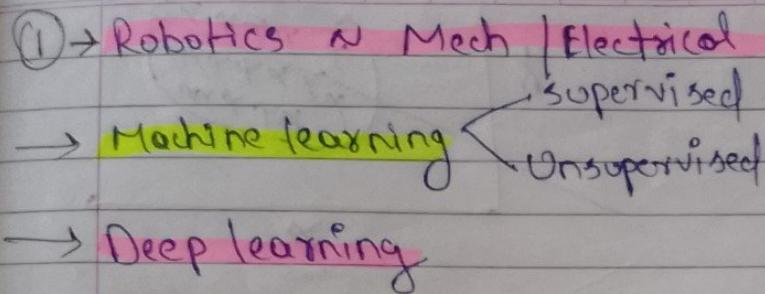


Artificial Intelligence



IOT

Alexa

smart work
smart watches
smart cards

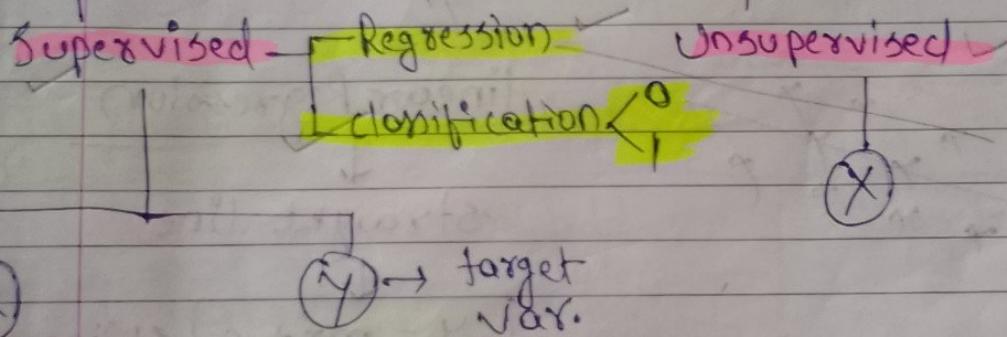
Reinforcement learning

Real Time Processing → Kafka

Machine learning

Computer Vision

→ face detection
→ image
→ pixel



linear

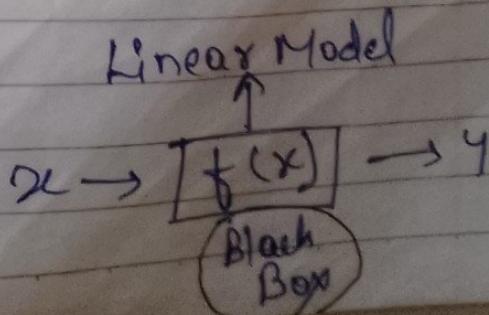
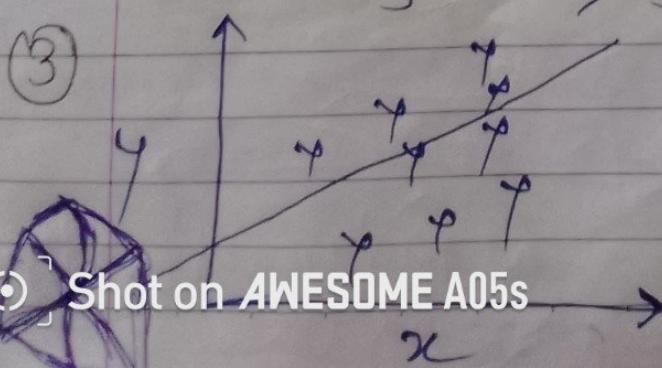
② Machine Learning

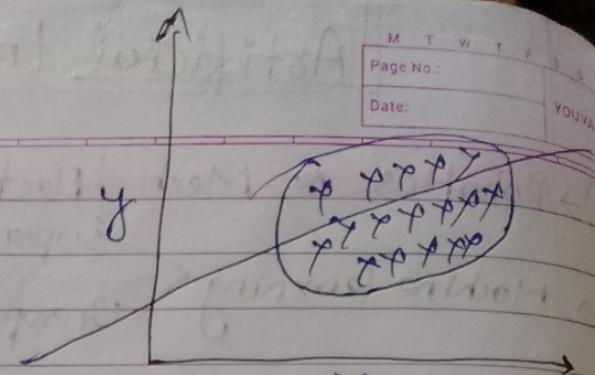
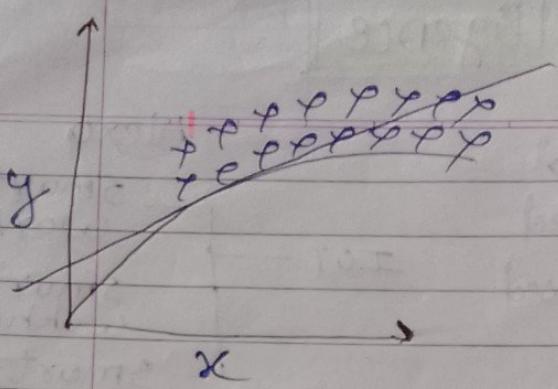
Predictor Independent x^2

Model $f(x)$

Target Dependent y

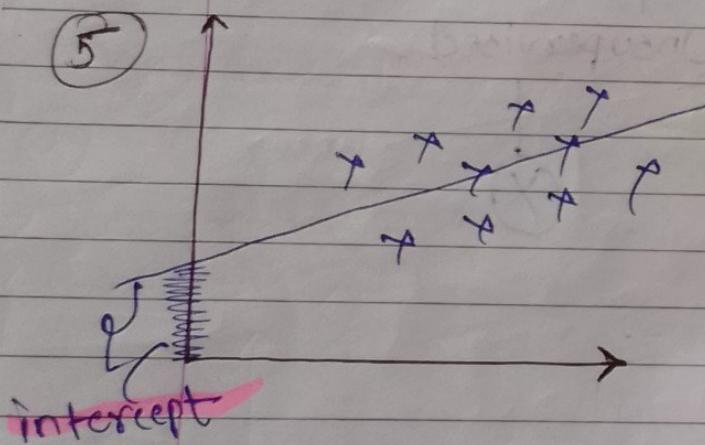
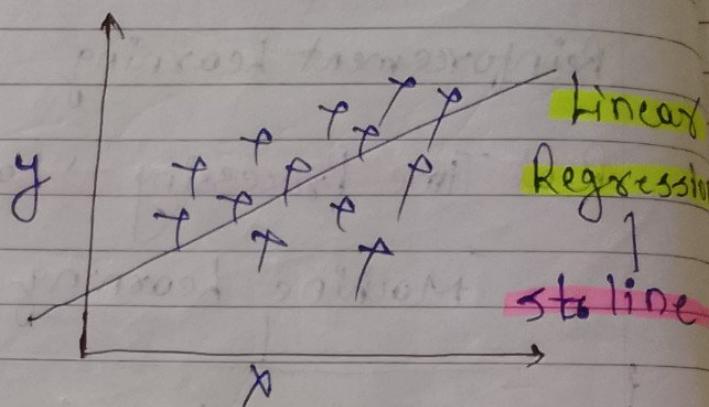
$$\begin{aligned} 1 &\rightarrow x^2 \rightarrow 1 \\ 2 &\rightarrow 2^2 \rightarrow 4 \\ 5 &\rightarrow 5^2 \rightarrow 25 \end{aligned}$$





(H)

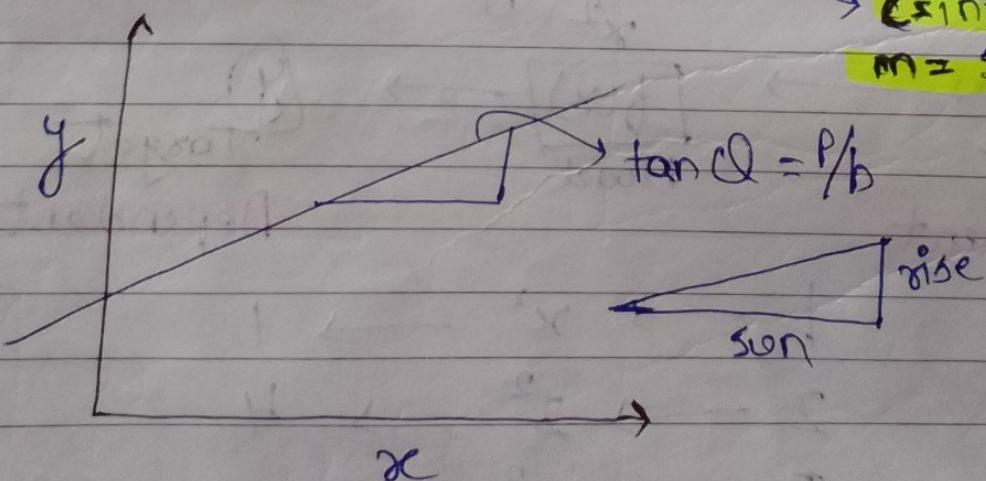
X	Y
1	2
2	4
3	6
4	8
5	10
6	12
7	14
8	16
9	18
10	20



linear Regression
↓
straight Line
↓

$$y = mx + c \quad \text{st.line. eqn}$$

$c = \text{intercept}$
 $m = \text{slope}$



(6)

Deterministic Model

Machine Learning

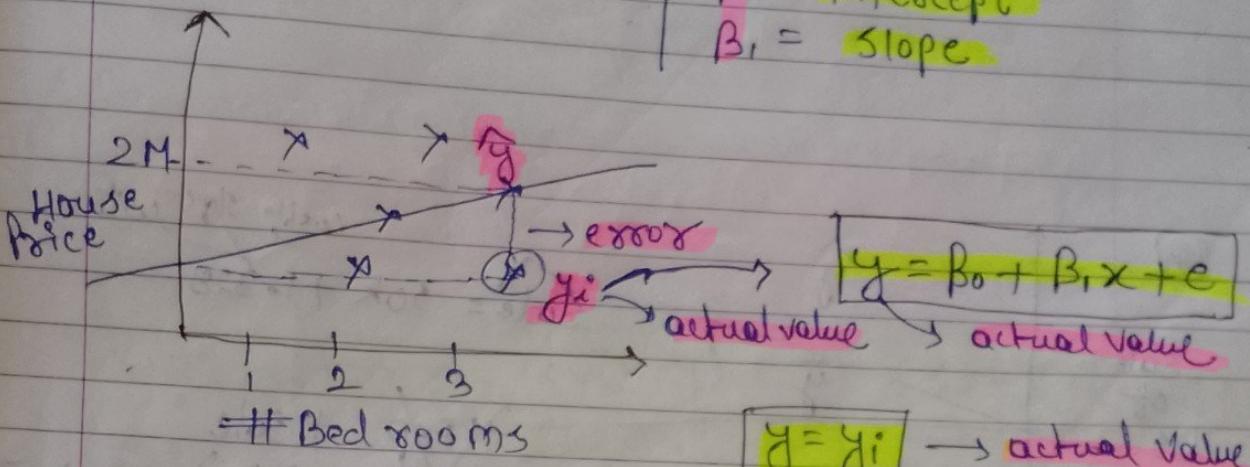
$$\hat{y} = \beta_0 + \beta_1 x$$

$$y = mx + c$$

\hat{y} = predicted value of y

β_0 = intercept

β_1 = slope



$$y = mx + c$$

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\text{error} = y_i - \hat{y}$$

Coordinate geometry equations.

machine learning equations.

(7)

YouTube spent
Lakhs (x_1)

35

40

45

TV add
spent
Lakhs (x_2)

15

20

21

Sales in Crores
 y

142

156

150

$\rightarrow \beta_1 x_1$

$\rightarrow \beta_2 x_2$

pred(y)

$$\text{sales} = 110 + 0.98yt + 0.32TV$$

↓
intercept

Multiple Linear Regression = $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

Sales

$$= 110 + 0.98yt + 0.32TV$$

$$yt = 0.98 \times 100000 = 98000$$

TV = 0.32 coeff

sns. lmplot

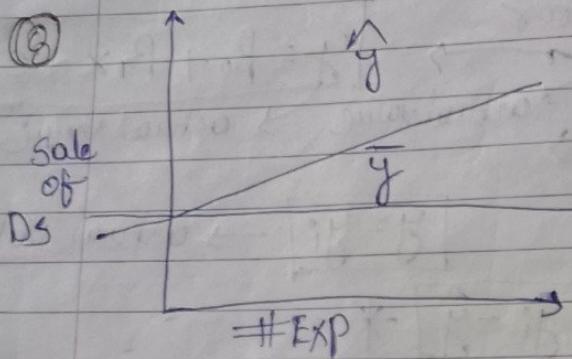
(8) Shot on AWESOME A05s

NULL MODEL

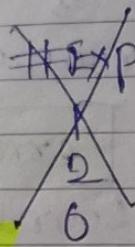
$$\hat{y} = \beta_0 = \bar{y}$$

Distance Matrix

* If my slope of the equation becomes zero then average sale in this problem $B_0 = 170$



$$\text{Sale} = 50k + 1.328$$



#Sale

50k

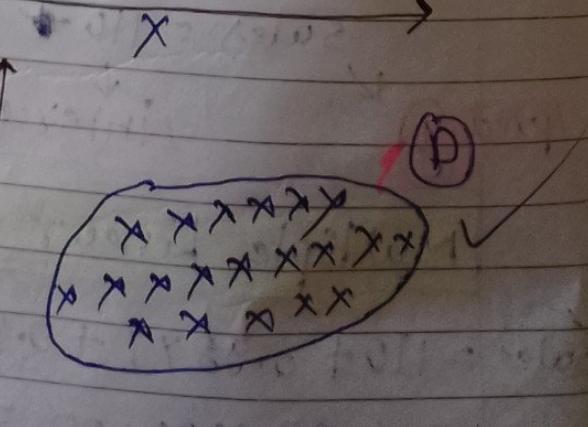
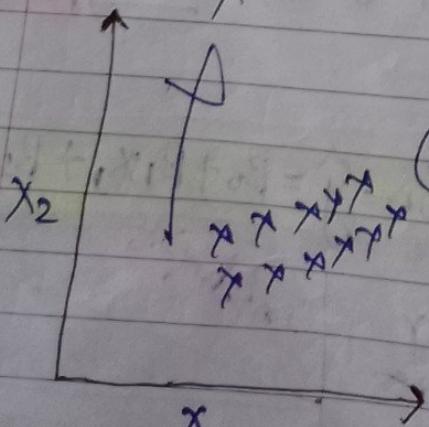
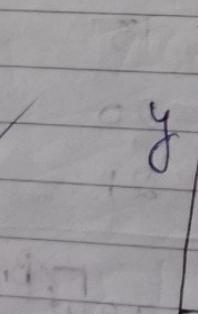
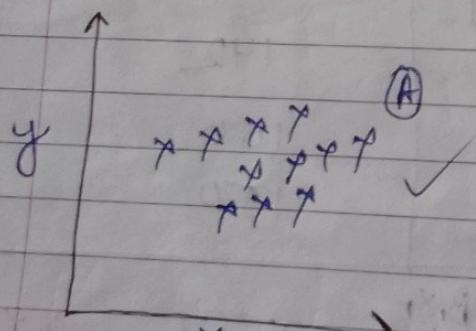
100k

90k

interval

* Base Model is always Model → is always worst model.

- Two predictors should not have a strong relation.
- Predictors should have strong relationship with target var



* Weak Predictors are strongly related to Multicollinearity
Shot on AWESOME A05s

(10)

$$\hat{y} = f(\text{youtube}, \text{TV add})$$

$$\hat{y} = f(\text{youtube}, \text{TV add})$$

M	T	W	T	F	S	S
Page No.:	6					
Date:	YOUVA					

$$\text{youtube} \quad \text{TV add} = 85\%$$

$$\text{Youtube v/s Sales} \approx 91\%$$

$$\text{TV add v/s Sales} \approx 89\%$$

(11)

$$\text{Covariance} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\text{Correlation} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

slope (B_1)

$$B_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Intercept} = \bar{y} - B_1 \bar{x}$$

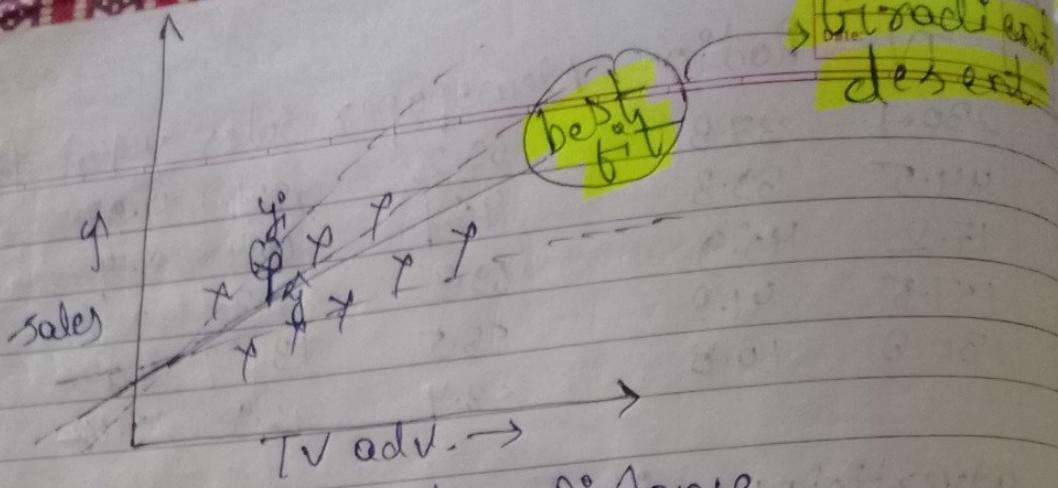
```
advertising = pd.read_csv('n/Downloads/Advertising.csv')
advertising.head()
```

	TV	Radio	Newspaper	Sales
0	230.1	37.8	65.2	22.0
1	44.5	39.3	45.1	10.4
2	17.2	45.9	69.3	9.3
3	151.5	41.3	58.5	18.5
4	180.8	10.8	58.4	12.9

About Linear Regression

- If I use only one Predictor ** and Target to build my linear Regression Model, then such a model is known as Simple Linear Regression
- Shot on AWESOME A05s

(13)



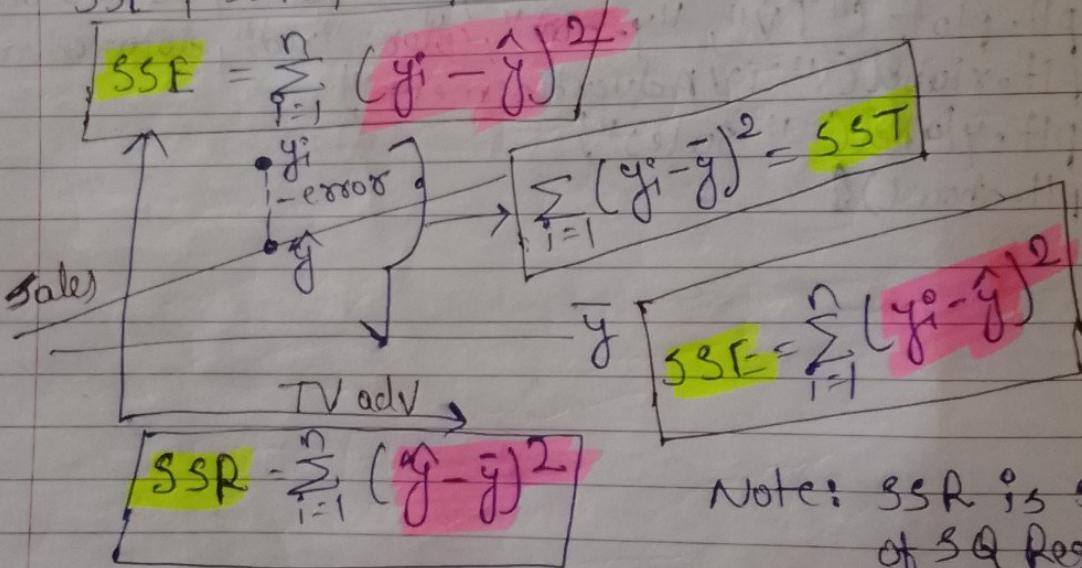
PCA — Perpendicular Distance

OLS — error is least, vertical Distance.

Ordinary least square Method

(14)

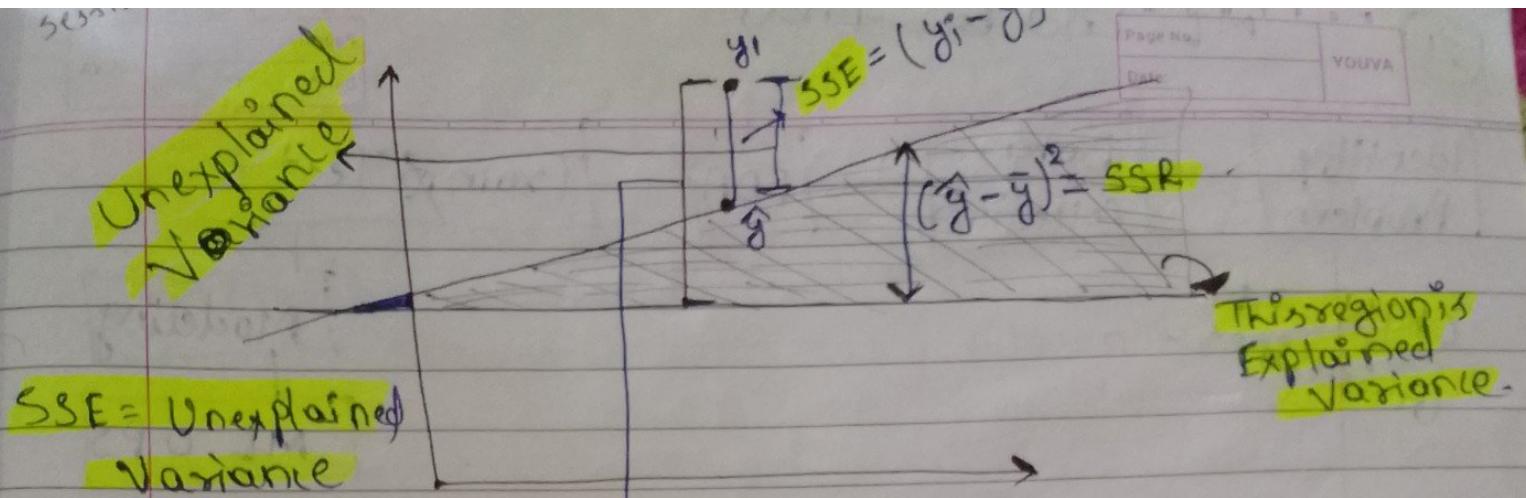
SSE , SSR , SST



Note: SSR is sum
of SQ Regression

$$R^2 = 1 - \frac{SSE}{SST} \text{ OR } R^2 = \frac{SSR}{SST}$$

OLS — fit the line in such a way that the errors are the least



$SSE = \text{Unexplained Variance}$

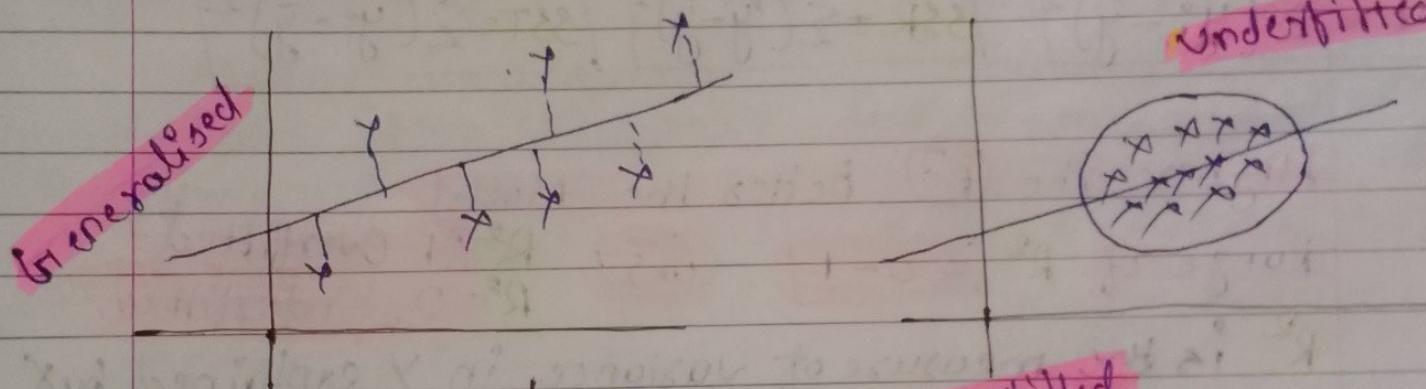
ML EQ

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + e$$

Reduce

Increase

Underfitted

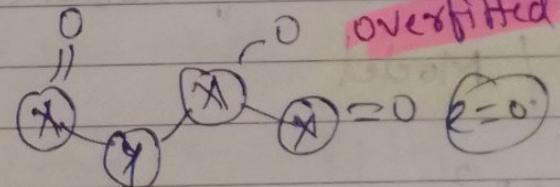


train error = high

test error = high

overfitted

train error = 0
test error => high



① Tree based ^{Model} a lot used in IT industry

①

Train error = 0

Test error = 12

M₁

Overfitting

②

Train error = 15

Test error = 15.83

M₂

Underfitting

③

Train error = 8

Test error = 7.65

M₃

Generalised

④ Shot on AWESOME A05s कम होगा वही Best Model hoga.

(21)

ML End To End

Date:

TO YOU

Identify Problem

Data Collection

EDA

Train & Test

Modeling

ML OPS

CRISP-DM Guidelines

(22)

$$R^2 = 1 - \left(\frac{SSE}{SST} \right) \text{ or } \frac{SSR}{SST}$$

$$SSE = \sum (y_i - \hat{y})^2 \quad SSR = \sum (\hat{y} - \bar{y})^2 \quad SST = \sum (y_i - \bar{y})^2$$

Higher the R^2 better the ModelRange of $R^2 = 0-1$

0.5

 $R^2=1$ Overfitted $R^2=0$ Underfitted R^2 is the measure of variance in Y explained by X (input)

LR / Model

→ R^2 → ADJ R^2

→ F stat / P value

→ Predictor P value

→ Durbin Watson

→ condition No.

→ JARQUE BER

Assumptions.

Note

while adding new features (new column in the data) & check abd k liye ki column Hamare model k liye sahi hai ya nahi, and Toh R^2 hamisha increase hoga new columns (new features add krye pr but agar adj R^2 increase hoga tbhi ye prove hoga ki NP new column / new features Hamara model acha hai

()

Shoton AWESOME A05s

① R^2 = Higher the R^2 better the model ✓

② $\text{Adj } R^2$ = $\text{Adj } R^2 \uparrow$, $R^2 \uparrow$ Good model
= $\text{Adj } R^2 \downarrow$, $R^2 \uparrow$ Bad model ✓

③ Predictor Pvalue < Alpha (0.5) Significant Predictor

④ Predictor Pvalue > Alpha (0.5) Not Significant Predictor

⑤ Fstat Pvalue < Alpha (0.5) Significant Model

Fstat Pvalue > Alpha (0.5) Not Significant Model

If all conditions met after Then we will called Our Model is good Model

R Squared

SSE

$$\text{SSE} = \text{np.sum}(\text{advertising}[\text{residual}]^{\star\star 2}) \# y_i - \hat{y}$$

SST

$$\text{SST} = \text{np.sum}((\text{advertising}["Sales"] - \text{advertising}["Sales"].mean())^{\star\star 2})$$

SSR = SST - SSE

$$r^2 = 1 - (\text{SSE} / \text{SST})$$

print("R Squared:", r_squared)

R Squared: 0.6118750508507

$$\text{Adjusted R Sq} = 1 - (1 - R^2) * (n - 1) / (n - k - 1)$$

n = advertising.shape[0] # No. of Rows in the Data

k=1 # Since it is Simple Linear Regression

$$\text{num} = 1 - ((1 - r^2) * (n - 1))$$

$$\text{deno} = n - k - 1$$

$$\text{adj-R2} = 1 - (\text{num} / \text{deno})$$

print("Adjusted R2:", 0.6099148238341623)

Shot on AWESOME A05s

F Stat & P Value

$$(24) \quad F \text{ stat} = \frac{MSR}{MSE}$$

$$k=1 - dfm$$

$$df = n - k - 1 - dfd$$

$$\frac{SSR}{k} = MSR$$

$$MSE = \frac{SSE}{df}$$

$$F \text{ test} = \frac{MSB}{MSW}$$

$$w = n - k - 1$$

(2.0) adj 1.0 > value

$$G_{17/16} \quad 6$$

$$1 \quad 2 \quad 5$$

$$2 \quad 3 \quad 6$$

$$3 \quad 4 \quad 7$$

$$n-1 + n-1 + n-1$$

$$3-1 + 3-1 + 3-1$$

⑥

$$SSR = SST - SSE$$

$$MSR = SSR/k$$

$$MSE = SSE/(n-k-1)$$

$$F \text{ stats} = MSR/MSE$$

print("F stats:", fstats)

fstats: 312.144994372713

PValue

import scipy.stats as stats

1 - stats.f.cdf(fstats, k, n-k)

1.1102230246551565e-16

Ho: the Regression model is not a significant model

Ha: the Regression model is a significant model.

Shot on AWESOME A05s

Day 2
Section 2

|| *+ T Test and P Value ||

* H_0 : There is No relationship b/w the X and Y

* H_a : There is a Relationship b/w the X and Y

** $t_{\text{stats}} = \text{coefficient} / \text{standard_error}$

$t_{\text{stats}} = 0.0475 / 0.003$
15.833333333333334

model.params
Intercept 7.032594
TV 0.047537

dtype: float64

model.params represent
the slope and
intercept

model.params[1]

Standard Error

model.bse

Intercept 0.457843

TV 0.002691

dtype: float64

model.bse = standard error

print("T TestStats:", model.params[1]/model.bse[1])

T TestStats: 17.66762560087555.

PValue * 2

df = n - k - 1

pvalue = (1 - stats.t.cdf(t.stats, df)) * 2

print("PValue:", pvalue)

T TestStats: 17.66762560087555

PValue: 0.0

Model. params = coefficient
Model. SSE = S.E.

M T W T F
Page No.:
Date: 10/10/2023

$$SE = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n-2}} \rightarrow SSE$$

Standard

Error

$$\sqrt{\sum (x_i - \bar{x})^2}$$

$$SSE \rightarrow num = \sqrt{\frac{SSE}{n-2}}$$

STANDARD ERROR

$$num = np.sqrt(SSE/(n-2))$$

$$xi = advertising["TV"]$$

$$xbar = advertising["TV"].mean()$$

$$deno = np.sqrt(np.sum((xi - xbar) ** 2))$$

$$se = num/deno$$

print("Standard Error:", se)

print("SE from Model:", model.SSE[1])

Standard Error: 0.0026906071877968703
SE from Model: 0.0026906071877968703.

Linear Regression
Degree of Freedom = n - k - 1 = 8

Confidence Interval

M T W T F S
Page No.:
Date: YOUVA

$$CI = \bar{x} \pm z(t) \times \frac{s/\sqrt{n}}{\text{slope}}$$

Confidence Interval of Slope

stats.t.interval(0.95, loc=modd.params[1], scale=modd.bse[1], df=df)
(0.04223071603269882, 0.05284256483334068)

Confidence Interval of Intercept

stats.t.interval(0.95, loc=model.params[0], scale=model.bse[0], df=df)
(6.129719268805535, 7.0935467829449855)

Revision All Formulas Till the Date

$$\text{① Slope} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\text{② Intercept} = \bar{y} - \beta_1 \cdot \bar{x}$$

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

$$\text{③ R-squared} = 1 - \left(\frac{SSE}{SST} \right) \text{ OR } \left(\frac{SSR}{SST} \right)$$

$$SST = \sum (y_i - \bar{y})^2$$

$$SSR = \sum (\hat{y} - \bar{y})^2$$

$$\text{④ Adjusted R}^2 = \frac{1 - (1 - R^2) \times (n - 1)}{(n - k - 1)}$$

Note: df Reg. = n - k - 1

$$\text{⑤ Ftest} = \frac{MSR}{MSE}$$

$$MSR = \frac{SSR}{k}$$

$$MSE = \frac{SSE}{df}$$

$$\text{⑥ Shot on AWESOME A05s} = \frac{\text{Coefficient}}{\text{SE}}$$

Linear Regression Assumptions

M T W T F S S
Page No.:
Date: YOUVA

- ① **Linearity** → **Rainbow Test**
- ② **Normality:** Residuals should be normal
→ **JARQUE BERA Test of normality**
- ③ **Multicollinearity:** Predictors should not be strongly correlated.
→ Correlation (Heat map)
→ Condition NO ($100 - 1000$) **MILD COLLINEARITY**
→ VIF range
- ④ **HETERO SKEDASTICITY**
Breusch Pagan Test.
- ⑤ **Auto Correlation of Errors DW - Test**

~~The Linear Regression Assumptions~~

Date:

YOUVA

- * **Linearity**: Statistical test for testing Linearity
is **Rainbow Test**
- * **Normality**: Statistical test for testing Normality
is **Jarque-Bera Test**
- * **Multicollinearity**: Statistical test for testing.
Multicollinearity is Variance inflation factor. This
is in addition to Correlation Plot and
Condition Number
- * **Heteroskedasticity**: Unequal Variance, so the
assumption is that there should be equal
variance and if that is not found then the
regression is not a good model. The
statistical test is **Breusch-Pagan Test**. There
is another test to check the Heteroskedasticity
(Equal variance) known as **Goldfeld-Quandt**
Test
- * **AutoCorrelation of errors**: The errors should
be correlated. The test to verify is **Durbin
Watson Test**

of ③ assumptions check one by one

Linearity Check - ①

M	T	W	T	F
Page No.:				
Date:				

H_0 : Data Model is Linear
 H_1 : H_0 is False

Our model
is linear
or not

import statsmodels.api as sma

```
teststats, pvalue = sma.stats.linear_rainbow(mod)
print("PValue:", pvalue)
PValue: 0.718500416483391
```

Conclusion: Since the Pvalue > 0.05, we fail to Reject.
Meaning that the Linearity is Present

// Normality check - ②

Normality check
Homoskedasticity
Residuals

H_0 : Residuals are normal

H_1 : Residuals are not normal

sma.stats.jarque_bera(mod.resid)

Conclusion: Reject the H_0 . It means that the Residuals are not Normal

significanceResult (statistic=151.2414204, pvalue=1.439934 - 33)
, 0.0

Pvalue = 0.0
Residuals are not normal

Residual Plot

pattern, etc
Date: 10/10/2023
Page No.: 10
Name: YOUVA

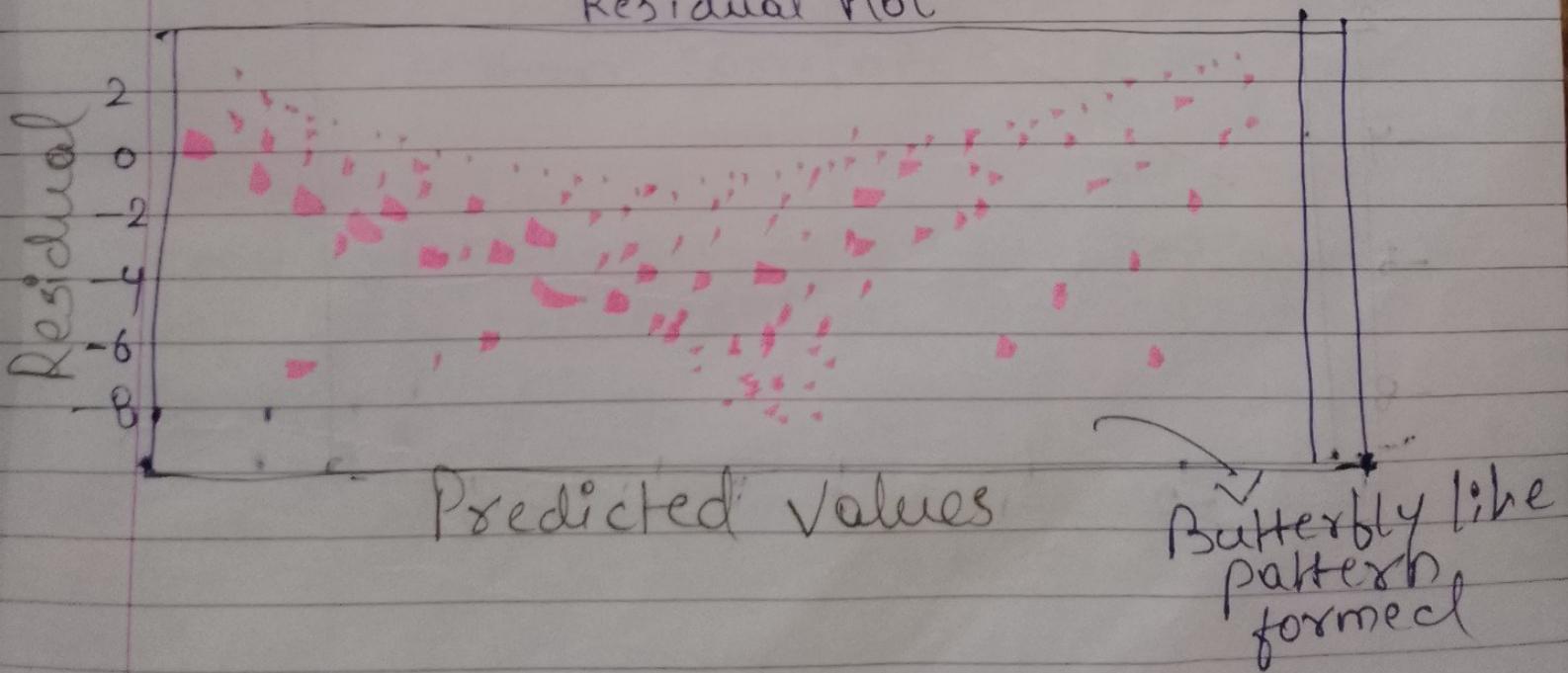
* A Residual Plot is a plot of Model's Fitted Values (\hat{y}) and Model's Residual Values.

* The important aspect of a Residual Plot is that it should be absolutely random and should not exhibit any shape.

* If there is some pattern in the Residual Plot, then it means that most likely the errors are related and Regression is not a good fit.

```
    sns.residplot(model.fittedvalues, model.resid, color="m")
    sns.set_style("whitegrid")
    plt.scatter(model.fittedvalues, model.resid, color="y")
    plt.xlabel("Predicted Values")
    plt.ylabel("Residual")
    plt.title("Residual Plot")
    plt.show()
```

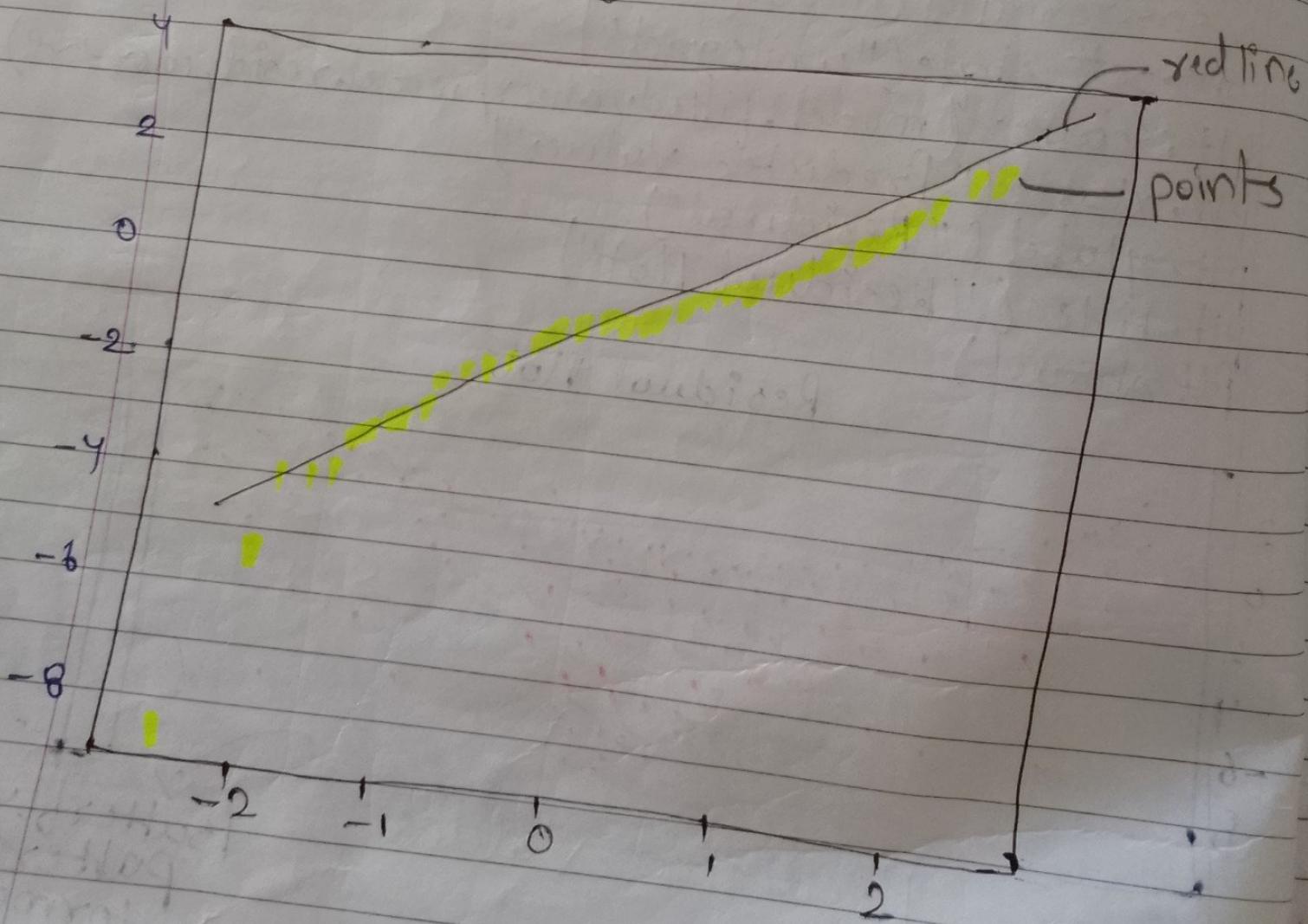
Residual Plot



ProbPlot also known as Q-Q Plot

stats.probplot(model.resid, plot=plt)
plt.show()

- # The points should stick to the red line and if that happens
- # then we assume normality. However, if the Residuals
- # are not normal
- # then the data points will deviate from the red line
- # Larger the Deviation, highly likely it's that the data is not normal



Multicollinearity

(3)

Predictor Variables should not be highly correlated amongst each other

If the prediction are Highly Correlated then is this a problem?

Yes, it is a problem because the predictors have the same prediction power and it would become difficult to identify which of those predictors are responsible in generating the prediction.

The ideal hypothesis is that the predictors should be the of independent of Nature

CORRELATION PLOT

if two predictors are highly correlated amongst each other then we will drop that predictor which less than correlation with the target

Condition Number

Condition Number is another way to find the the multicollinearity in the variables.

if the

- * CN Number is < 100 : No collinearity whereas if
- * CN Number is between $100 - 1000$: Mild collinearity
- * ≥ 1000 : Severe collinearity

Variance Inflation Factor:

- * So in this deployment, it takes all the predictions and removes the target
- * Shot on AWESOME A05s make one predictor or target, and using remaining predictors, it will build

VIF Value = 1: No collinearity
 VIF Value > 5: Severe collinearity

`sns.heatmaps(advertising.drop(['y-hat'], axis=1).corr(), annot=True)`

	TV	Radio	Newspaper	Series
TV	0.055	0.057	0.78	1.0
Radio	0.055	1	0.35	0.8
Newspaper	0.57	0.35	1	0.6
Series	0.78	0.58	0.23	0.4

Q. What is difference between Condition No. and VIF.

VIF

from statsmodels.stats.outliers_influence
import variance_inflation_factor

M	T	W	T	F	S	E
Page No.						
Date:						YOUVA

Drop the Target

```
X = advertising.drop(["Sales", "y-hat", "residual"],  
axis=1)
```

make a list and store the values of VIF

```
vif-list = []
```

```
for i in range(X.shape[1]):
```

```
    vif-list.append(variance_inflation_factor(X,  
values, i))
```

```
pd.DataFrame({ "Features": X.columns, "VIF": vif-list })
```

	Features	VIF
0	TV	2.486772
1	Radio	3.285462
2	Newspaper	3.055246.

CONCLUSION Since the VIF value is less than 5, these predictors are good to go, we can use all the predictors are good to go. we can use predictors to build the model.

However, if any of the predictor(s) had a high VIF value, we would have concluded that there is presence of High collinearity in that variable, and that feature should be excluded from the model building.



Shot on AWESOME A05s

HETROSKEDESITICITY

- * The assumption under Heteroskedasticity is Equal Variance of Residuals. In other words, it tells us that the sample drawn from the population has equal variance or not.
- * The shape of Heteroskedastic pattern is either a open funnel or closed funnel.

This represents that there is presence of unequal variance in the Data/ Model

WHAT CAUSES HETROSKEDESITICITY

presence of outliers/ Extreme points in the data is very case of unequal variance

way to Remove Heteroskedasticity
Remove the Outliers.

What's the Effect

Standard error will either increase or decrease. This will lead to Predictor being significant or not significant

✓ Confidence Interval Region will also narrow down or become wide. Net, the model will loose reliability

5

AUTOCORRELATION OF ERRORS

DURBIN-WATSON TEST

- * If there is Autocorrelation of errors, then the Metric that will be severely impacted is the Standard Error.
- * If DW Test statistic = 2 No Autocorrelation
- * If DW Test statistic is between 0 to 2 Positive Autocorrelation
- * If DW Test statistic is between 2 to 4 Negative Autocorrelation

Note: The Rule of Thumb is that DW value between 1.5 to 2.5 is acceptable

from statsmodels.stats.stattools import durbin_watson
durbin_watson(model.resid)

2.08364805294407

→ NO AUTOCORRELATION

Day 4 Session 3

FEATURE SELECTION METHODS

- (A) When we build a model, we include all the features to build our model.
- (B) One way to find if the feature is significant or not is by applying Statistical Tests during the EDA stage.
- (C) The Machine Learning way is to use a wrapper method known as Sequential Feature Selector coming from Mlxtend Library.

How SFS works?

THE CONCEPT IS CALLED FORWARD SELECTION OR BACKWARD SELECTION TECHNIQUE

Note: There is another approach known as Step wise Regression Approach

How SFS works?

(A) Forward Selection Approach: It will start with 1 predictor and build the model. It is a blank slate initially and will start with 1 predictor and capture the Adjusted R².

(B) Then, it will add another feature, build the model and measure the Adjusted R² square. If the Adj R² score has increased, it will keep the predictor else remove the predictor and try the remaining predictors.

(C) This process will keep on going till all the predictors are not exhausted. In the end, we will get the final list of the predictors.

Shot on AWESOME A05s

x_1	x_2	x_3	x_4	y	① Build The Model using all features $y = b(x_1, x_2, x_3, x_4)$ Pvalue.
—	—	—	—	—	

- ② $[x_1 | y] \rightarrow$ LR Model $\rightarrow R^2 \rightarrow 0.65$ Statistical Approach
- ③ $[x_1 | x_3 | y] \rightarrow$ LR Model $\rightarrow R^2 \rightarrow 0.78$

$$[x_1 | x_2 | x | y] \rightarrow R^2 = 0.65$$

④ The Similar Approach works for Backward Elimination Technique where the model will be build using all the predictors and then the predictors will be removed recursively to find the impact on Adjusted R². The final list of the important features will be generated in the end.

अपने इन प्रैदूरी feature selection को कैसे करें है?

```
# ! pip install mlxtend
from mlxtend.feature_selection import SequentialFeatureSelector
```

```
from sklearn.linear_model import LinearRegression
```

FORWARD SELECTION APPROACH \rightarrow 2 दिए गए ताजे उत्तरों से पहले फॉरवर्ड सेलेक्शन करना है।

```
SFS = SequentialFeatureSelector(LinearRegression(), forward=True, scoring='r2')
```

CV=5 पाच बार अपेक्षा दिए गए हैं।

(forward Selection approach की ये लिखा कृति कैसे है ये फॉरवर्ड सेलेक्शन करना है।)

It will (SFS) will be applied on X and y

```
sfs.fit(X, y)
```

Extract the Best features...

```
print(sfs.k_feature_names_)
```

⑤ Shot on AWESOME A05s

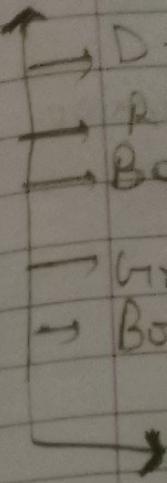
Backward fit

sts = SequentialFeatureSelector (estimator = lr, forward=True,
 false_scoring = "f1", cv=5,
 k_features = "best")

Day
Forecast

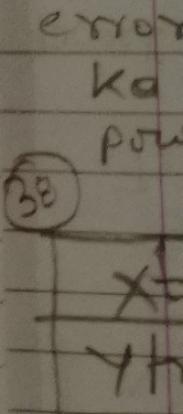
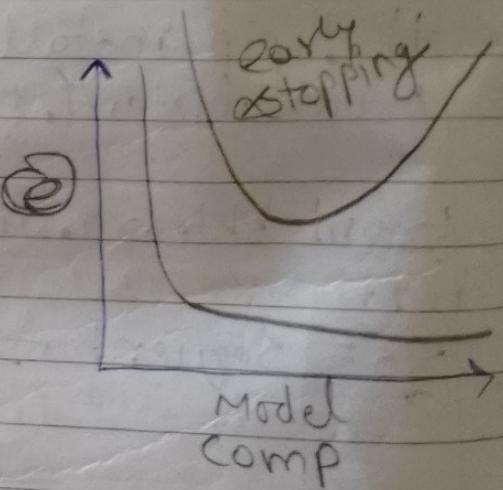
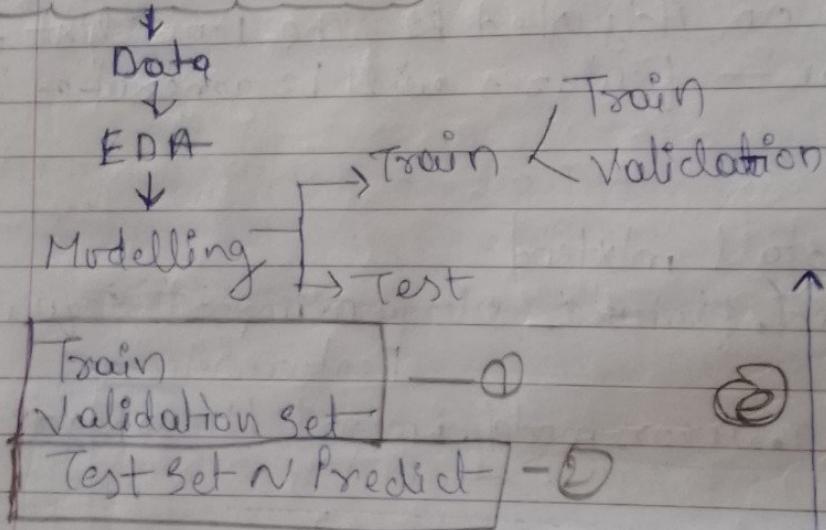
sts.fit(X, y)
print(sts.feature_names_)

['T1', 'RH_1', 'T2', 'RH_2', 'T3', 'RH_3', 'T4', 'RH_4',
 'T5', 'RH_5', 'T6', 'RH_6', 'T7', 'RH_7', 'T8', 'RH_8',
 'T9', 'RH_9', 'T-out', 'Press_mm_hg',
 'RH-out', 'Windspeed', 'Visibility']



Cross Validation

Problem statement



Cross Validation → cross validation means machine
ont ko baar den karne aadhi karne se error ko
azat marni hain but den time aik baar increase ho
vai bhi vishay graph mein koi 2nd 3rd 4th kar.

Day 6
Session 1

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

→ D-Tree

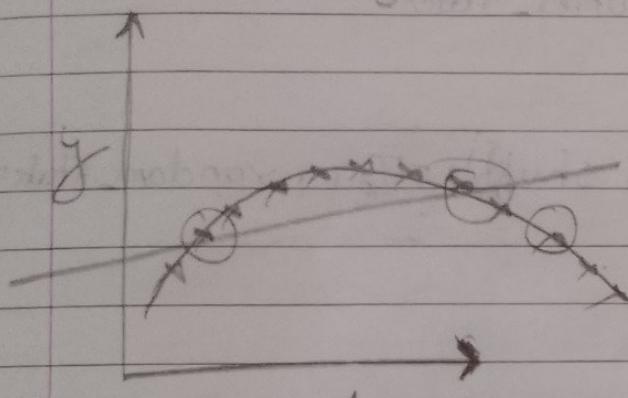
→ RF

→ Bagging

→ Gradient

→ Boosting

→ Overfitting All High Variance Problems



linear Regression

Model Exist

Underfitted Model

$$y = mx + c$$

$$x^{1.2} / x^{1.5} / x^2$$

अगर Linear Regression model mein error तोड़ता है तो यह predictors का power असे होता है और means $y = mx + c$ में x की power असे होती है कि यह curve बन जाता है।

(3)

Data

		80-20	
X-train	X-test	Y-train	Y-test

Cross Validation

k fold CV

epoch	c=1	T ₁	T ₂	T ₃	T ₄	Test

epoch	c=2	T ₁	Test	T ₂	T ₃	T ₄

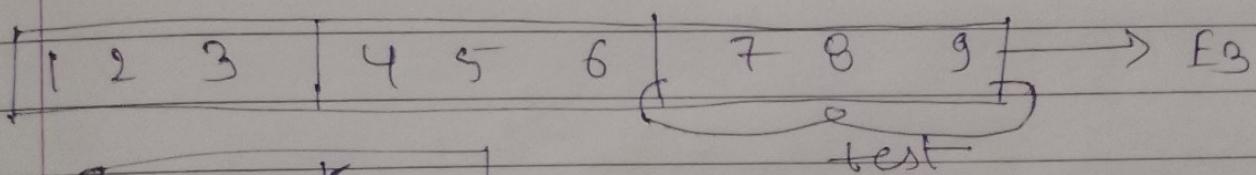
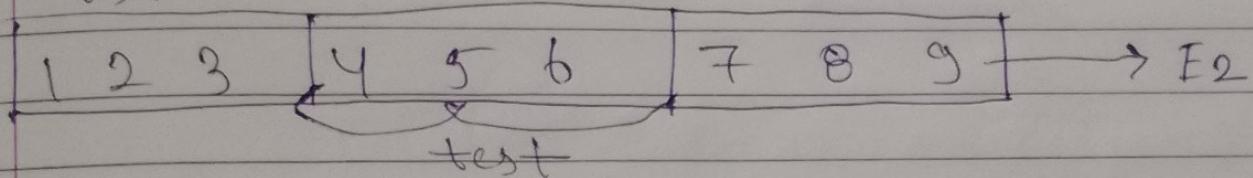
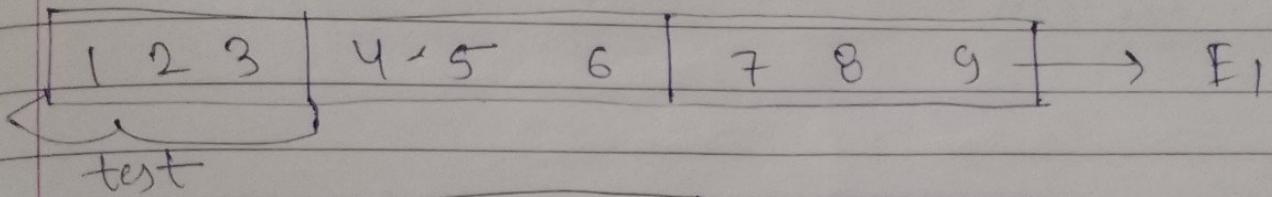
epoch	c=3	Test	T ₁	T ₂	T ₃	T ₄

Shot on AWESOME A05s

Cross validation

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

3 - Fold Cross Validation



$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

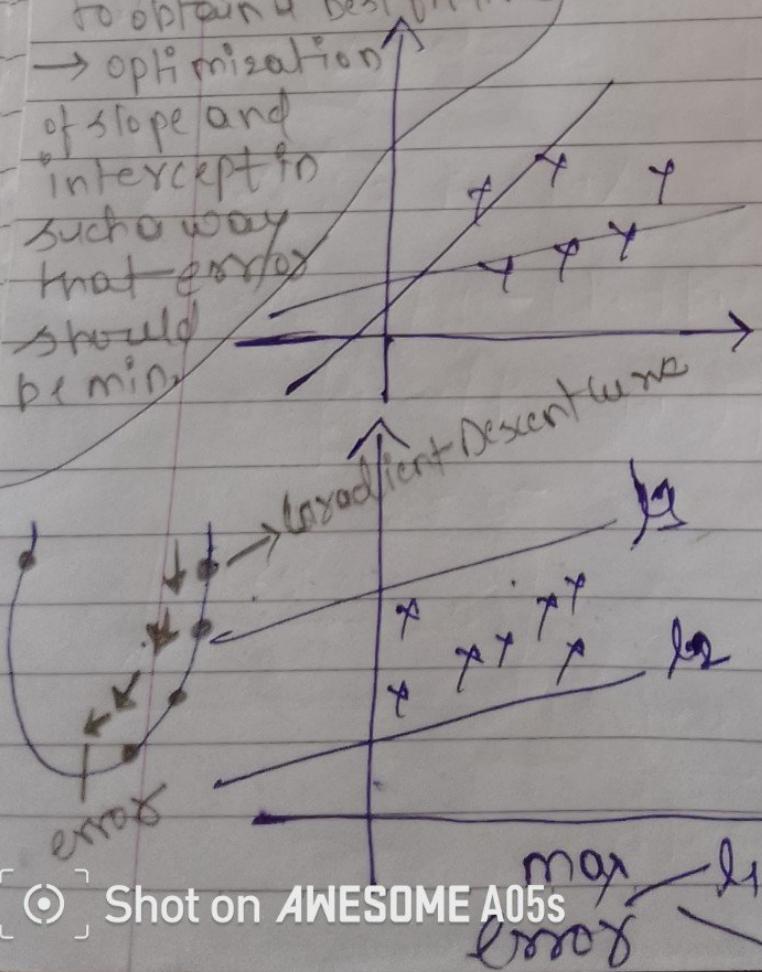
cross validation solves the problem
of → test overlapping.

Week ② Day ① Session ①

→ Job of gradient descent is to reduce the error in such a way that we are able to obtain a best fit line

→ Optimization of slope and intercept in such a way that error should be min.

Gradient Descent



OLS - stat model
sklearn

`OLS("y-x", data).fit()`

$$\beta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\hat{y} = \beta_0 + \beta_1 x_i$$

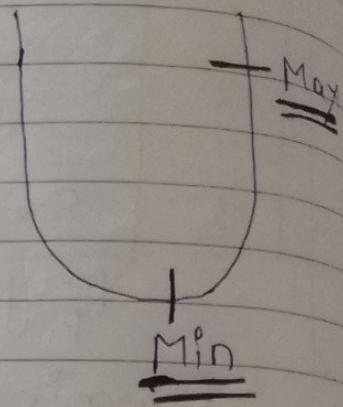
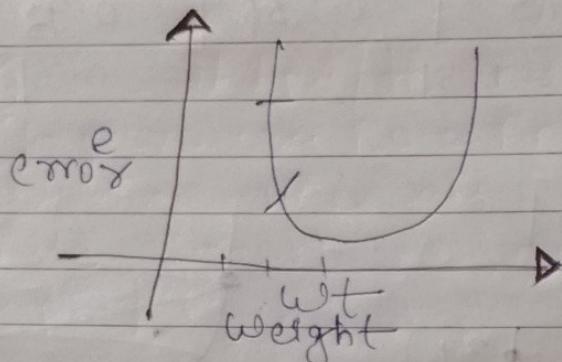
↳ predicted \Rightarrow Residual

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

(53) Convergence

→ when in 2 or more than 2 iterations, the error reduction is minimal, that is the point of convergence

So, goal is to reach from maxima to minima



(56)	cat
Auto Saler Cat	
Data	
X	Y
M.T	0
AT	1
	8
	16

(54)

Gradient Descent

→ Batch Gradient Descent
(Underfitted)

→ Stochastic Gradient Descent
(Overfitted) (Row by row)

→ Mini-Batch Gradient Descent
[Best of three]

Batch \leftrightarrow 1000 Rows.

\rightarrow 100 Rows

100 Batches

(55) Sklearn \rightarrow Grid Search

$$SSE_{0.01} = ?$$

$$SSE_{0.001} = ? \text{ min.}$$

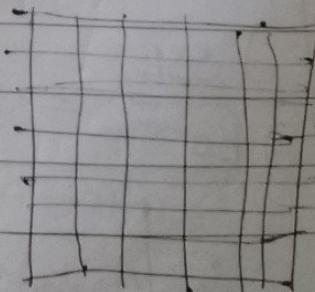
$$SSE_{0.1} = ?$$

$$SSE_1 = ?$$

$$\text{Rate} = \eta [0.01, 0.001, 0.1, 1] \\ CV = 5$$

Grid Search (estimator =
param_grid = params)

$$LR = 0.001$$



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

$$(x_1) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$w_1$$

$$f(x) \xrightarrow{\text{Pred}} \hat{y}$$

$$w_2$$

$$w_3$$

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

R Tanweer

MSE

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

RMSE

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

coefficient

Day 2
Session 2

Regularization

overfit

Ridge

Lasso

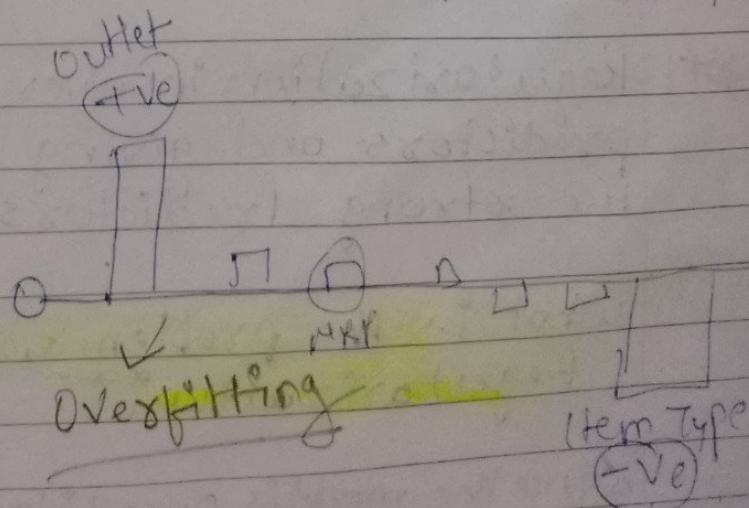
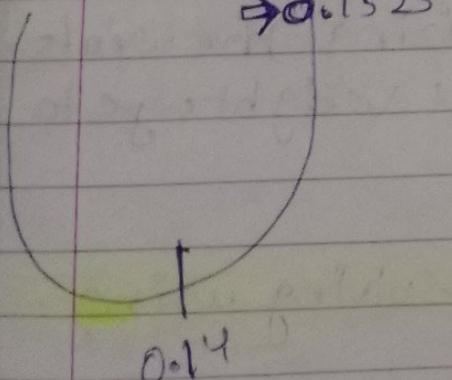
Elastic Net

$$SSE = \sum (y_i - \hat{y}_i)^2 + \lambda \cdot \sum_{i=1}^p (\beta_i)^2$$

$SSE = 0.14 \quad \lambda = 0.05 \cdot \sum_{i=1}^p \beta_i = 0.25$

$0.14 + 0.05 \cdot 0.25 \rightarrow 0.1525$

$\rightarrow RMSE = 0.14$



Overfitting

Item Type
-ve

Ridge

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

Ridge Regularization

$$\sum_{i=1}^n (y_i - \hat{y})^2 + \lambda \sum_{i=1}^p (\beta_i)^2 \Rightarrow \sum_{i=1}^p \beta_i = 0.25$$

Regularization Parameter

Regularization Parameter = Bias

(Q)

$$\text{Outlet ID} = 10 \quad \hat{y} = \beta_0 + 1.5 \times MRP + 10$$

$$MRP = 1.5$$

$$SSE + \lambda \sum (\beta_i)^2$$

$$\beta_1 = 10 \approx 5$$

$$\beta_2 = 1.5$$

$$= \frac{10}{OUTLET} \quad \frac{1.5}{MRP}$$

Q what is Regularization Parameter?

Dealing with

Dealing with

~~Dealing with~~ ~~Dealing with~~

Regularization \rightarrow penalizing the weak Predictor and giving the weightage to strong Predictors.

Regularization

*

Regularization involves penalizing the weak predictors and giving correct weightage to the strong Predictors.

**

what is the problem we are solving using Regularization??

Since the model is overfitted and gives weightage to the weak predictors. Thus we introduce

Shot on AWESOME A05s

★ How it works ★

It will introduce a Bias in the model and the equation for me said Model under Ridge Regression would be $SSE + \alpha \cdot (\text{sum of coefficient weights})^2$

WHAT REALLY HAPPENS IS WE INTRODUCE BIAS TO CONTROL THE VARIANCE

L2

$$\text{Ridge} = \sum_{i=1}^n (y_i - \hat{y})^2 + \alpha \cdot \sum_{i=1}^p (\beta_i)^2$$

weak predictors
ko kam weightage karega

no. of columns/features
Data small

$$\beta_1^2 + \beta_2^2 + \dots + \beta_p^2$$

L1

$$\text{Lasso} = \sum_{i=1}^n (y_i - \hat{y})^2 + \alpha \cdot \sum_{i=1}^p |\beta_i|$$

weak predictors
ko high weightage karega

when you have more than 100 features.

Plastic Net

$$\begin{cases} \text{Lasso} & \left\{ \begin{array}{l} aL_1 + bL_2 \\ a+b=1 \\ L_1-\text{ratio} = \frac{a}{a+b} \end{array} \right\} \\ \text{Ridge} & \left\{ \begin{array}{l} aL_1 + bL_2 \\ a+b=1 \end{array} \right\} \end{cases}$$

If $L_1-\text{ratio} = 0$

$$\frac{a}{a+b} = 0 \quad a=0 \quad \left\{ \begin{array}{l} a+b=1 \Rightarrow b=1 \\ a=0 \quad b=1 \end{array} \right.$$

Q. When L1 is 0, then what's the model about? +21?

$$L_1-\text{ratio} = 0$$

$$\therefore L_1-\text{ratio} = \frac{0}{a+b}$$

$$\frac{0}{a+b} = 0$$

$$a=0$$

$$n_1 + n_2 = 1 \Rightarrow b=1$$

Shot on AWESOME A05s

$$L_1-\text{ratio} = 1$$

$$\frac{0}{a+b} = 1$$

$$a=a+b$$

$$b=0$$

$$a+b=1$$

$$a=1$$

$$aL_1 + bL_2$$

$$L_1-\text{ratio}=0$$

L1-L2

else

As per Regularization Model, we are increasing Bias to Reduce the Variance (Overfitting)

So, when in a model we increase Bias to reduce variance or increase variance to reduce bias, such a relationship is called Bias Variance Trade off.

Ridge Regression

`ridge = Ridge(alpha=10)`

for

`xtrain`

`ytrain`

`xtest`

`ytest`

`yhat = ridge.fit(xtrain, ytrain).predict(xtest)`

`RMSE: 1021.20594`

`RMSE: 973.50693`

`RMSE: 931.671122`

`RMSE: 1017.5012`

`RMSE: 1000.1234`

Lasso Regression

`lasso = Lasso(alpha=50)`

for

`xtrain =`

`ytrain =`

`xtest =`

Shot on AWESOME A05s

`yhat =`

`lasso.fit(xtrain, ytrain).predict(xtest)`

```

from sklearn.model_selection import GridSearchCV
params = {"l1ratio": [0.1, 0.12, 0.2, 0.5, 0.65, 0.7, 0.75, 0.85], "alpha": [1, 2, 3, 4, 5]}
grid = GridSearchCV()
estimator = (net, param_grid=param_grid, scoring="rmse", cv=5)

```

grid.fit(X, y) # Fitted the Model

```

> GridSearchCV
> estimator: ElasticNet
    > ElasticNet?

```

```

grid.best_params_
{'l1ratio': 0.1}

```

cross validation iteratively
 splits the data set into two
 portions: a test and a training.
 The prediction errors from each of
 the test sets are averaged to
 determine the expected prediction
 error for the whole model

K-fold - Cross-validation

- ✓ Ⓛ Data is divided into 'k' folds or subsets
- ✓ Ⓛ The model is trained 'k' times, each time using k-folds for training and remaining fold for validation
- ✓ Ⓛ Average performance across all folds is often used over the model's generalization
- ✓ Ⓛ row by row

Leave-one-out Cross Validation

- ✓ Ⓛ 'k' is set equal to the number of data points.
- ✓ Ⓛ Each data point serves as a validation set exactly once, while the rest are used for training
- ✓ Ⓛ Provides a more thorough assessment but can be computationally expensive for large datasets.
- ⌚ Shot on AWESOME A05s

K fold splits the data into K folds and validation, whereas leave-one-out validation, whereas leave-one-out validates the model iteratively, leaving out one data point at a time, K-fold is computationally more efficient, while LOOCV provides a more exhaustive evaluation at a cost of increased computation.

Leave One Out

This will give 100% efficiency.

So we have never used it for modelling.

→ If no. of features are best, we have option of providing best no. of features.

Ex. min best ~~soam ho kuch nahi kota~~, tumhe manually karne hoga.

GridSearch CV

- ① set the learning rate range ✓
- ② train model with each learning rate ✓
- ③ evaluate the performance ✓
- ④ choose the learning rate that gives best result. ✓

Grid Search CV is a technique for finding the optimal parameter values from a given set of parameters.

Ay