

LINEAR REGRESSION

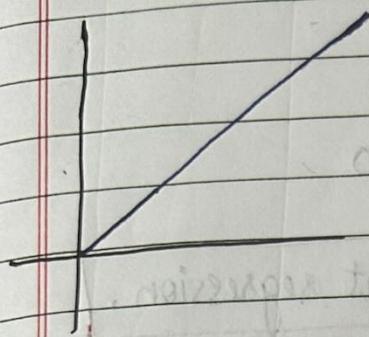
→ linear Regression Equation-

$$y = mx + c$$

standard equation

or

straight line equation



$$\frac{y_2 - y_1}{x_2 - x_1} \quad ? \text{ slope}$$

* I love Python → encoding

Machine learning sochta hai

$0 < 1 < 2 \rightarrow$ gives priority

0 | 1 | [010] ↑

⇒ If distributions are skewed → Non values with median
normal → Non values with mean

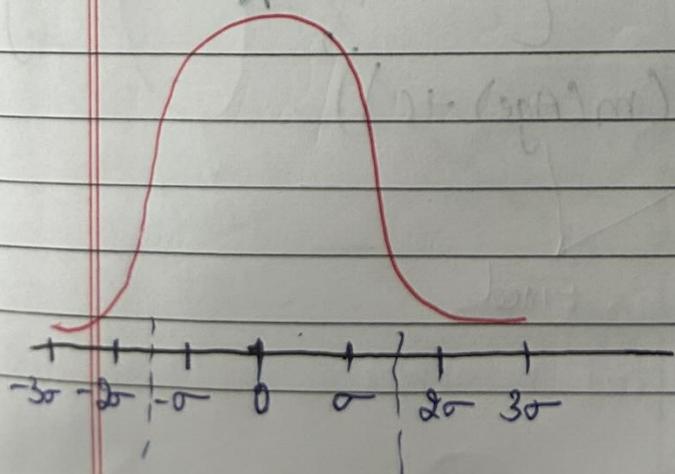
⇒ Scaling should be same for all columns.

⇒ Has column ko 0 & 1 ke beech mein like aata hai

→ off - 6 & 10 → 68%.

2 deviation & 21 R

→ 97%.



9/1/23

Date / /
Page No.

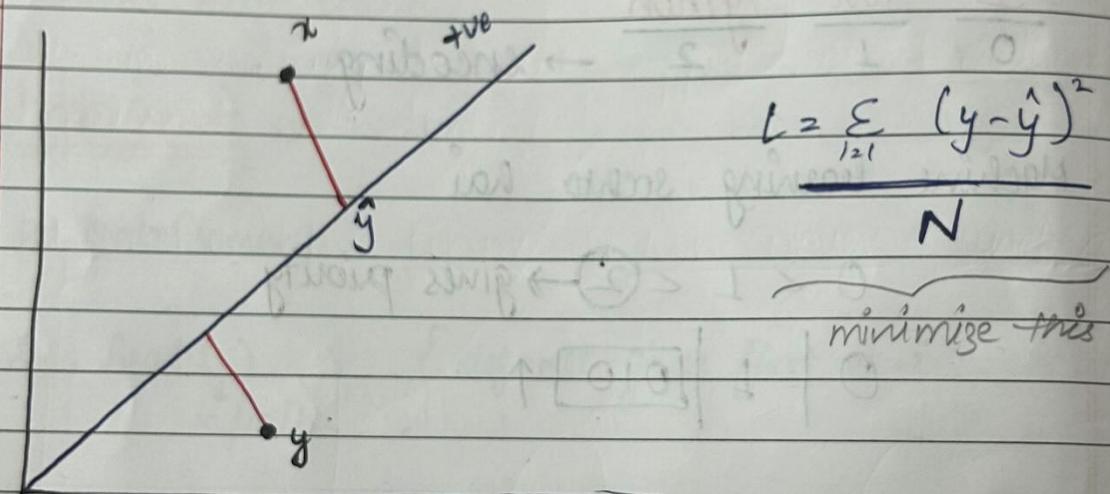
- * How much far the predicted value is far from the actual value → is loss

$$\text{Single Loss} = \sum_{i=1}^n (y - \hat{y})^2$$

choose line jiska loss minimum ho

- * Linear regression also called univariate regression.

10/1/23



Eg)	Age	Salary	Diff.
	21	89	
	22	88	
	23	70	

Parameters

$$\left\{ L = \frac{\sum (y - (m\alpha + c))^2}{N} \right\}$$

$$L = \frac{\sum_{i=1}^n (\text{Salary} - (m(\text{Age}) + c))^2}{N}$$

Fixed

Objective = Find m, c

Eg) Dataset (Plotted)

The gap is called residual error

choose the line having min error

$$(m, c = 0)$$

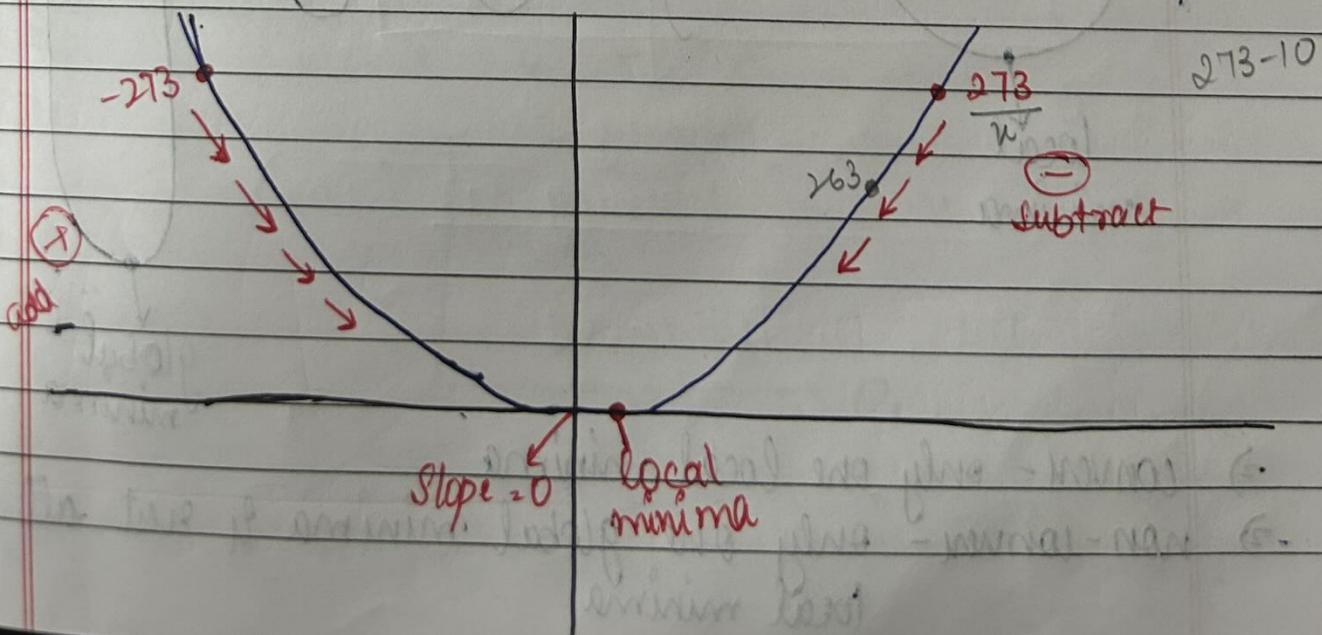
↓
machine
will start from
here & gradually
move up

→ L = Distance of all points from line

$$L = \sum_{i=1}^n (y - (mx + c))$$

minimise

Eg) $y = 4x^2$ { set slope = 10 }



Page No. _____

Technique which tells which direction we have to go, is GRADIENT DESCENT

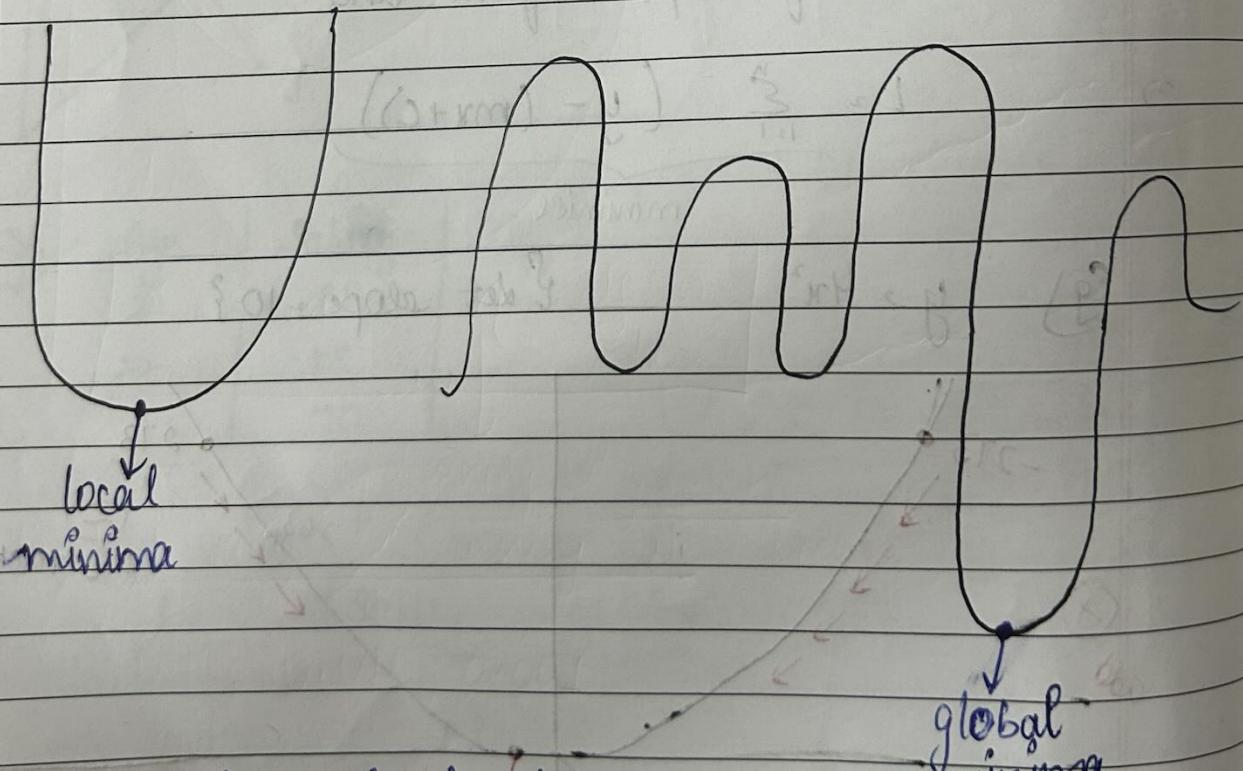
GRADIENT, DESCENT,
slope decrease

→ Two ways:

- ① Normal Method (Direct differentiation)
- ② Gradient Descent (Iterative method)

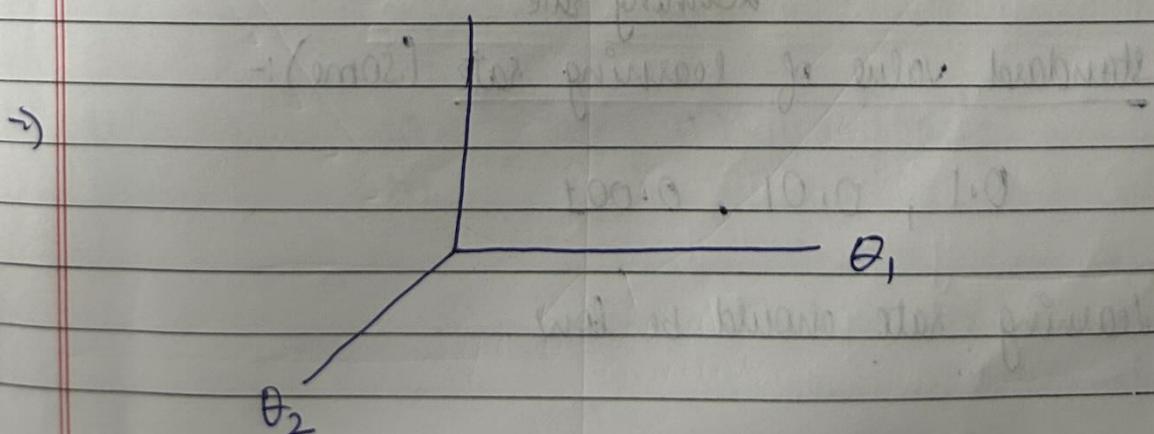
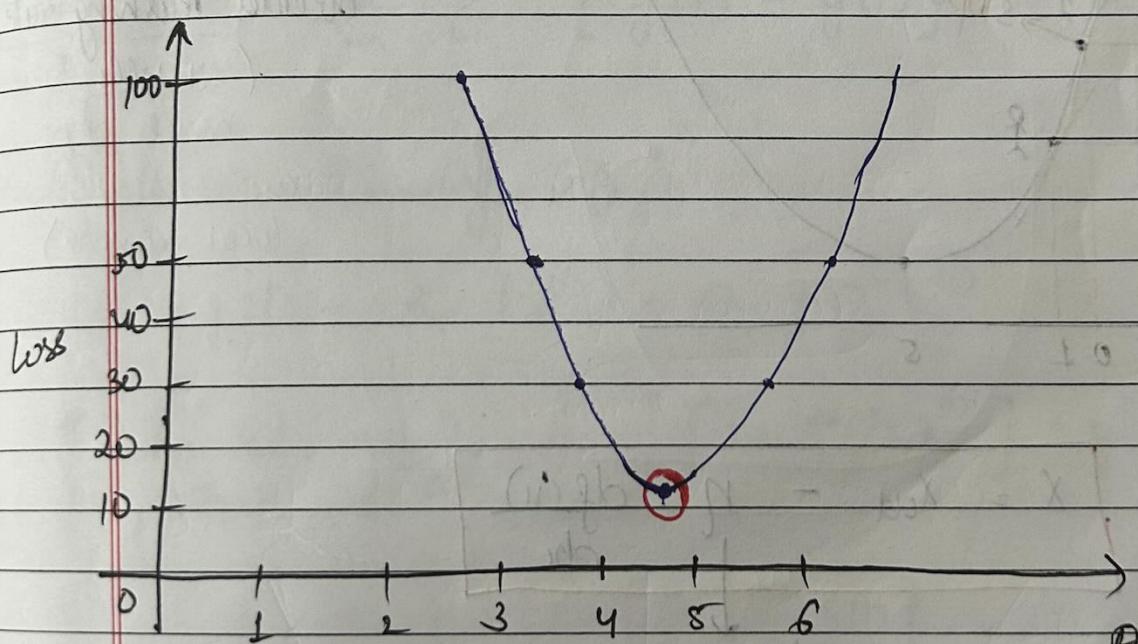
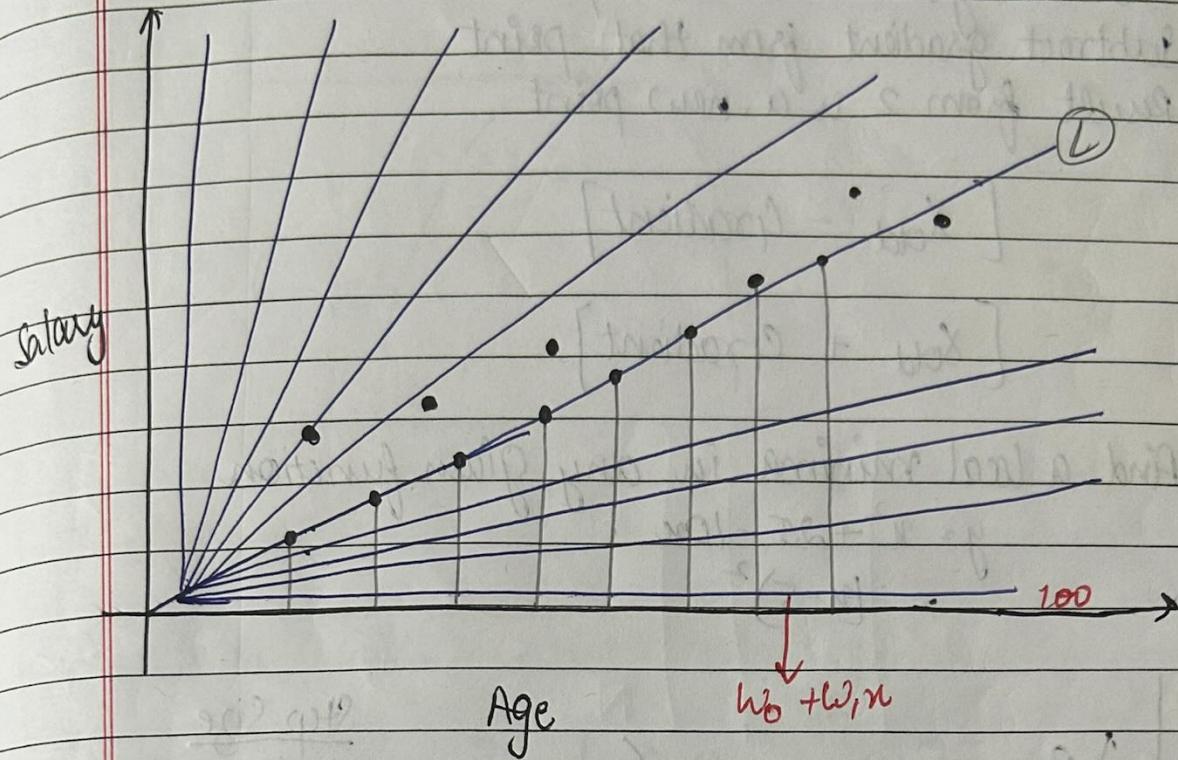
→ Machine learning does not support normal method, because it costs time.

→ Difference b/w Convex & Non-convex function



- convex - only one local minima
- non-convex - only one global minima & rest will all be local minima

2) Curvature of Loss Function



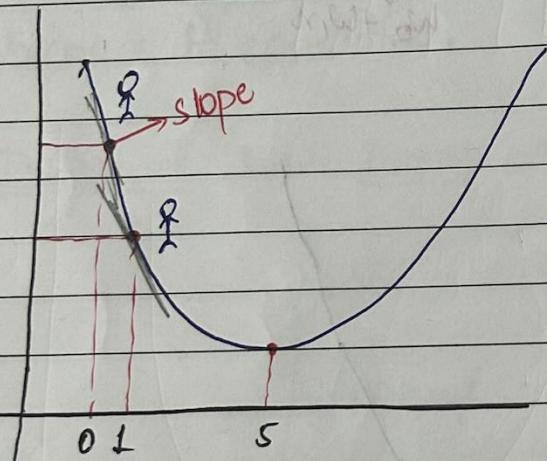
- Steps :-
- i) Find the gradient at that point
 - ii) Subtract gradient from that point
 - iii) Result from 2 is a new point

$$[x_{\text{old}} - \text{Gradient}]$$

$$[x_{\text{old}} + \text{Gradient}]$$

Ques) Find a local minima in any given function

$$\begin{aligned}y &= u^2 + 25 - 10u \\&= (u-5)^2\end{aligned}$$



Step Size
decides learning rate
(It tells how much steps are required to reach local minima)

⇒
$$X = x_{\text{old}} - \eta \cdot \frac{df(u)}{du}$$

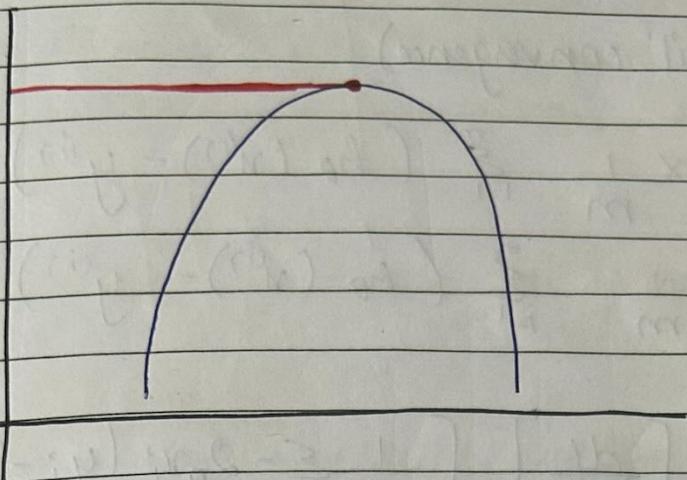
Learning rate

2) Standard value of learning rate (some) :-

0.1, 0.01, 0.001

2) Learning rate should be low

* for maxima, Gradient Ascent



→ Steps:

$$① \theta = .init()$$

$$② J(\theta) = \frac{1}{m} \sum_{i=1}^m [\hat{y}(i) - y(i)]^2$$

③ Update θ using GD

$$\hat{y}(i) = h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

$$*\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{m} \sum_{i=1}^m [(\theta_0 + \theta_1 x^{(i)}) - y^{(i)}]$$

$$= \frac{2}{m} \sum_{i=1}^m [\hat{y}(i) - y(i)]$$

$$(2) \frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{m} \sum_{i=1}^m [\hat{y}(i) - y(i)] x^{(i)}$$

Gradient Descent Algorithm

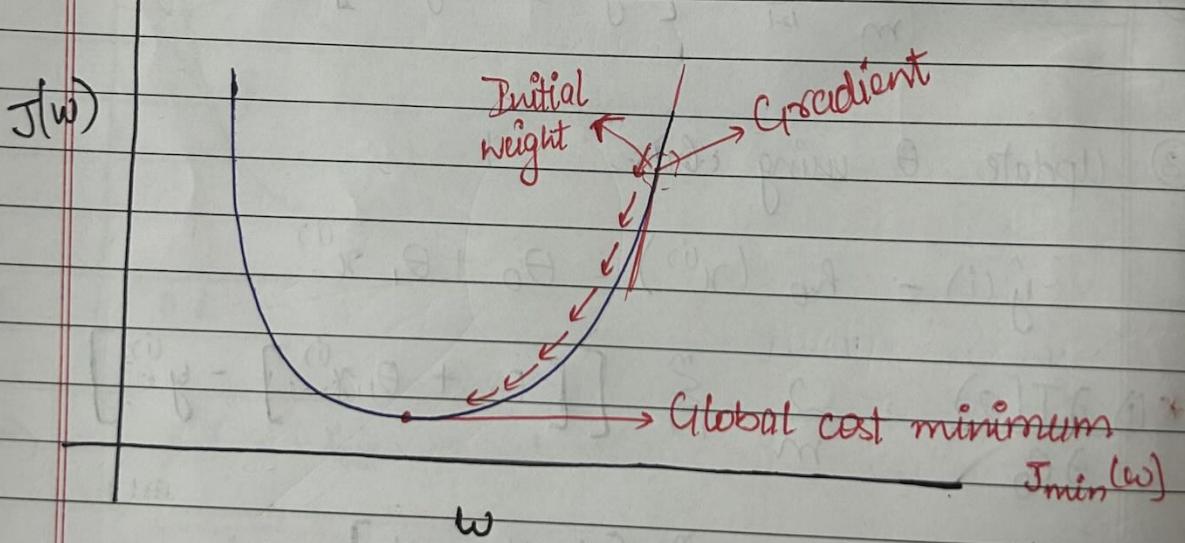
(repeat until convergence)

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

12/1/23

* $J'(m, b) = \begin{bmatrix} \frac{df}{dw} \\ \frac{df}{db} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N 2x_i (y_i - (wx_i + b)) \\ \frac{1}{N} \sum_{i=1}^N 2(y_i - (wx_i + b)) \end{bmatrix}$



Assumptions of Linear Regression

- i) Linear relationship
- ii) No multicollinearity
- iii) Normality of residual
- iv) Homoscedasticity
- v) No auto correlation

13/11/23

Date / /
Page No.

→ Linear relationship / Linearity

There should be linear relationship b/w independent & dependent variables, bcoz it suggests that a change in y due to change in x is constant, regardless the value of x .

{ Non-linear \rightarrow log, exponential to convert into linear }

→ No multicollinearity

Independent variables should not be correlated, absence of this phenomenon is known as multicollinearity.

Multi collinearity \rightarrow Tolerance $(1 - R^2)$ [$T < 0.1$, coefficient of determination exists]

↓
VIF $\left(\frac{1}{1 - R^2} \right)$ [VIF > 10 , exists]

{ variance inflation factor }

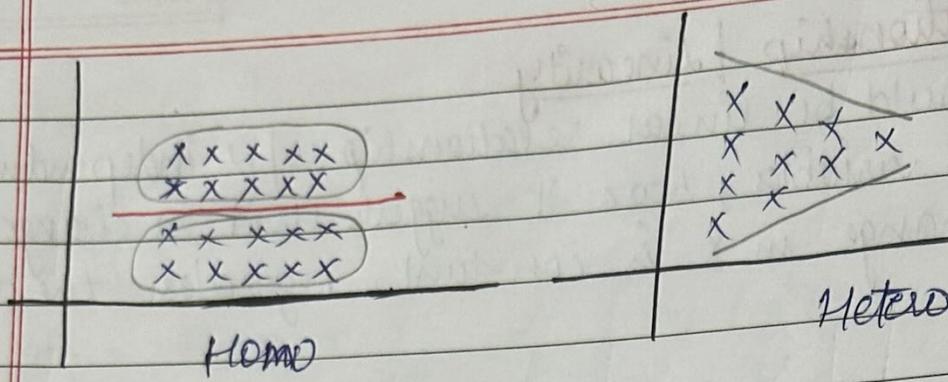
→ Normality of residual / Normal distribution of error Error terms must be normally distributed.

Normal distribution of error \rightarrow Analytical
(small dataset)
(large not well)
Graphical

Quantile Plot (to check normal distribution)

→ Homoskedasticity

Constant variance of the residuals



→ No autocorrelation

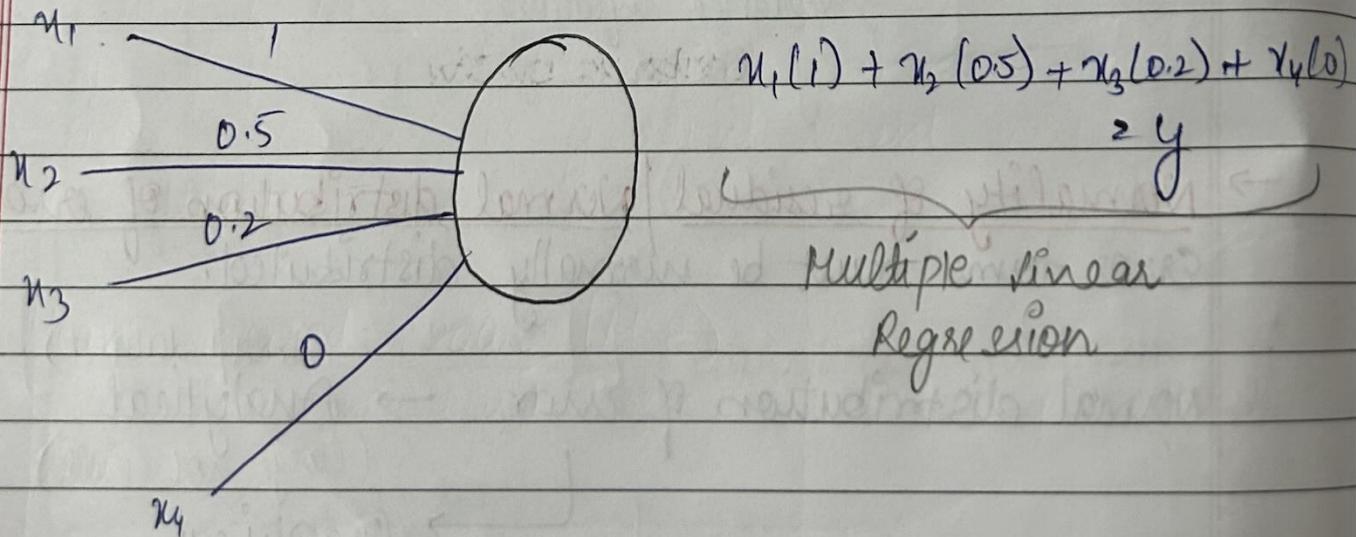
There should be no correlation b/w the residual (error) terms.

16/11/23

Ques) How to check performance of the model?

→ Evaluation Metrics

* Weights - gives importance to inputs.
(for feature selection)



The incoming weights to the neuron is a measure to what the neuron cares about.

Types of Gradient Descent

Gradient Descent

Batch Gradient Descent

- computes on whole training data
- slow & computationally expensive
- not for huge training dataset

Stochastic Gradient Descent

(stochastic - random)

computes on single training sample

fast & less computationally expensive

can be used for huge training data.

* GD

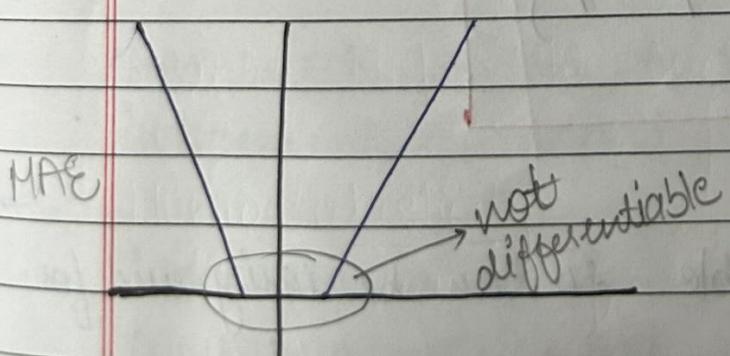
Batch Stochastic Mini-Batch
(does in batches)

* Faster - mini-batch gradient descent

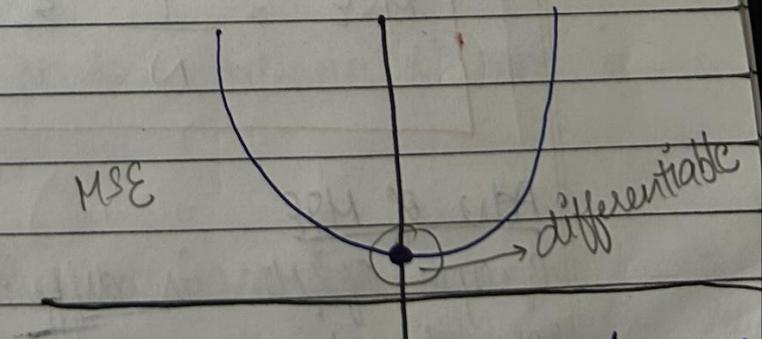
Evaluation Metrics

Mean Absolute Error (MAE)

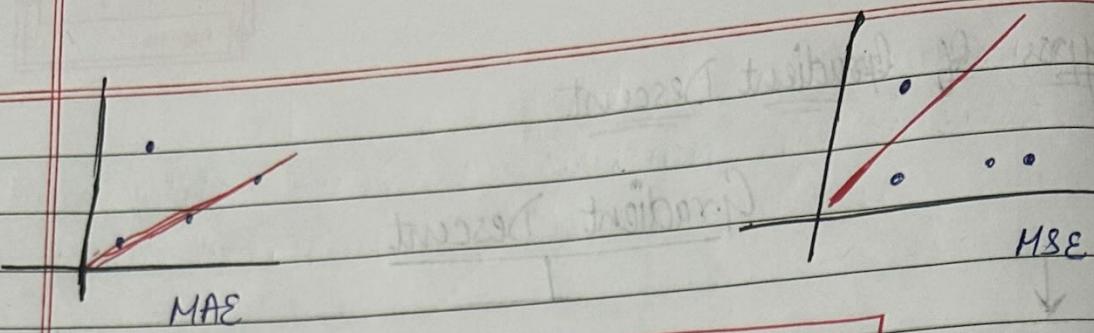
MAE calculates diff b/w actual & predicted values.



Modulus
Outlier won't affect or have an impact on the line



line will be distorted bcoz of outliers



$$\boxed{\text{MAE} = \frac{1}{N} \sum |y - \hat{y}|}$$

Divide by total no. of data points sum of Absolute value of residual
 Actual output Predicted output

Adv of MAE

- gives same unit as output variable
- It is most robust to outliers

Disadv of MAE

Graph of MAE is not differentiable, so we have to apply various optimisers like gradient descent which can be differentiable.

Mean Squared Error (MSE)

It represents the squared distance b/w actual & predicted values.

$$\boxed{\text{MSE} = \frac{1}{N} \sum (y - \hat{y})^2}$$

Adv of MSE

Graph of MSE is differentiable, so can be easily used for loss function.

Disadv of MSE

- output unit is squared unit. (Eg $m \rightarrow m^2$)
- it penalises the outliers most & calculated MSE is bigger,
∴ it is not robust to outliers.

→ Root Mean Squared Error (RMSE)

It is square root of mean squared error

$$RMSE = \sqrt{MSE}$$

$$RMSE = \sqrt{\frac{1}{N} \sum (\hat{Y} - Y)^2}$$

Adv of RMSE

same unit as output variable

Disadv of RMSE

not robust to outliers

* mostly preferred in deep learning

→ R Squared (R²) | Coefficient of Determination

It tells performance of the model, not the loss in absolute sense

MAE & MSE depends on the content, whereas R² score is independent of content.

With the help of R-squared we have baseline model to compare a model which none of the other metrics provides. The same we have in classification problems which we call

a threshold which is fixed at 0.5. So, basically R^2 calculates how much regression line is better than a mean line.

→ also known as Goodness of fit

$$R^2 = \frac{RSS}{TSS}$$

RSS → sum of square of residuals

TSS → total sum of squares

R^2 → coefficient of determination

→ $R = 1$ (Perfect line)

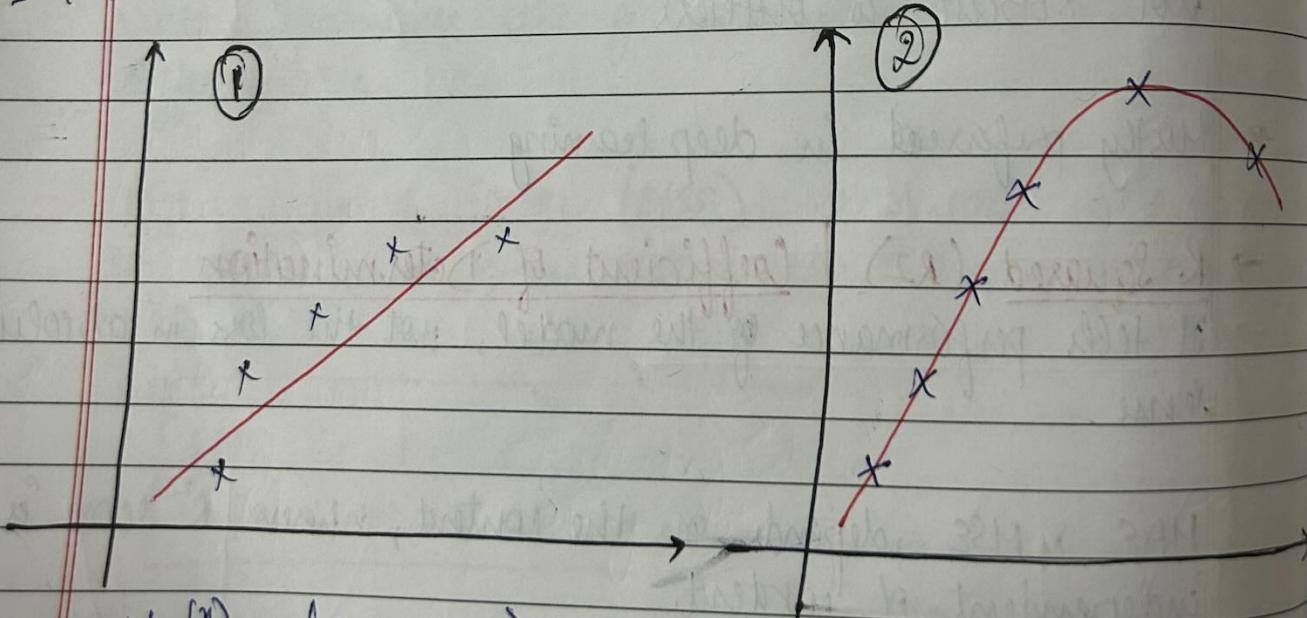
$R = 0$ (Worst line)

$R < 0$

$R > 0$

{ R should be near }
 to $\frac{1}{R} \rightarrow R \rightarrow 1$

~~17/1/23~~

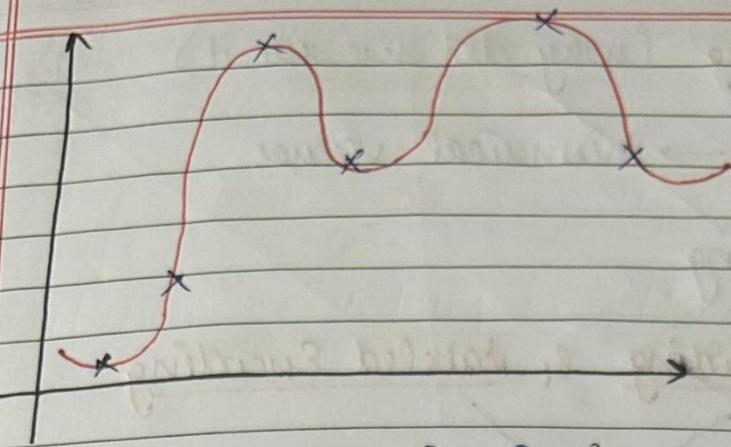


$$M\theta^{(n)} = (\theta_0 + \theta_1 u)$$

H^+ hypothesis is giving me straight line which is underfitting

$$M\theta^{(n)} = \theta_0 + \theta_1 u + \theta_2 u^2$$

Here we are increasing the order of polynomial

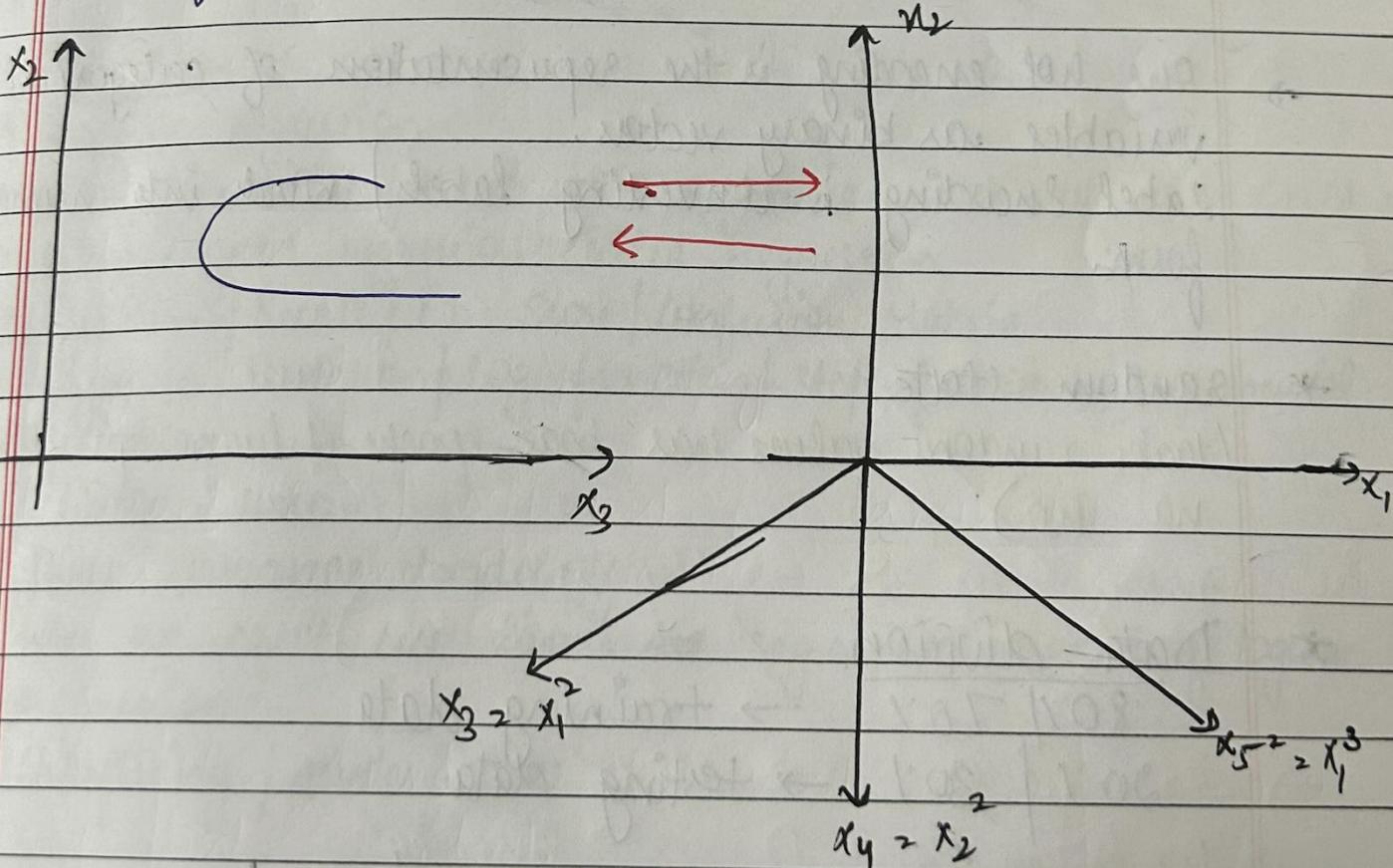


$$M_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \theta_5 x^5 + \theta_6 x^6$$

Order of Poly = 6 (Perfect)

→ Original data (Quadratic) = $(u, v_i) \in R^2$

→ Transformed data (linear) = $f(u) = (x_0, u^2, u^3, \dots, u^n)$



* numpy > dataframe
(faster)

* Weights tells me importance of features
- If tells us that feature is not imp.

* One hot Encoding (why did we use it)

\downarrow
Categorical values \rightarrow numerical values
 \uparrow
Labelled Encoding

\Rightarrow ~~Both~~ One hot Encoding & Labeled Encoding

\swarrow \searrow
Categorical to vectors

- \Rightarrow Jahan pe order hai / meaning hai / comparison hai / much hai / greater - less hai, wahan pe one hot encoding use hoti hai
- \Rightarrow One hot encoding is the representation of categorical variables as binary vector.
Label encoding is converting labels / words into numeric form.

* random state

(Toaki random values har baar epochs chlaane pr change na ho)

* Data division

80% / 70% \rightarrow training data
20% / 30% \rightarrow testing data

- \Rightarrow Due to one hot encoding, parameters / inputs gets increased in number.
 \therefore we get multiple m values.