

Libraries

Seaborn

Matplotlib

Plotting

Pandas

Numpy

Data

Manipulation

Scikit-learn

Tensorflow

PyTorch

XGBoost

Statsmodels

→ Data
prediction

To install kernel

pip install ipykernel

To install libraries

PIP install numpy, Pandas,

Scipy, seaborn

Libraries name.

Extensions to install in

VS code

→ Text explorer UI

→ Python

→ PyLance

→ Jupyter

→ Material icon theme

→ Python extension pack

jis pc par salikne se uski
information ac g;

→ Intellicode.

Exploratory Data analysis

- 1) Identify patterns, trends and relationships.
- 2) spot errors or unusual data points.
- 3) Potential explanations for observations.
- 4) Data cleaning and transformation.

(1) Numpy

Numerical computing. Support for arrays, matrices etc.

(2) Pandas

It provides series data (one dimensional) or data frame (2-D)

1) Data cleaning

2) Data transformation

3) Handling Time series data

3) Matplotlib

Line plots, bar charts, histograms,
scatter plots etc.

Graphs

2) Histogram

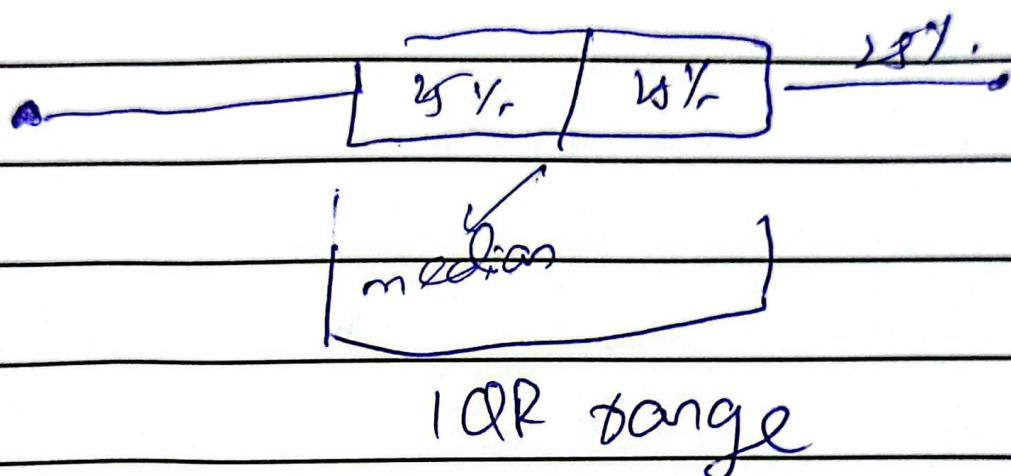
→ use when data is numerical.

→ use to check normal distribution
of data.

2) Box plots

used for numerical data

2) we have a clear understanding about outliers, data range about lowest and highest and we have an idea about skewness.



3) [Scatter plot]

To observe and show relationship between two numeric variables. Use to find co-relational relationships.

4) Heatmap

Dark colors indicate stronger co-correlation, while light colors indicate weaker co-correlations.

- 1) Positive co-correlations are usually represented by warm colors such as red or orange.
- 2) negative co-correlations are usually represented by cool colors, such as blue or green.

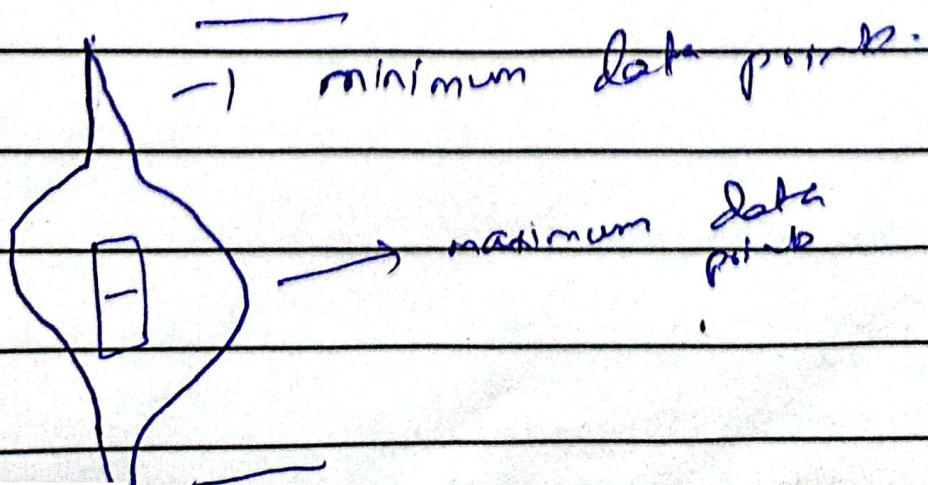
5) Pair plots

It combines histogram/ KDE and scatter plots

used because others are univariate i.e if species are focus, use pair plot.

6) Violin plots

Used to observe distribution of numeric data.



script-18

Outliers removal

- 1) If data is normally distributed then Z-score.
- 2) If data is not normal then IQR.

One hot encoding is not

- ↪ suitable when data is high dimensional it further increases dimension -

2) Label encoding.

- 3) Frequency → if many categories
 { then we use this.

for big value data & then
use frequency.

1) Ordinal encoding

jo orders hain, usi orders pe encoding

Feature Engineering

→ 1) Creating new feature by combining existing feature.

2) Encoding categorical variable

3) Feature selection.

i) Filter method → filters unnecessary data.

ii) wrapper method + forward selection or backward elimination

one by one variable add us in prediction models check user give

when target is numeric, then regression
when target is categorical, then classification

iii) Dimensionality reduction

↳ PCA, T-SNE -

iv) Regularization method

↳ Ridge and LASSO regression

etc

v) Deep learning Based method

Random Forest, Decision trees.

(4) Feature scaling

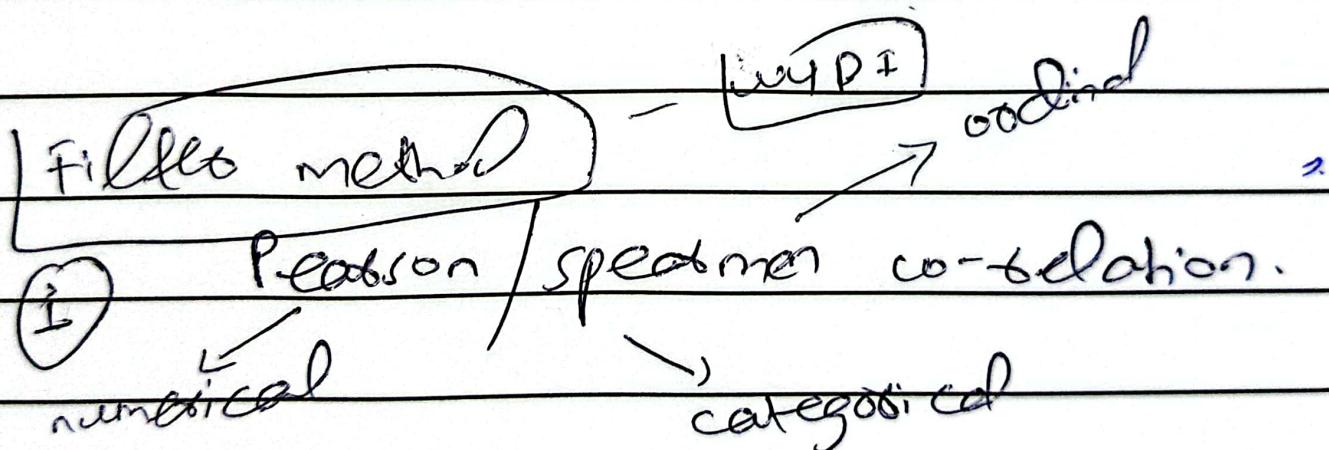
(5) Handling missing data

Scale sb column ko kyun ge

but transformation sift un main
jin main issue ho.

~~feature selection vs feature extraction~~

Domain knowledge has to
feature selection.



⇒ Between features correlation

ni hua chahiye agar hogi
to multi-linearity.

if more than 2 variables
then 2 or + test
more than 2 variables

Written by FAISAL JILLANI SATTI

faisal.qau2018@gmail.com

feature able

numerical

② **ANOVA** → analysis of variance
mean testing

③ chi-square test → it features
are categorical

co-variance tells direction

in which variables vary. but
correlation tells the amount
between 0 and 1

⇒ For supervised learning
we prefer that variables
are independent -

⇒ variance main histogram ka
ic jiski variability kam hai
usko consider klein and jip main
zada usko drop.

$R^2 = 0.34$ means 34%.
Written by FAISAL JILLANI SATTI
due w+21

faisal.qau2018@gmail.com

one by one step
include no h
deletion gE Real

Wrapper method

elimination

- i) forward selection
- ii) backward

iii) Recursive feature elimination

iv) Exhaustive feature selection.

Features extraction

distance
measures
(co-selection)

j) features all have w information
defy how we combine w many
hair. \Rightarrow Dimension reduction

- i) Principal component analysis (PCA)
when components are linearly
related : means straight line.

Distance measures

i) Euclidean

ii) Manhattan

ii) T-SNE

non-linear relationship

infandata

Tree distributed standardized
neighbourhood embedding

iii) Linear discriminant analysis LDA

$\frac{1}{1}$
linear relationship

Time series Data

regular

trend, seasonal, residual

model \rightarrow additive, multiplicative

perplexity means average
points Vega
idea value 5 to 30

high dimensional to 2D main idea
has

B2ed:

Null-hypothesis and alternate hypothesis

LOA

To check association / ~~indep~~ between two variables.

If not associated then independent.

class

Scree plot

→ If plot is straight then importance of all is equal.

→ jahan bend ae ga, top se vahan
ki certain kamyais

Data Science

Machine Learning

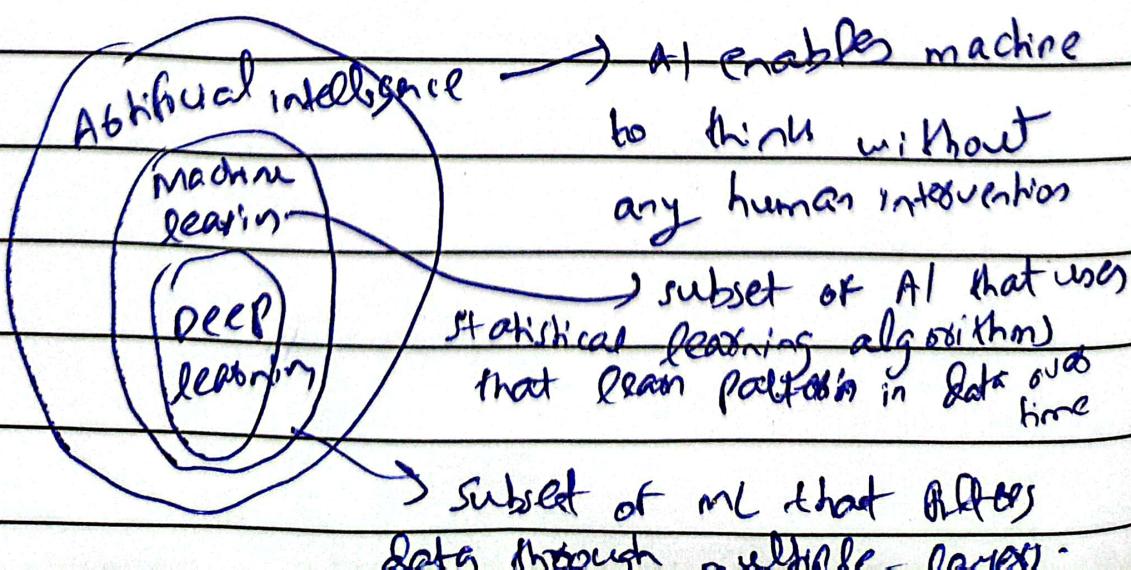
we humans learn from experiences

and then implement it for the future.

So in machine we learn (learn) from data (for prediction / classification)

→ Machine Learning

It is a branch of artificial intelligence that focuses on building systems that can learn from and make decisions based on data.



Errors → Due to some natural factors

Bias → Due to human errors

↳ What am I doing is what has thought machine

If target variable is numerical
then regression and if
categorical then classification.

Decision Tree and Random Forest

Random Forest is combination of
Decision Trees.

For machine learning tasks, we need
to split data in two sets.

Training dataset and testing

dataset: $(70, 30)$ vs $(80, 20)$

Overfitting of model

not reliable information because trained on noise (irrelevant data)

i.e if training and testing errors are low then good fitted -

i.e if training error is low ~~then~~
and testing error is high
then overfit.

Causes

- i) complex model.
- ii) Too many features
- iii) Insufficient training data.

Solution

- i) Simplify model.
- ii) Regularization (L_1 shrinkage method and L_2 is shrinking as well as

Hyper parameter tuning Ridge and LASSO

Written by FAISAL JILLANI SATTI

feature & domain
collinearity meth.

feature selection method.

iii) cross validation

5 to 10 fold

aggregates from all test and training
in each bracket & then then second
phase main way shuffle user.
one bracket select user for training
and one position for test and
so on.

Understanding of model

Training and testing both error
high.

of

Bias-variance trade off

Bias → underfitting

Variance → overfitting

Captures noise
and fluctuations

Balancing bias and variance-

Workflow of ml

monitor and update model and data

retrieve data

deploy to production

clean and explore

validate/ evaluate
model.

develop and train model

prepare/ transform

Model performance on training data is called calibration

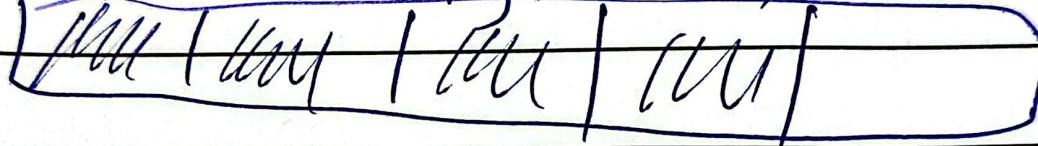
→ Model performance on test data is called validation

Validation type

k-fold cross validation

Training

Test



If k is 5 then we divide data into 5 part and changing testing one part in every iteration.

Monte carlo cross validation

k-fold main hm ne fix kya tha kisra part train and test hm but is main randomly



randomly hash per 4 pack Kohi
per 3 etc.

③ Bootstrapping → Data increase
or sampling with replacement

overall data main ra randomly
sample pic hotay hash.

~~for sample th koi~~ han jb

consistency check koi ho ya

data karo to enhance koi
chart.

Evaluation measures (Regression)

① R² beta hai ki model kaha
acha fit hoa hai. (or it)
forse.

②	MSE	Mean Square error
	RMSSE	Root mean square error
	MAE	Mean absolute error.

↓
ham ne jo prediction ki hain
and jo actual data tha uskiy
~~se~~ jo ham ne test data man
predictor ko poshae kia tha
uskiy satr difference vo kaha
aa gaya hai.

$$MSE = \frac{y - \bar{y}}{n-1} \quad \begin{matrix} \approx y = \text{actual value} \\ \bar{y} = \text{predicted value.} \end{matrix}$$

Evaluation measures (classification)

Confusion matrix

		0	1
		TN	FP
Actual	0		
	1	FN	TP

Sensitivity / Recall

How much were correctly identified as positive to how much were actually active.

Specificity

how much were correctly classified as negative to how much were

actually negative.

Precision

How much ~~are~~ correctly classified as positive out of all positives

F1 score tells overall performance of binary classification model.

PCA

co variance matrix

$$\begin{pmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{pmatrix}$$

diagonal = variances

off-diagonal = co-variance

Eigen values and eigen vectors

from determine what basis the
principal components sufficient
basis.

just eigen value greater than
one has us component no
certain principal basis.

Linear regression

Dependant variable (y) → target
Independent variables (x) → feature.

Simple linear regression → we have
one dependant and one independent

Multiple Linear Regression

when features are more and target value is one.

High Dimensional Data \rightarrow

no of features $>$ no of observations

Multiple linear regression

β_0 = intercept

β_p = slope

ϵ = residuals or error.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

intercept \rightarrow agot effect n' aa
 data + pe to y intercept
 mean advertising n' wi less
 phis be much na much sales
 ho oae hor.

OLS

| ordinary least squares regression

minimize sum of squares of ~~error~~

| Ridge regression

$$\beta = \arg \min \left(\sum y_i - \beta_0 - \sum \beta_i x_{ij} + \lambda \sum \beta_j^2 \right)$$

λ = ridge penalty for multi-collinearity

in ridge, lambda (hyper-parameter)
as ridge penalty gets higher,
coefficients get closer to zero
mean variations lower.

If λ
we check for optimal value of λ
Re MSE kann ae.

if lambda = 0, then linear
regression

Lambda means dm bias put when ge

test size = 0.3 means 30% for test
random-state = 42 means seed point

OLS Regression results

Adj R-Square \rightarrow irrelevant info
ni beta so ye
wala De�na ha -

F-statistic

check significance
that model is good fitted
 ϕ not

$P(F\text{-stat}) < 0.05$, there is
significant contribution.

AIC, BIC should be low.

Akaike information criterion

Bayesian information criterion.

$|P| > |t|$

if less than 0.05 it
means significant

against values bracket u bracket
main aa sac, it means
no significant, null hypothesis
not rejected

SKEWNESS

if 0 then data is normal
greater than 200 positively skewed.
less than -200 negatively skewed

KURTOSIS

leptokurtic if > 3 means heavy
tail distribution.

If $= 3$ then mesokurtic means normal

if < 3 phlogistic means flat.

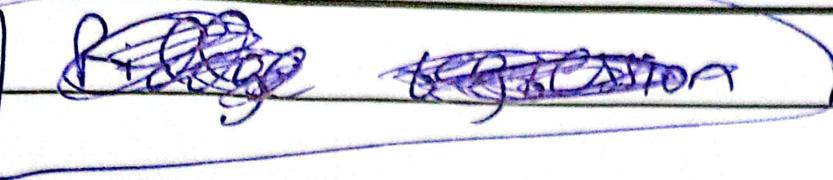
Durbin - Watson Test (cont)

Tells about auto correlation
means all observation will
be dependent to nth term.

if $= 2$, there is no auto correlation
if > 2 high auto-correlation
if < 2 negative correlation

(Durbin - Watson)

if P value is less than 0.05
then normally distributed
then non-normality



Steps for EDA

check missing values and
impute.

- 1) Investigate each variable
- 2) If variables are quantitative
then check their distribution.
if categorical then label the
data like label encoding, one
hot encoding etc and see their
pattern by using box plot,
pie charts.
- 3) Scale and standardize the
features for avoiding any
influence of high or low
magnitude.
- 4) After scaling and standardizing
next step is transformation
(if required). when we

, quantitative

plot histograms or pair plots. if there is any non-normality then we need to transform that.

- 5) Then plot data again to see.
- 6) if not normalized then apply other transformation methods.
- 7) Then remove outliers ^{feature selection} and β
- 8) Then modelling.

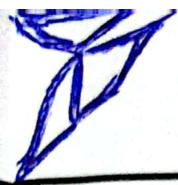
\Rightarrow L2 regression.

Ridge regression

excess shrinkage

Lambda (~~alpha~~ ridge penalty) add us tay hasn't to re-define coefficients. takay predicted actual line k agreeb aajae.

- i) used mostly when large data
- ii) and when multi-collinearity is a serious issue.



β_1 = slope coefficient

β_0 = intercept

initiative

f jo unimportant features hotay hain
hen unkey beta co-efficient ko zero k
qaseeb hotay hain.

y alpha → determines the weighting
to be used.

selection
L¹

→ Lambda / alpha jb zero hoga to
means linear regression because
 B^2 bhe zero hojae ga.

not
use

Ridge main Lambda penalty
dalray k bd the beta co-efficient
kahi exactly zero ni hotay, blukh
zero k qaseeb hotay hain.

at

Ridge cv se best alpha find
hotay hain.

LASSO Regression

Least absolute Shrinkage

and selection operators.

→ The benefit is it also provides features selection, with the help of ridge, we cannot select features.

→ LASSO shrink unimportant features co-efficients to zero.

→ It also help to reduce multi-collinearity

⇒ Target variable ko scale ni karey

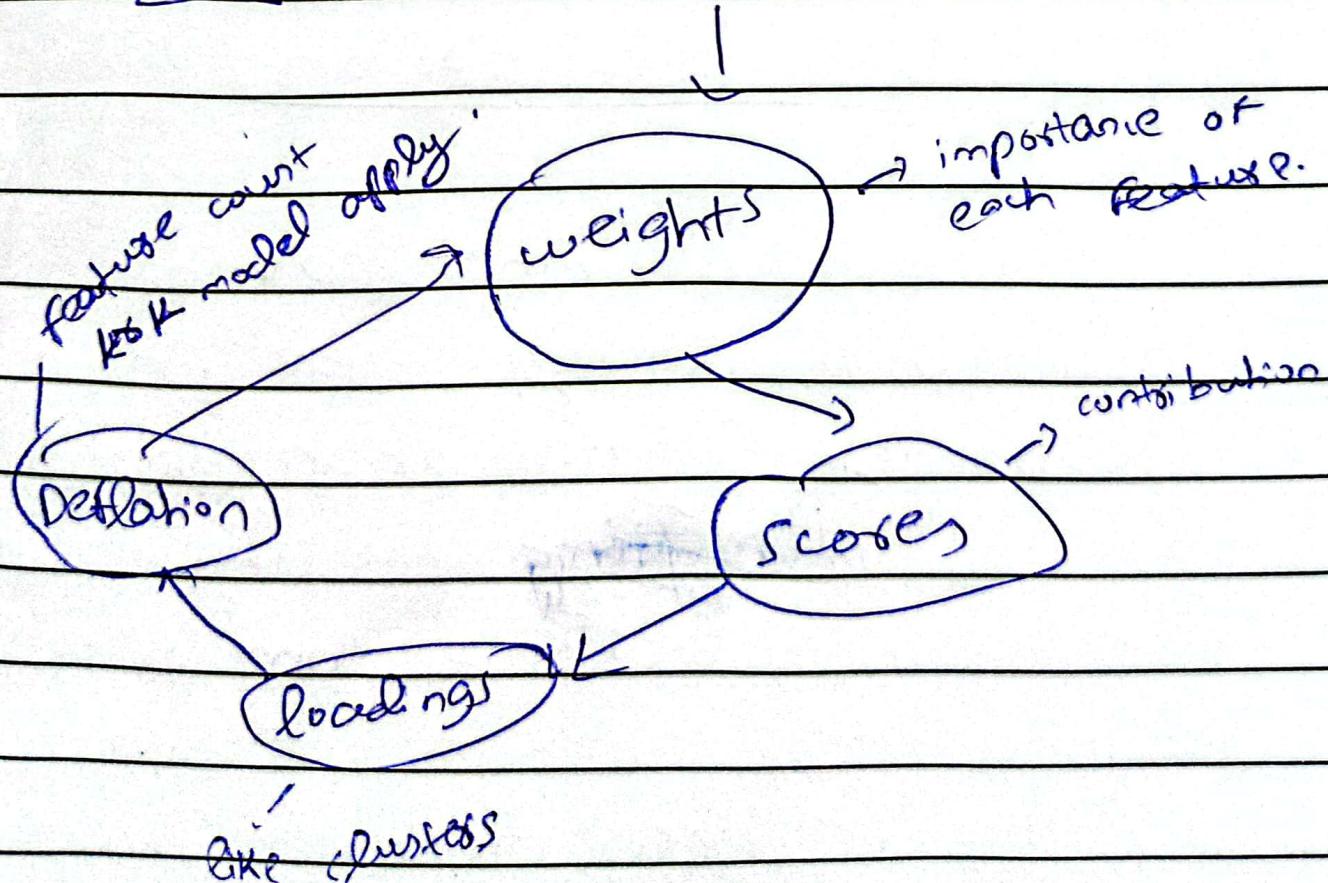
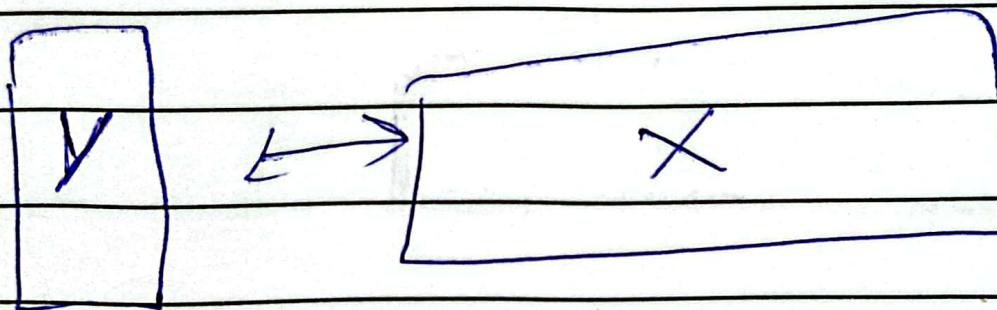
⇒ unimportant features ko jib

or hm drop waley hain
to automatically ~~feature~~

~~select~~ dimension reduction
hoga hain.

Partial Least squares regression (PLS)

purpose same as LASSO, Ridge, but it is used for High dimensional dataset.



pta chol jata hai utnay component

for MSE ham hai, means utnay

component PC maximum information

mil ga hai.

- > PLS generates principal components, hence reduces dimensions.

I Classification

II Logistic regression

supervised machine learning algorithm widely used for binary classification tasks.

such as identifying whether an email is spam or not, presence and absence of disease, but also be used for multi-class.

→ It can be used for image classification.

Linear Discriminate Analysis

- Classification algorithm and also reduces dimensions.
(Supervised)
- PCA reduces dimensions on the basis of linear correlations.
- LDA reduces dimensions on the basis of class labels ~~—~~
correlations both.
- LDA is for that type of tasks where linear relationship exists means jis no straight line se separate vr sken. For non-linear QDA/RDA is used.

PCA pre-processing main apply
kta has and classification ni
whta and unsupervised learning
main ada hai. LDA is
supervised.

- ROC curve tells us about accuracy.
- Suitable for low-dimensional data

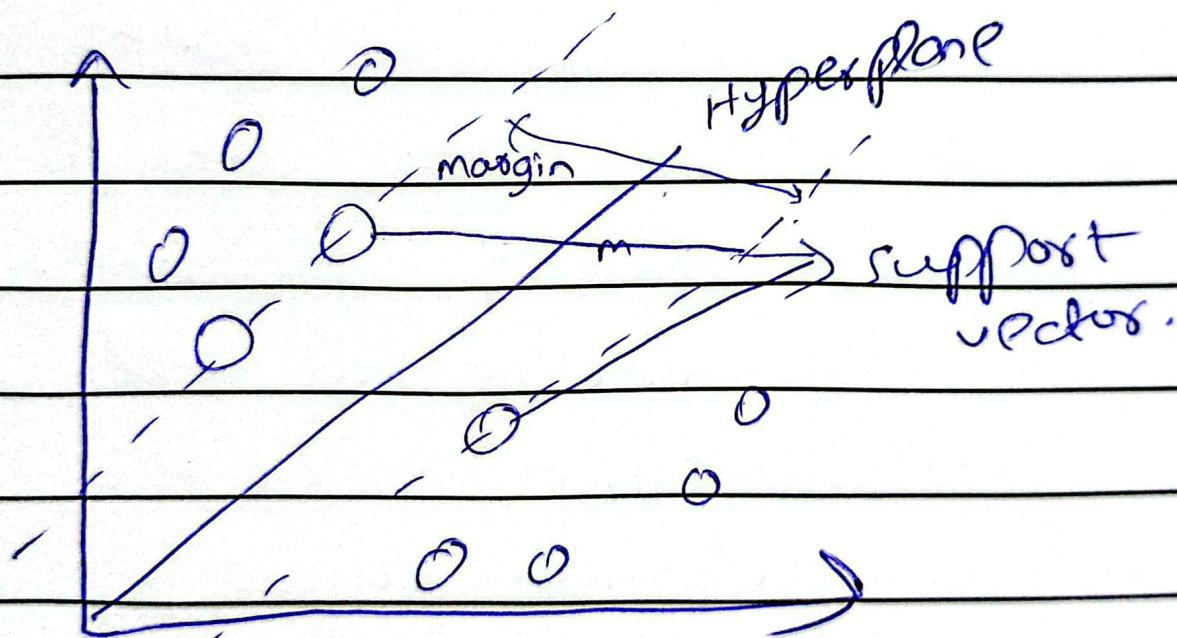
Support vector Machine SVM

- classification algorithm.

- SVR for regression.
- Low dimensional data ko project kota hai high-dimensional space pe

Kernel trick se non-linear decision boundaries ko handle karta hai.

* Binary and multi-class classification dono ko kernel karta hai.



(Common Kernel Trick)

Linear, polynomial, radial basis function (RBF) and sigmoid.

Decision Function

- Classifies new data points based on their position relative to hyperplane.

KNN

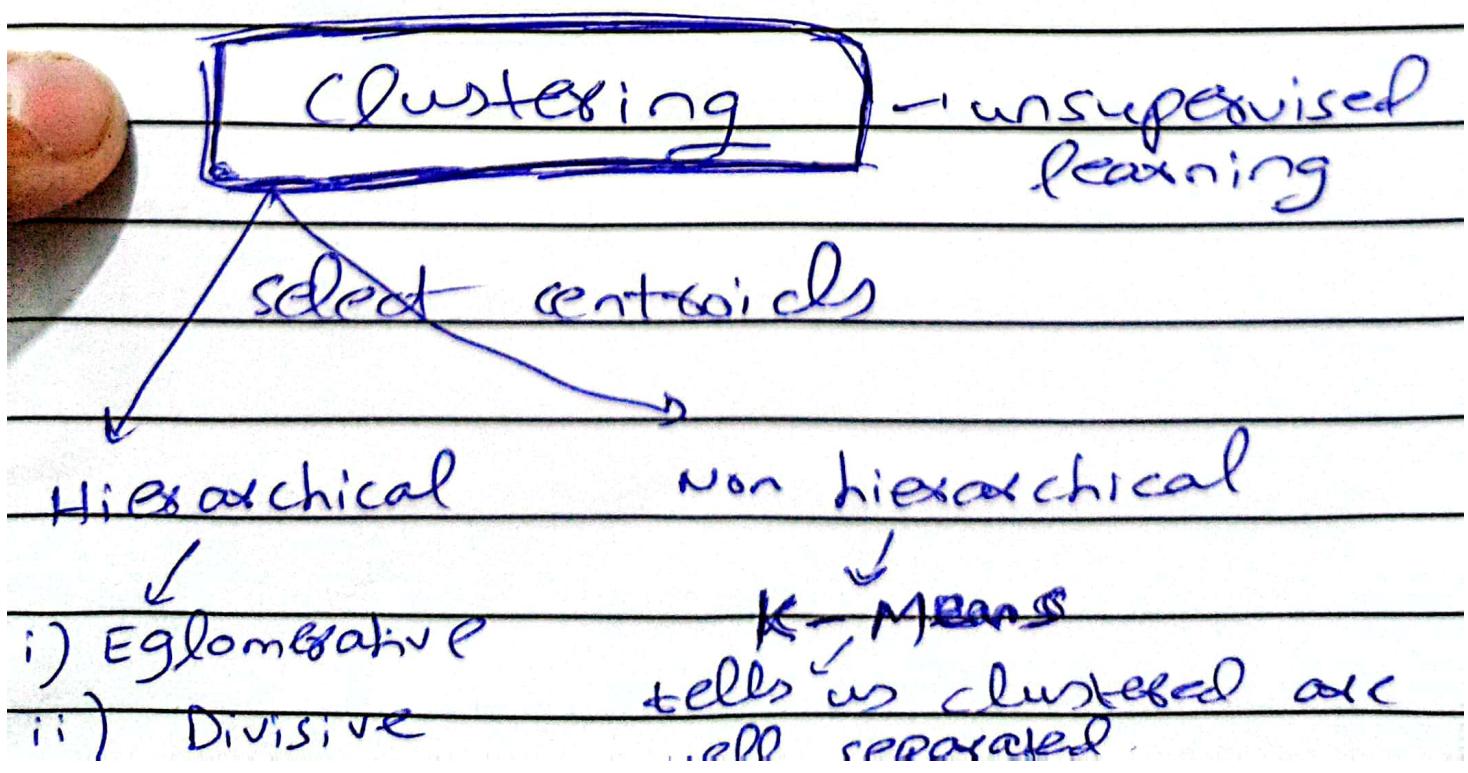
Supervised Learning

stores all available data

- and classifies new data point based on similarity. It can be used for regression as well as classification but mostly it is used for classification purpose.

It is also called lazy learner algorithm because it does not learn from training set immediately instead it stores the dataset and at the time of classification, it performs an action on dataset.

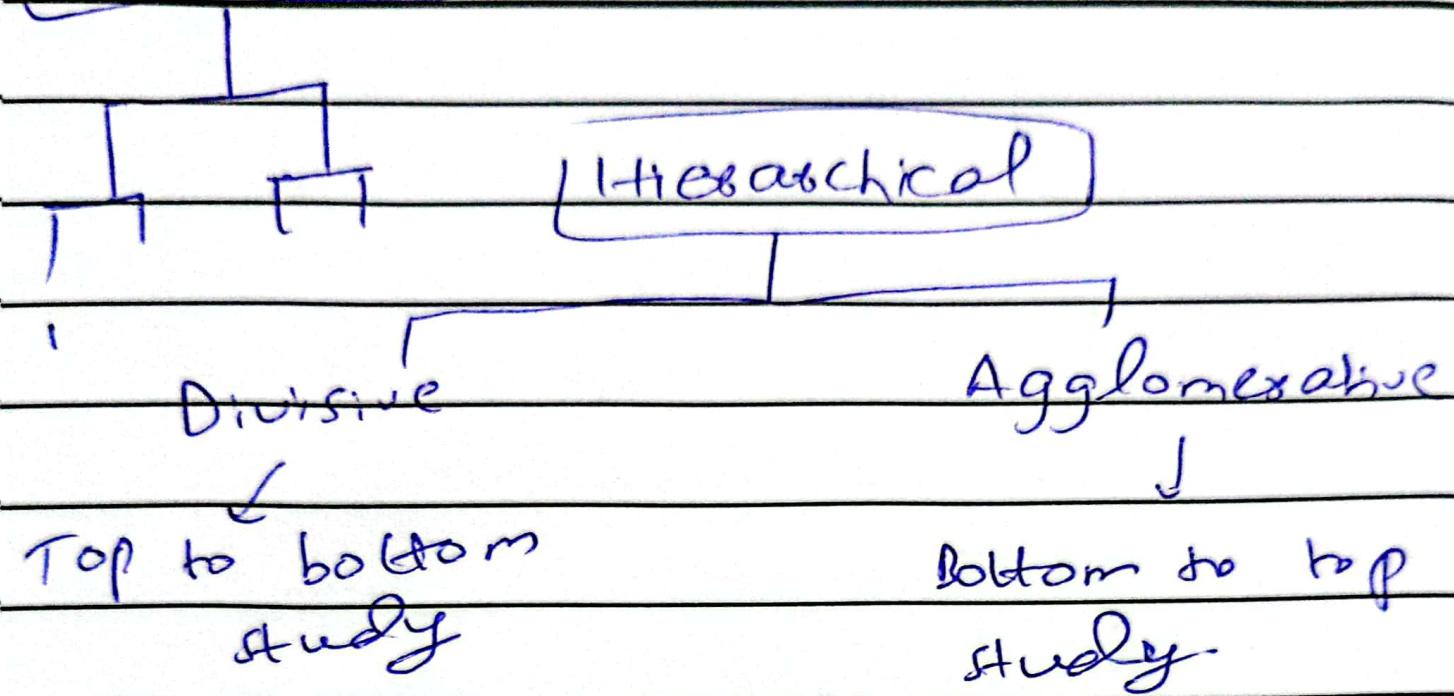
- Step 1) Select the number k of neighbours.
- Step 2) calculate Euclidean Distance of k -numbers of neighbours.
- Step 3) Take the k -nearest neighbour as per calculated Euclidean Distance.
- Step 4) Among these k -neighbours, count number of data points in each category.
- Step 5) Assign new data points to category for which the number of neighbours is maximum.



Dendrogram

Written by FAISAL JILLANI SATTI

faisal.qau2018@gmail.com



Neural Networks

→ Deep learning

both for regression and classification

→ inspired by human brain

TPU's → Tensor processing unit
advanced form of GPU.

RNN → when time oriented
image data

CNN → image datasets, video datasets, classification tasks.

- Deep learning model pre-processing whole variety task.
- Perception main probabilities calculate.



$$\sum f_i$$

Sum of probabilities
so decide

Sigmoid and ReLU most popular activation functions.

Activation function in context of neural networks is a mathematical function applied to output of neuron.

The purpose of activation function is to introduce non-linearity into

- the model, allowing it to learn and represent complex patterns in data. without non-linearity, neural network will behave like a linear regression model, regardless of number of layers.
- > The activation function decides whether a neuron should be activated or not by calculating weighted sum and further adding bias to it.
 - > In neural network, we would update the weights and biases of neurons based of the error at the output. This is known as back-propagation.
 - Activation functions make back-propagation possible.

→ Colab or kaggle notebook for free if our RAM is low

→ Base theorem.

Public Datasets

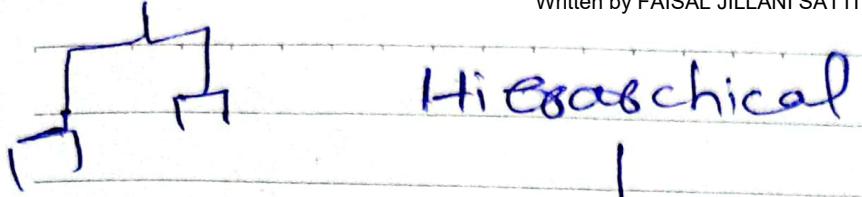
- i) Image → Microsoft COCO
- ii) Text → SQuAD.
- iii) Video → ~~Youtube~~ YouTube 8M
- iv) Audio → Google Audio set.

Libraries

- i) Scikit-learn → Machine Learning
- ii) Tensor Flow → Google.
- iii) PyTorch → Facebook 2016
- iv) Keras → 2019 most popular.

Existing architecture

- i) Image classification → Resnet.
 - ii) Image segmentation → Uonet
 - iii) Text classification → Best (Transformer)
 - iv) Image transformation → Pix2pix
 - v) Object detection → YOLO
 - vi) Speech generation → wavenet.
- Radial kernel → if polynomial structure
- RBF → for complex tasks.
- Designed for quantitative datasets but can handle image datasets
- need large data to train well.



Hierarchical

Divisive

Agglomerative

Top to bottom study

Bottom to top study

Convolutional neural networks (CNN)

Computer vision
algorithm

Computer vision is mostly tasks
main use here is CNN,

- Image recognition.
- object detection.

- Transfer Learning .

Multiclass \rightarrow softmax.
Binary class \rightarrow sigmoid.

Gross - entropy

act node se
ketni information
aa sae hai
and wo waha
important hai

Date

ID main usay li 28 orat ni ph
like ANN.

LP main data de leha hai.

- ~~position~~
- RNN mostly use usay hash

Homework

- i) padding layers.
- ii) Filters → edge detection (3×3)
- iii) pooling layers.
max, min, average. mostly

$$(n-m+1) \quad (n-m+1) \quad \begin{matrix} 6 \times 6 - \text{original} \\ \downarrow \\ n \end{matrix}$$

$$(6-3+1) \quad (6-3+1) \quad \begin{matrix} \text{new filters} = 3 \\ \downarrow \end{matrix}$$

$$\boxed{4 + 1}$$

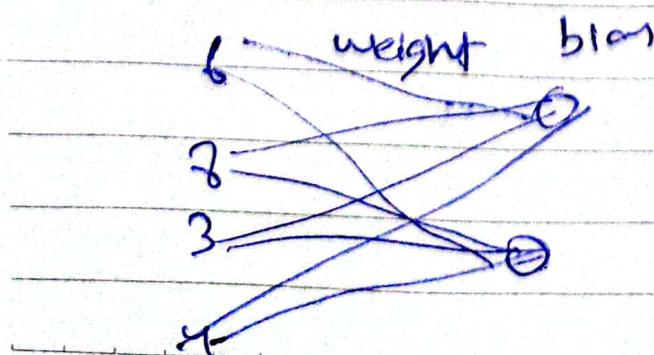
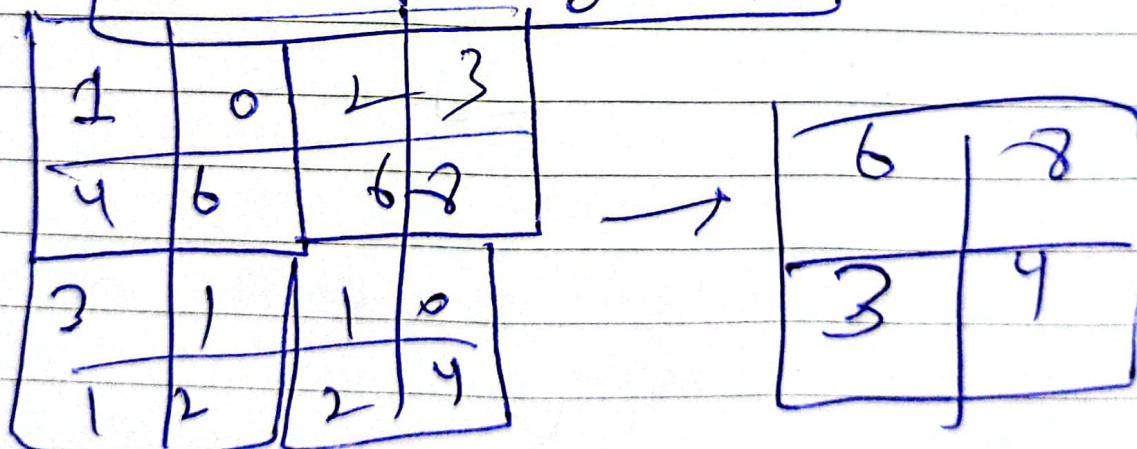
so filters se dimension kam ho ja raha.

side val values size one time ah
 ha so who cater wonay h
 liye extra layer lgakay ha
 that is known as padding layer
 pooling main usay han

strides :- determines how many
 squares or pixels out
 after step when they move
 across the image.

note kitha wna ha, jump kitha
 wna ha.

Max Pooling

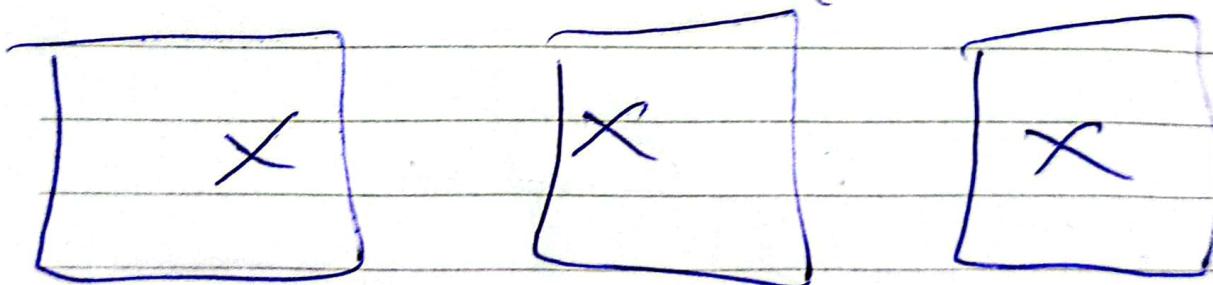


CNN

four layered concepts.

- 1) Convolution. → Filters. Feature map.
- 2) ReLU → Activation function.
- 3) Pooling.
- 4) Fully connected layer.

Pooling layers add because not
on same grid ~~feature~~



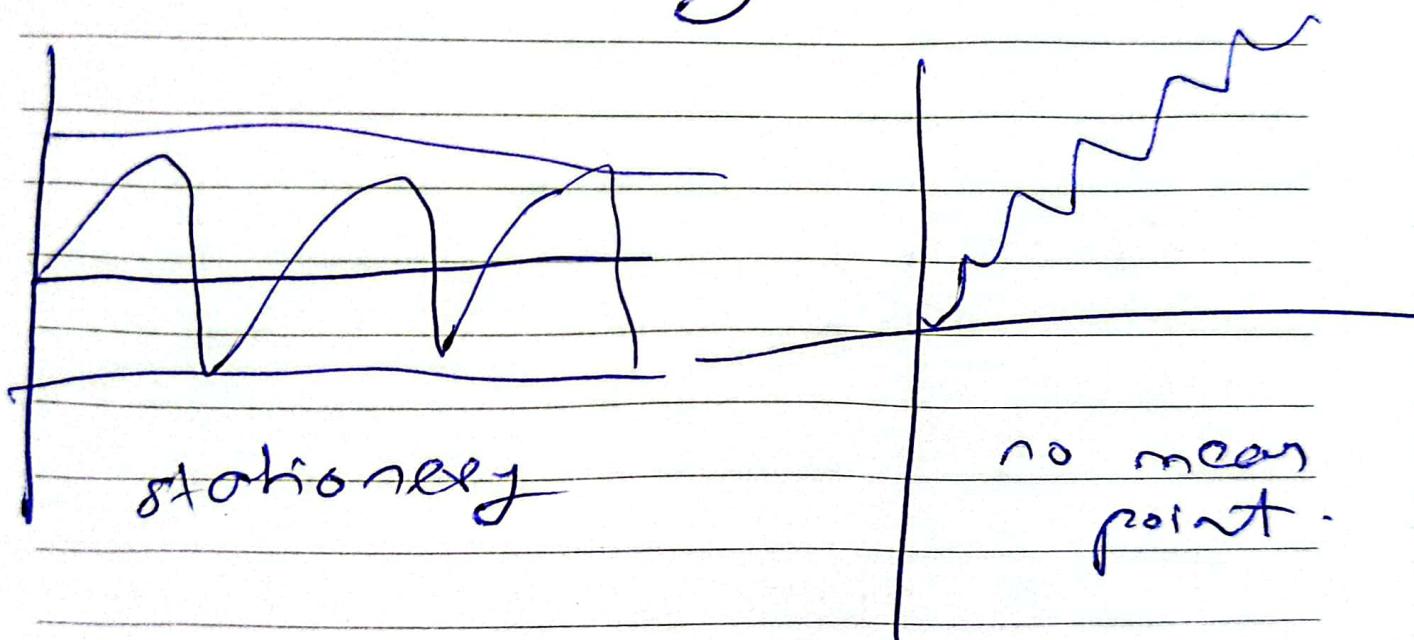
→ MNIST dataset available in
Python library.

- Automatic driving car
- Facial recognition

DENSE \rightarrow Folds layers in CNN
 Dropout \rightarrow To reduce dimension.

Time-series analysis.

For TSA, we need to make data stationary.



Trend, seasonality, cycle, irregular in sb ko remove kro kia aagay chalna hai means Data Stationary.

interpolation and extrapolation.

\nearrow
 within existing values

TSA main line plot e bhatay hain.

Date

Date no stationary, may use
smoothing techniques.

→ Holt's linear → use when data
have trend.

.. Holt's winters seasonal →
when both trend and
seasonal.

$$\begin{aligned} AR &\rightarrow P \\ I &\rightarrow d \\ MA + Q \end{aligned}$$

10 steps in

Recurrent Neural network

Time series, Date both
image and numeric.

when Date is sequential.

Email may be RNN work in
that ha.

bank e use utay ha.