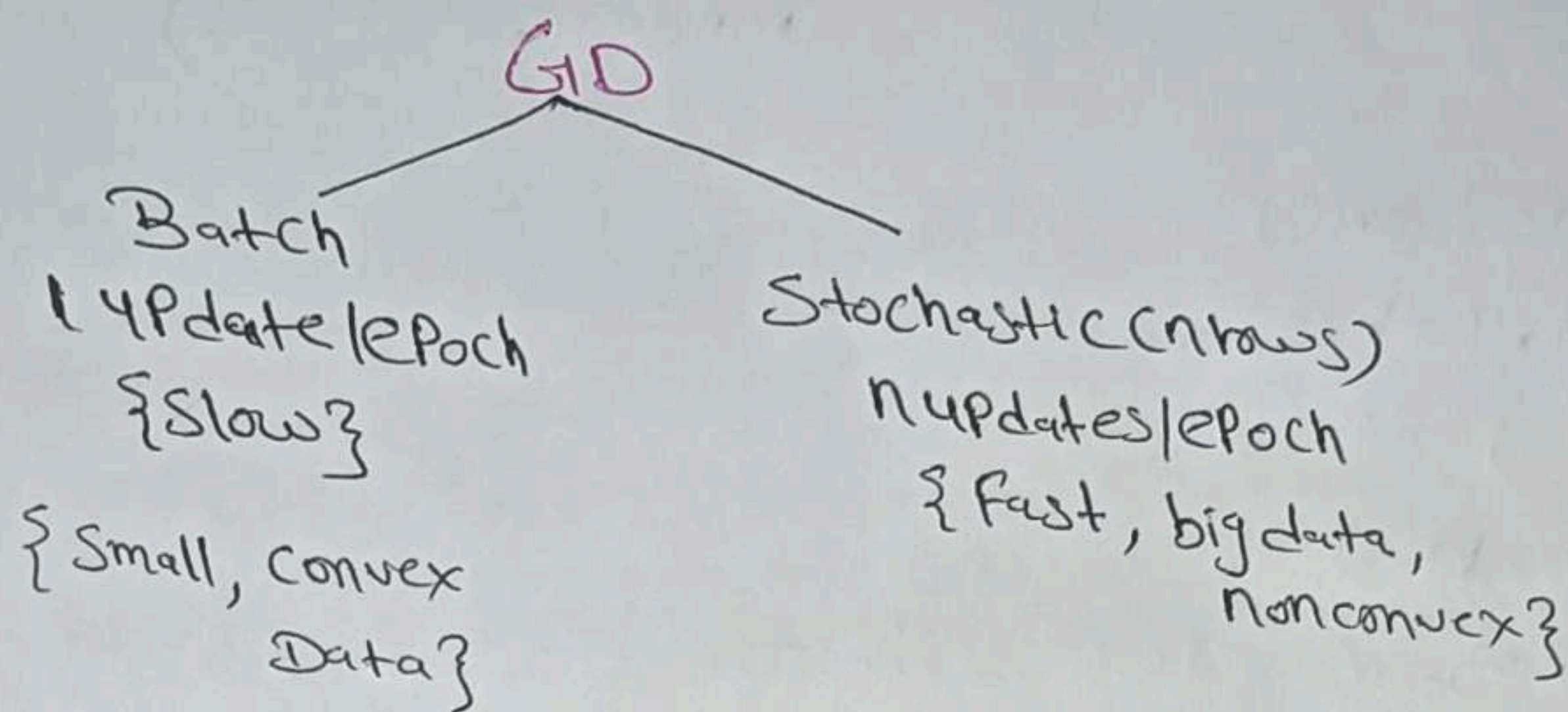
```
        self.coef_ = np.ones(X_train.shape[1])
```

```
        for i in range(self.epochs):
```

```
            for j in range(int(X_train.shape[0]/  
                               self.batch_size)):
```

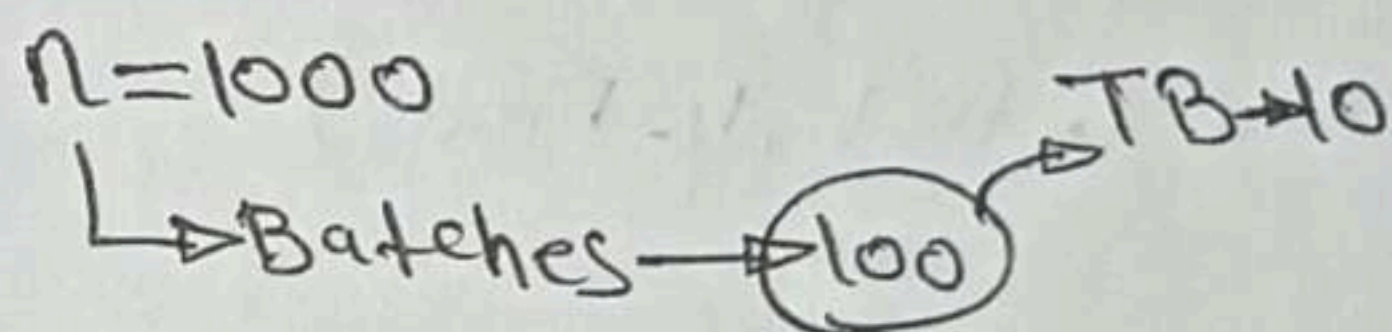
```
                idx = random.sample(range(X_train.shape[0]),  
                                     self.batch_size)
```


* Mini-Batch Gradient Descent *

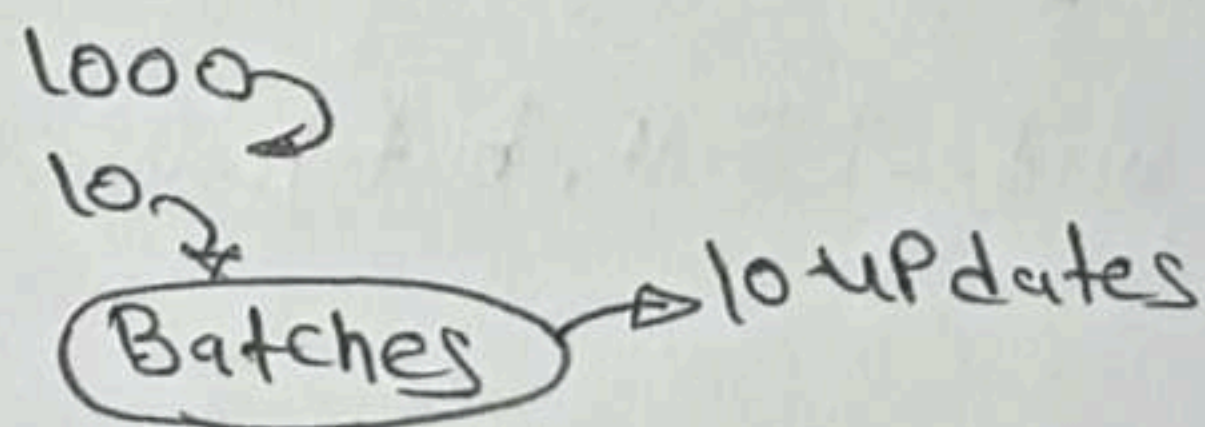


minibatch GD: Combination of Batch & Stochastic GD

Batch \rightarrow Group of rows



Batches updates/epochs



* Coding Part *

```
from sklearn.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)
```

\rightarrow 442, 10
(442,)

*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 SGDRegressor

```
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

batch-size = 35

for i in range(100)

idx = random.sample(range(x_train.shape[0]), batch_size)

Sgd.Params_.fit(x_train[idx], y_train[idx])

Sgd.Coeff_

• $\rightarrow [49.19, -67.84, 338.57, 247.97, 25.30, -24.71,$
 $-155.45, 116.19, 312.91, 133.36]$

Sgd.intercept_

$\rightarrow 148.61$

y_pred = Sgd.predict(x_test)

r2_score(y_test, y_pred)

$\rightarrow 0.4271$