

★ Process of Data Science

- ① Refinement of Problem Statement
→ AIM
- ② Data Acquisition
- ③ Data preparation / preprocessing
- ④ EDA → Exploratory Data Analysis
/ feature selection
- ⑤ Model Planning → Finalizing Model

- ⑥ Visual Communication
- ⑦ Deployment

- HW
- D- Kaggle download - Titanic DB
review & discuss in next class

similar wanted questions with doubts etc (iv)
(iv) 7/7/2023 . 7th ←

- ★ Know your data ★
- i) Size of data
→ df.shape
 - ii) How data look like
→ df.head()
 - iii) Detail of data
→ df.info()
 - iv) Calculate Null value
→ df.isnull().sum()
 - v) How data look mathematically
→ df.describe()
 - vi) To check duplicate
→ df.duplicated().sum()
 - vii) To check the correlation between columns
→ df.corr()

* EDA univariate

I) Categorical data

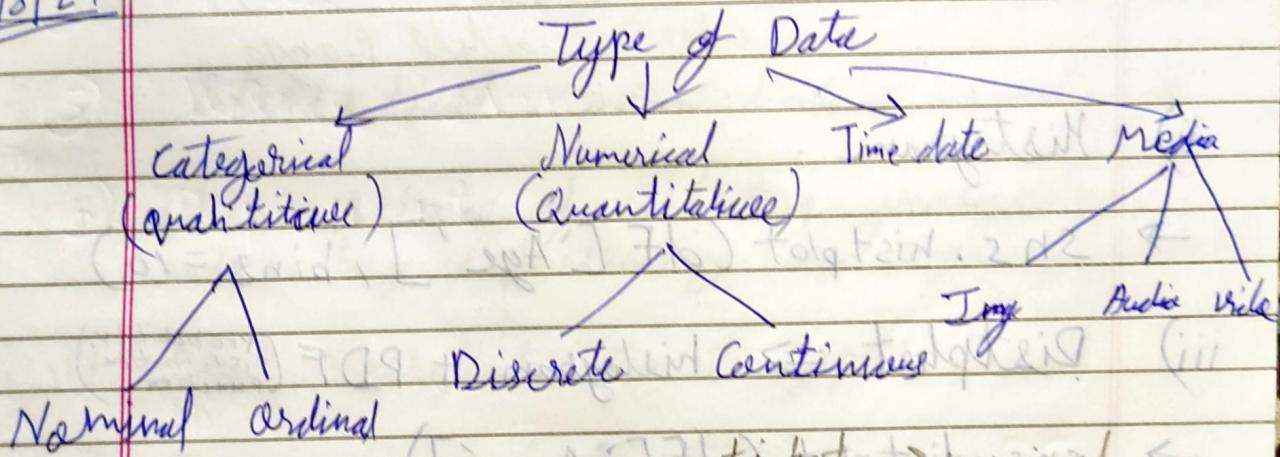
i) Pie chart

`df['Age'].value_counts().plot(kind='pie', autopct = '%.2f')`

ii) Count Chart

`sns.countplot(df['Sex'])`

7/8/24



P_ID → Nominal Parch → Nominal

Survived → Nominal Fare → Continuous

Pclass → Ordinal cabin → Nominal

Name → Nominal embark → nominal

Sex → Nominal

Age → Discrete Continuous

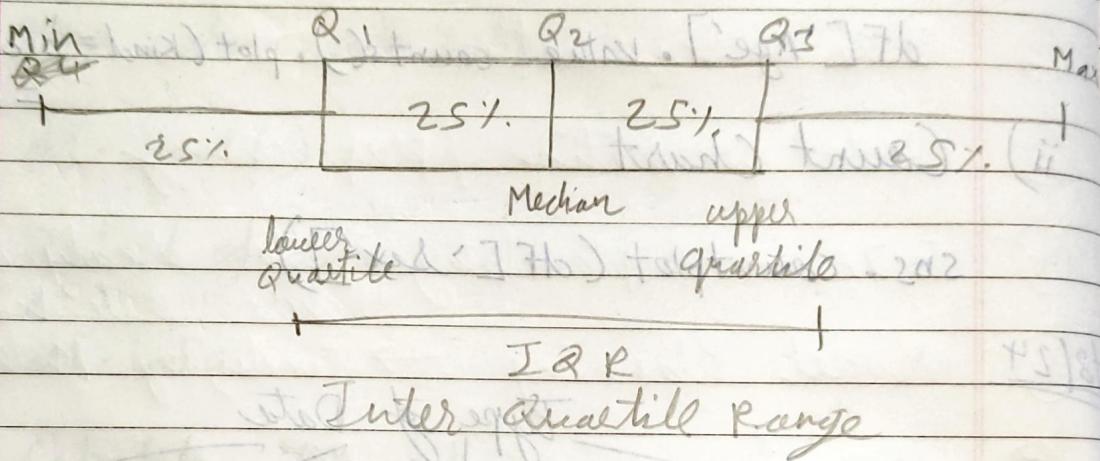
SibSp → Discrete

PW → Car Detail → Aggyle datatype

II) Numerical data

i) Box Plot

→ sns.boxplot(df['Age'])



ii) Histogram

→ sns.histplot(df['Age'], bins=10)

iii) Distplot \Rightarrow histogram + PDF (Probability Distribution Function)

→ sns.distplot(df['Age'])

★ EDA Bivariate / Multivariate

i) Heat Map

YH DIFF probability vs likelihood

- Probability - It is measure of likelihood of an event occurring
- Random Variable - variable
- Probability Distribution - description of possible value of each random variable
- Conditional Probability

$$P\left(\frac{B}{A}\right) = \frac{P(A \cap B)}{P(A)} = \frac{P(B)}{P(A)} P(A|B)$$

- Statistics
- i) It is a science of collecting, analyzing, interpreting, presenting and organizing data
 - ii) Classified into 2 types
 - a) Descriptive Statistics
 - b) Inferential Statistics

- iii) Also helps in hypothesis testing
- iv) Correlation -
- v) Regression - it is a model, the relationship between dependent variable and one or more independent variables

~~12/8/24~~

* Measures of Central Tendency .

$$i) \text{ Mean } (\mu \text{ or } \bar{x}) = \frac{\sum A_i f_i}{\sum f_i} = \frac{(A_1) f_1 + (A_2) f_2 + \dots + (A_n) f_n}{(f_1) + (f_2) + \dots + (f_n)}$$

* Measure of Dispersion .



i) Mean

- central tendency

$$\bar{X} = \frac{\sum_{i=1}^N X_i}{N}$$

ii) Variance

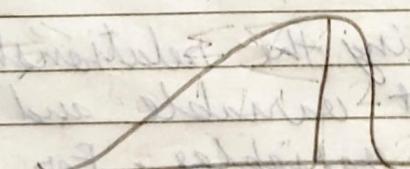
- Dispersion

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N}$$

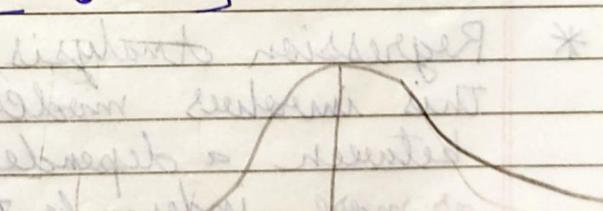
iii) Skewness

- Symmetry (positive or negative)

$$\text{Skew} = \frac{1}{N} \sum_{i=1}^N \left[\frac{(X_i - \bar{X})^3}{\sigma} \right]$$



-ve Skewness



+ve Skewness

iv) Kurtosis

- Shape (Tall or flat)

$$\text{Kurt} = \frac{1}{N} \sum_{i=1}^N \left[\frac{(x_i - \bar{x})}{\sigma} \right]^4$$

- Mesokurtosis - An excess kurtosis of 0

- Platykurtosis - A -ve excess kurtosis.
- Leptokurtosis - A +ve excess kurtosis

13/8/24

* ANOVA - Analysis of variance

This involves comparing means across multiple groups to determine if there are any significant differences. For ex, comparing the mean heights of individuals from different regions.

* Regression Analysis

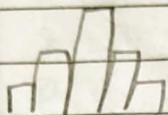
This involves modelling the relationship between a dependent variable and one or more independent variables. For ex, predicting the sales of a product based on advertising expenditure.

* Chi-Square tests

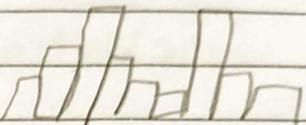
This involves testing independence

* Shapes of Histogram

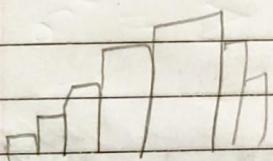
i) Symmetric



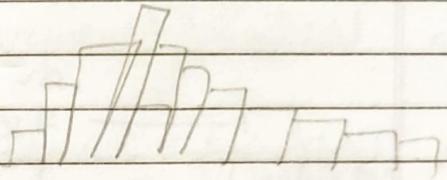
ii) Bimodal



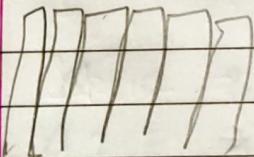
iii) Left Skew



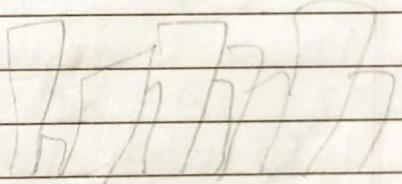
iv) Right Skew



v) Uniform



vi) No pattern



Page No.	
Date	14 P 24

★ ANOVA

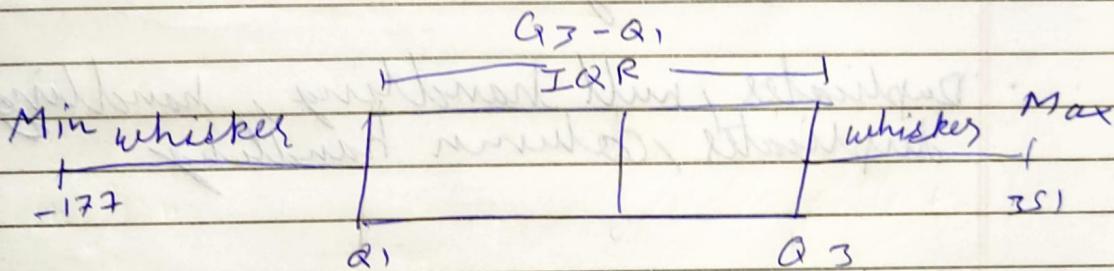
★ Boxplot example

1	6	-20
2	10	30
3	11	33
4	31	35
5	35	37
6	50	40
7	65	45
8	75	50
9	100	100
10	105	300
11	201	500
12	700	700
13	1000	

$$\text{IQR} = 160$$

$$\text{Min} = -125$$

$$\text{Max} = 449$$



$$Q_1 = 25\% = 21$$

$$\text{Min} = Q_1 - 1.5 \times \text{IQR}$$

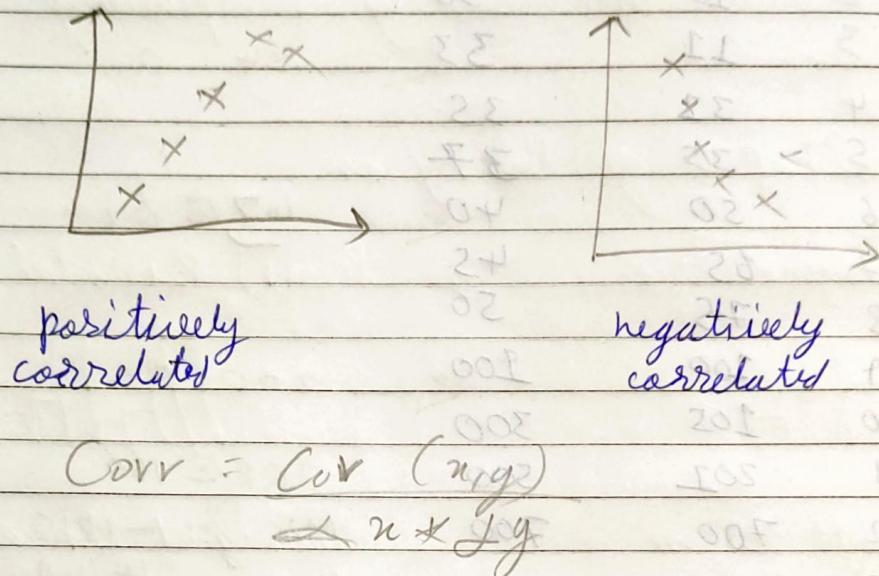
$$Q_2 = 50\% = 65$$

$$Q_3 = 75\% = 153$$

$$\text{Max} = Q_3 + 1.5 \times \text{IQR}$$

$$\text{IQR} = 132$$

★ Correlation



★ Data Cleaning

- Duplicates, null handling, handling duplicates, column handling

★ Pandas Profiling

! pip install pandas_profiling ↴

from pandas_profiling import ProfileReport

prof = Profile Report (df)

prof.to_file (output_file = 'output.html')

★ Confusion Matrix

		Predicted	
		0	1
Actual	0	FP	TN
	1	TP	FN

Actual Predicted
 Same \rightarrow True 1 \rightarrow Positive
 Differ \rightarrow False 0 \rightarrow Negative

		Predicted
Actual	0	1
0	Y	Y
1		

$$P = 10 \quad N = 0$$

0	1	6	0
1	1	4	0
1	1		
0	1		

$$TPR = \frac{TP}{P} = 0.6 \text{ (Recall)}$$

$$TNR = \frac{TN}{N} = 0$$

$$FPR = \frac{FP}{P} = 0.4$$

$$FNR = \frac{FN}{N} = 0$$

$$\text{Accuracy} = \frac{TP + TN}{T} = \frac{6+0}{10} = 0.6$$

$$\text{Precision} = \frac{TP}{TP + FN} = \frac{6}{6+0} =$$

Q

y	\hat{y}			
0	0	A	P	1
1	1	1	3	3
0	0	1	1+1+1	1+1+1
1	0	0	1	3
1	0	0	1	1+1+1
1	0			
0	0			
1	1			
0	1			
1	1			

$$TPR = \frac{TP}{P} 75\% \text{ (Recall)}$$

$$TNR = \frac{TN}{N} 50\%.$$

$$FPR = \frac{FP}{P} 25\%.$$

$$FNR = \frac{FN}{N} 50\%.$$

$$T=10$$

$$\text{Accuracy} = \frac{TP + TN}{T} = 60\%.$$

$$\text{Precision} = \frac{TP}{TP + FN}$$

Q	y	\hat{y}	$\textcircled{1}$	$\textcircled{2}$	$\textcircled{3}$
0	0	$A \setminus P$	0	1	$\frac{1}{2}$
1	1	$3 - 0$	1	2	0
2	2	$4 - 1$	0	$1+1+1+1$	0
2	2	$3 - 2$	0	1	$1+1$
1	1				
0	1				
1	1				
2	1				
0	1				
1	1				

$Q \quad y \quad \hat{y}$

1 1
1 0
1 0
1 1
0 1
0 0
0 1
0 0
1 0
0 0

		P		N	
		TP	FP	FN	TN
		1	2	2	3
		0	1	2	3

$$\text{Recall} = TPR = \frac{TP}{\text{Actual P}} = \frac{2}{5} = \frac{TP}{TP+FN} = 40\%$$

$$TNR = \frac{TN}{\text{Actual N}} = \frac{3}{5} = \frac{TN}{TN+FP} = 60\%$$

$$FPR = \frac{FP}{\text{Actual P}} = \frac{2}{5} = \frac{FP}{FP+TN} = 40\%$$

$$FNR = \frac{FN}{\text{Actual N}} = \frac{3}{5} = \frac{FN}{FN+TP} = 60\%$$

$$\text{Precision} = \frac{TP}{\text{predicted P}} = \frac{TP}{TP+FP} = \frac{2}{4} = \frac{1}{2} = 50\%$$

$$\text{Accuracy} = \frac{TP+TN}{T} = \frac{5}{10} = \frac{1}{2} = 50\%$$

Q

$$\begin{array}{c} X \\ \hline 123 \\ 567 \end{array} \Rightarrow \begin{array}{c} X \\ \hline 1-2-3 \\ 5-6-7 \end{array}$$

~~Q1~~ Hypothesis Testing - Unit 1

Q1 Population of Gandhinagar is more than 1Cr

H_0 : population of gandhinagar is greater than ~~<~~ 1 Cr

H_1 : population of gandhinagar is less than 1 Cr

Q2 Salary of Ram and Shyam is not same

H_0 : salary of Ram & shyan are not same

H_1 : salary of Ram & Shyam are same

Q.3 Value of Ared is more than 100 Km²

H_0 : Value of area is greater 100 Km²

H_1 : Value of area is less than 100 Km²

Q $H_0: \bar{S} = S$

$H_a: \bar{S} \neq S$

Q $H_0: Y = 10$

$H_a: Y > 10 \text{ or } Y < 10$

Q

$H_0: A \neq B$

$H_a: A = B$

13/09/24

Significant value (α)

$$CI = 1 - \alpha$$

(I \Rightarrow SY.)

$$\alpha \Rightarrow S. O\% = 0.05$$

Ex Ex $H_{ia}: \bar{U}_{ia} = 100$

$$n = 30$$

$$x_1 = 50 \quad n_1 = 30$$

$$x_2 = 60 \quad n_2 = 30$$

$$\bar{x} = 100$$

$$G_1 = 20, G_2 = 50$$

$$G = 20$$

degree of freedom
= $n - 1$

Page No.

Date

H₀: similar or equal

for
one sample

t-value =

$$\frac{\bar{x} - \mu}{\text{standard error}}$$

$$= \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$\textcircled{1} \quad \frac{110 - 100}{\frac{20}{\sqrt{30}}}$$

(t-value)
110

$$= \frac{10}{20} \sqrt{30}$$

$$= 2.73$$

$$\textcircled{2} \quad \frac{95 - 100}{\frac{20}{\sqrt{30}}}$$

(t-value)
95

$$= -1.36$$

for
Two Sample
Independent

t-value =

$$\frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$\text{DOF} = n_1 + n_2 - 2$$

for
Two Sample
dependent

t-value =

$$\frac{\bar{x}_c - 0}{\frac{s}{\sqrt{n}}}$$

$$\text{DOF} = n - 1$$

C.I. 90% 95% 99% 99.9%

1 tail 0.12. 0.5 0.01

2 tail 0.05% 0.025 0.005

	True	False	
Rejected	X	✓	
Not Rejected (Accepted)	✓	X	

17/1/24

Type I error : when H_0 is true we may reject it

$$P(\text{Reject } H_0 \text{ when it is true}) = P(\text{Reject } \frac{H_0}{H_1}) = \alpha$$

Type II error : when H_0 is false we may accept it

$$P(\text{Accept } \frac{H_0}{H_1}) = \beta$$

$$Q2 \bar{X} = 5g$$

$$\bar{X} = 5.05g$$

$$\sigma = 0.1$$

$$n = 12$$

$$H_0: \mu = 5 \checkmark$$

$$H_1: \mu > 5 \times$$

$$\text{Sol } t_c = \frac{\bar{X} - \bar{U}_0}{S / \sqrt{n}}$$

$$t_c = \frac{5.05 - 5}{0.1 / \sqrt{12}}$$

$$= \frac{\pm 0.05}{0.1 / 2\sqrt{3}}$$

$$= \frac{\pm 0.05 \times 2\sqrt{3}}{0.1}$$

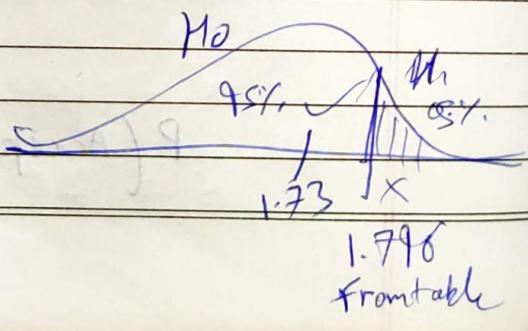
$$= \pm 0.5 \times 2\sqrt{3}$$

$$= \pm \frac{8 \times 2\sqrt{3}}{10}$$

$$= \pm \sqrt{3}$$

$$t_c = \pm 1.73$$

Accept H_0



$$\text{Q2 } \mu = 42 \text{ min}$$

$n = 12$ journeys

$$\bar{x} = 50 \text{ min}$$

$$\sigma = 15 \text{ min}$$

$$CI = 5\%$$

$$\text{Sol } t_c = \frac{50 - 42}{15/\sqrt{12}}$$

$$= \frac{8}{15/2\sqrt{3}}$$

$$= \frac{8}{15} \times 2\sqrt{3}$$

$$= 1.84$$

Reject H_0

Failed to re

$$Q.3 \quad n = 25$$

$$U_0 = 2.45 \quad H_0: M = 2.45$$

$$\bar{X} = 2.65 \quad H_1: M > 2.45$$

$$S = 0.35$$

$$CI = 99\%$$

Sol

$$t_C = \frac{2.65 - 2.45}{0.35 / \sqrt{25}}$$

$$= \frac{0.2 \times 5}{0.35}$$

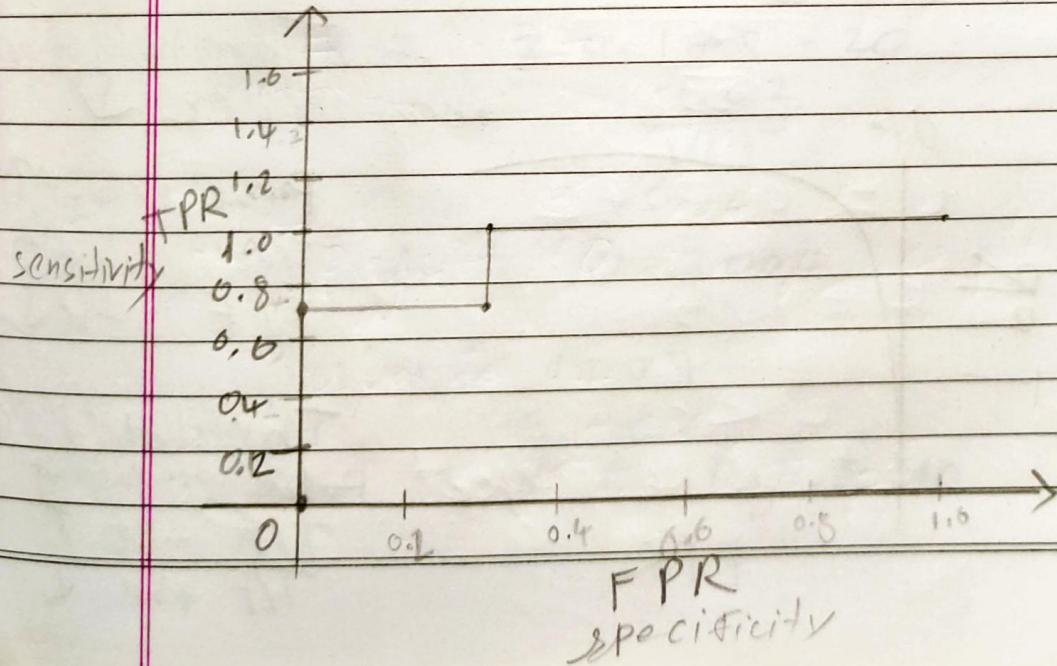
$$= \frac{2 \times 8}{3.5 \times 2}$$

$$t_C = 2.85$$

Score	Prediction	"	"	"	"	"	"	Y
(1)	(0.7)	(0.6)	(0.5)	(0) Y				
0.25	0	0	0	0	1	1	0	0
0.45	0	0	0	0	1	1	0	0
0.55	0	0	0	1	1	1	1	1
0.67	0	0	1	1	1	1	0	0
0.82	0	1	1	1	1	1	1	1
0.95	0	1	1	1	1	1	1	1
FP	0	0	1	1	3			
TP	0	2	2	3	3			
TPR	0	0	$\frac{1}{1+2} = \frac{1}{3}$	$\frac{1}{1+2} = \frac{1}{3}$	1			
FPR	0	$\frac{2+1}{2+1} = \frac{2}{3}$	$\frac{2}{2+1} = \frac{2}{3}$	$\frac{3}{3+0} = \frac{3}{3}$	1			

$$FPR = \frac{FP}{FP + TN} = \text{Actual } 0$$

$$TPR = \frac{TP}{TP + FN} = \text{Actual } 1$$



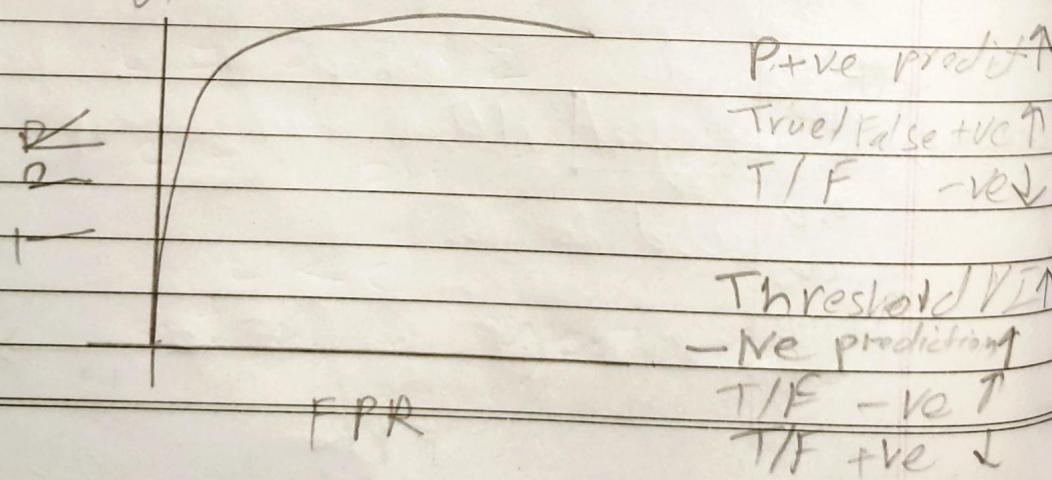
★ ROC

The Receiver operating Characteristics

$$97 = \frac{99}{97 + 98}$$

$$97 = \frac{99}{97 + 98}$$

Threshold value
dec ↘



Ex 2

$$\text{Sol } \bar{x} = 20.175$$

$$\mu = 20$$

$$\sigma = 3.02$$

$$n = 12$$

$$H_0: \mu = 20$$

$$H_1: \mu > 20$$

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$t = \frac{20.175 - 20}{\frac{3.02}{\sqrt{12}}}$$

$$t = 0.2007$$

$$c = 1.79$$

$$t < c$$

\therefore Fail to Reject H_0

Ex 1

Sol

$$\bar{x} = 21$$

$$\mu = 20$$

$$n = 25$$

$$\sigma = 7$$

$$H_0 : \mu = 20$$

$$H_1 : \mu > 20$$

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$t = \frac{21 - 20}{\frac{7}{\sqrt{25}}}$$

$$t = 0.7142$$

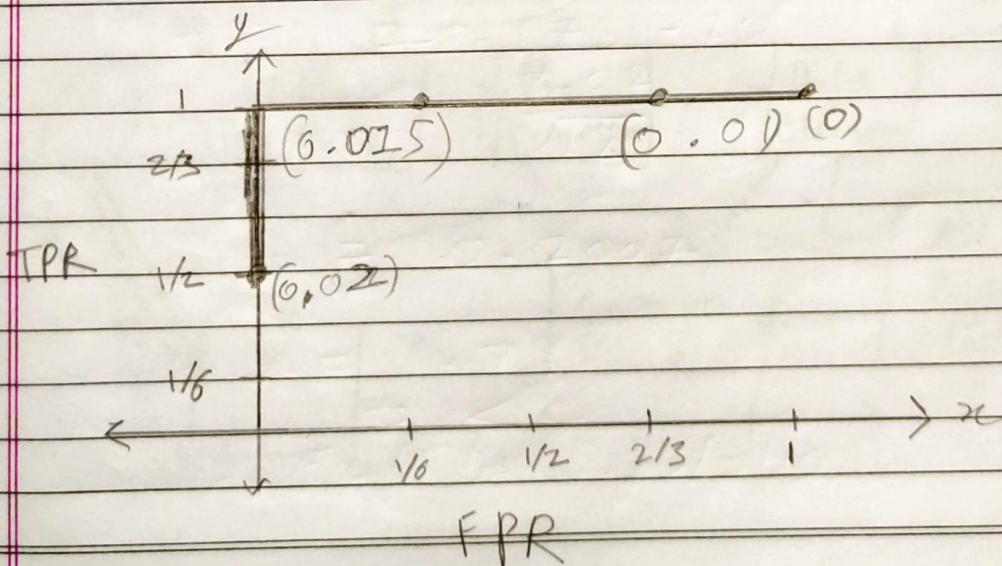
$$c = 1.711$$

$$t < c$$

\therefore Fail to reject H_0

Ex ROC

Instance	Cut-off	0.020	0.015	0.01	0	Actual
1	0.008	0	0	0	1	0
2	0.011	0	0	1	1	0
3	0.021	1	1	1	1	1
4	0.009	0	0	0	1	0
5	0.014	0	0	1	1	0
6	0.015	0	1	1	1	0
7	0.012	0	0	1	1	0
8	0.015	0	1	1	1	1
TP		21	2	2	2	
TN		6	6	2	0	
FP		0	1	4	6	
FN		1	0	0	0	
TPR		1/2	1	1	1	
FPR		0	1/6	2/3	1	



Quiz

- 1) FP
- 2) Both TP + Ve \downarrow
- 3) Both T/F - ve T
- 4) A
- 5) F + ve \downarrow
- 6) No
- 7) Recall \downarrow if precision \uparrow
- 8) False

* EDA (Exploratory Data Analysis)

Q Is EDA & DA the same

Sol DA - descriptive, diagnostic, predictive and prescriptive

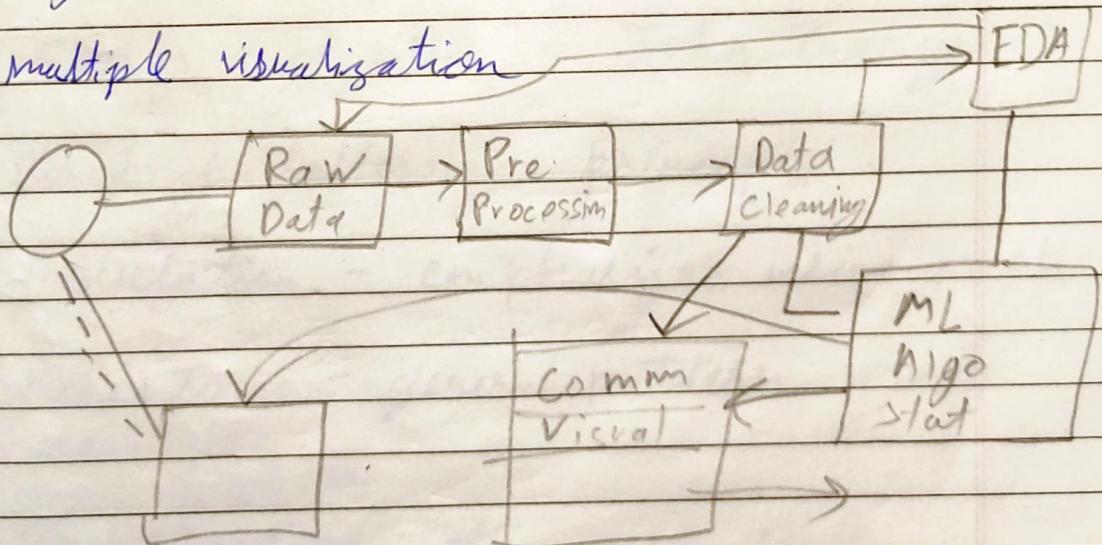
EDA - synonymous with descriptive analysis

where one explores hidden relationships and patterns in the data available data

Q How one should write an EDA report

Sol Thorough explanation of dataset's variables

- their correlation
- any preprocessing perf on ds for ml algo
- multiple visualization



Page No.

Date

Types of EDA

EDA

Univariate

graphical non graphical

Multivariate

graphical non graphical

* Tools of EDA

1. Python

2. R

3. MATLAB

* Philosophy of EDA

Father of EDA - John Tukey

Based on following principles

1. Relation - emphasizes using graph
2. Resistance - general pattern which can
3. Reexpression: to new scale

Page No.			
Date			

① Curiosity and

Unit 3

* Basic ML Algos

1) Association Rule mining

- unsupervised learning technique
- tries to find rel^h between / among the variables of dataset
- relation / association
- Ex - Market Basket analysis, Web usage mining, etc
- working
 - If and Else statement concept

• Measure of associations between multiple sets of data

i) Support

$$\text{supp}(X) = \frac{\text{Freq}(X)}{T}$$

ii) Confidence

$$\text{confidence} = \frac{\text{Freq}(X, Y)}{\text{Freq}(X)}$$

iii) Lift

$$\text{Lift} = \frac{\text{Supp}(X, Y)}{\text{Supp}(X) * \text{Supp}(Y)} \quad \left(\because \text{Supp}(X, Y) = \frac{\text{Freq}(X, Y)}{T} \right)$$

'Q Ex

$B \rightarrow C$

i) Supp = $\frac{6}{9}$

ii) Confidence = $\frac{4}{6}$

iii) Lift = $\frac{4/9}{6/9 \times 6/9} \leq \frac{4}{36/9} = 1$

ID	Itemset
T1	A, B
T2	B, D
T3	B, C
T4	A, B, D
T5	A, C
T6	B, C
T7	A, C
T8	A, B, C, E
T9	A, B, C

- If lift = 1

- If lift > 1

- If lift < 1

* Bo MLA

- Machine Learning (ML)
- subset of AI
- involves training of computers to learn from data without being programmed
- Identify patterns, make predictions
- performs tasks that usually require human intervention

* Key Components of ML Algo

- 1 Data - raw material . unstructured (ex - csv)
- 2 Features - relevant characteristics
- 3 Model -
- 4 Algorithm -

Types of ML Algs

1. Supervised Learning

- Regression
 - predicting a continuous numerical value
(ex - house price)
 - i) Linear
 - ii) Logistic
 - iii) Decision Tree
 - iv) Random Forest
- Classification
 -
 - i) Decision Tree
 - ii) Random Forest
 - iii) Neural Networks

2. Unsupervised Learning

- Clustering - grp data pts into similar clusters
 - i) K-means clustering
 - ii) hierarchical clustering
 - iii) DB SCAN

3. Reinforcement Learning - learn through trial and error, interacting with an environment and receiving rewards or penalties.

Page No.

Date

Machine Learning

Supervised
Learning

Unsupervised
Learning

Reinforcement
Learning

Classification Regression

Clustering

* Types of Regression Algo's

- ① Linear
- ② Random Forests
- ③

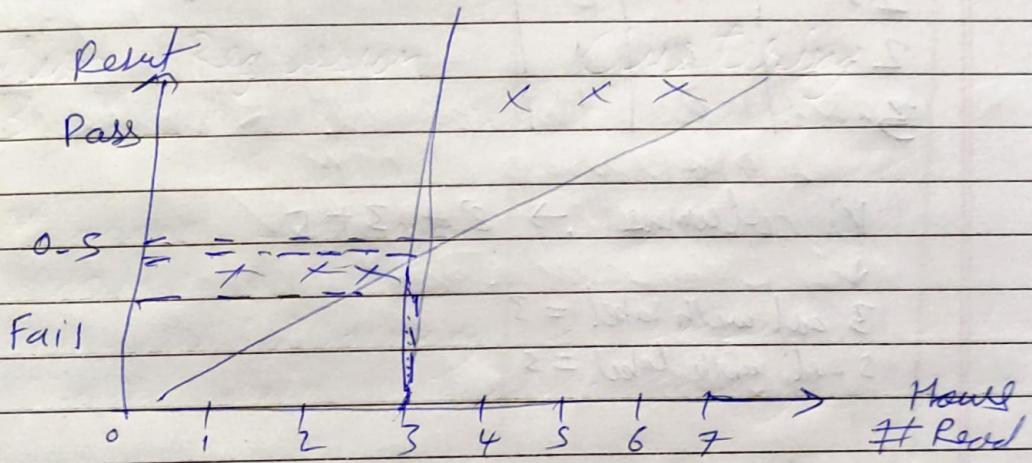
Usage

- 1
- 2
- 3

① Linear Regression

Mean Square Errors

# hours	Result
Read	Fail
1	Fail
2	Fail
3	Pass
4	Pass
5	Pass
6	Pass
7	Pass



$$y \text{ or } g = h(x) = m x + c \\ = \beta_0 + \beta_1 x$$

$$g(x) = \frac{1}{1 + e^{-x}} \quad \text{sigmoid function}$$

logistic Regression (x)

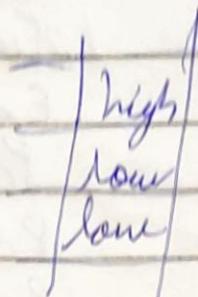
$$= \frac{1}{1 + e^{\beta_0 + \beta_1 x}}$$

* Encoding Categorical Data

① Label Encoding

Low = 1

High = 2



② One-hot Encoding (N categories adding new n columns)

	Female	Male	
1.	1	0	
2.	0	1	
3.	0	1	

old

new

10 columns $\rightarrow 8 + 3 + 5$



3 col with label = 3

5 col with label = 5

id name grade			\Rightarrow	id name A B		
1	X	A		1	X	1 0
2	X	B	\rightarrow	2	X	0 1
3	Z	A		3	Z	1 0

* Model Evaluation

- | | |
|---------------------|----------------------------|
| ① R-squared Value | close fit
(nearer to 1) |
| ② Mean-Square Error | lower MSE |

a) Example 1

	old - x = new (w/o op)
1) 10 col in which 2 col header	8 + 5 + 5
5 Sub category 10 old	18 (i)
2 columns with 3 category	17 + 3 + 3 + 10
1 col with 10 columns	new (w/o op)

30/9/24

* Classification Techniques

1. logistic regression
2. K nearest neighbors (Knn) (K=odd)

② KNN

If K is too low

If K is very high

hyperparameters = K - no. of nearest neighbors

- * ~~how~~ how to decide K value?
 - i) Bootstrap
 - ii) Square root of N
 - iii) cross validation

* Application

- Null Handling \rightarrow simple Imputer (preprocessing)
- Classification

* Advantage

- Easy to implement
- less hyperparameter
- Adjustable

* Disadvantages

- Scale
- Curse of Dimensionality
- overfitting / Underfitting

2/10/24

(3)

Naive Bayes

→ Theorem - Bayes theorem

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

i) Prior

ii) Conditional

* Disadvantages

- i) Basic assumption like independent event are not real life/world data
- ii) Zero frequency for a new keyword, gives zero probability as outcome

* Advantages

- easy & efficient
- effective with large no. of features
- can handle high dimensional data
- less complex
- scales well
- overfitting

* Applications

- Spam Email Filtering
- Text Classification

Example

$w_1 \quad w_2 \quad w_3 \quad \text{Spam}$

Sale offer Fast Y

Sale Click Quick Y

Test Offer Money N

Predict $w_1 = \text{Sale}$

$w_2 = \text{Money}$

$w_3 = \text{Money}$

$w_1 = \text{Sale}$

$w_2 = \text{Offer}$

$w_3 = \text{Fast}$

① Prior probability

$$P(Y) = 2/3$$

$$P(N) = 1/3$$

② Find Conditional Probability

w_1	Y	N
Sale	1/2	0/1
Test	0/2	1/1

w_2	Y	N
Offer	1/2	0/1
Click	1/2	0/1

w_3	Y	N
Fast	1/2	0/1
Quirk	1/2	0/1
Money	0/2	1/1

$$(3) P(Y/w_1 = \text{Sale}, w_2 = \text{Money}, w_3 = \text{Money})$$

$$= P(Y) * P(w_1 = \text{Sale}) * P(w_2 = \text{Money}) * P(w_3 = \text{Money})$$

$$= \frac{2}{3} * 1 * 0 * 0 = 0$$

$$P(N/w_1 = \text{Sale}, w_2 = \text{Money}, w_3 = \text{Money})$$

$$= \frac{1}{3} * 0 * 0 * 1$$

$$= 0$$

Example 2

$$P(Y) = 5/10$$

$$P(N) = 5/10$$

	Y	N
Red	<u>$\frac{3}{5}$</u>	<u>$\frac{2}{5}$</u>
SUV	<u>$\frac{1}{5}$</u>	<u>$\frac{3}{5}$</u>
domestic	<u>$\frac{2}{5}$</u>	<u>$\frac{3}{5}$</u>

$$P\left(\frac{Y}{(C=\text{Red}, T=\text{SUV}, O=\text{domestic})}\right) = \frac{5}{10} \times \frac{3}{5} \times \frac{2}{5} = \frac{3}{125}$$

$$P\left(\frac{N}{(C=\text{Red}, T=\text{SUV}, O=\text{domestic})}\right) = \frac{5}{10} \times \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} = \frac{9}{125}$$

i)* Entropy

$$E = - \sum_{i \in C} P_i \log_2 P_i \quad \frac{\log_{10}^x}{\log_2^x} = \log_2^x$$

Ex

Play TT

Yes

Yes

No

No

No

No

$$C = \{ \text{Yes}, \text{No} \}$$

$$P(\text{Yes}) = \frac{2}{3} = \frac{1}{3} \log_2 \frac{1}{3} = \frac{1}{3} (-0.584) \\ = -0.584$$

$$P(\text{No}) = \frac{1}{3} = \frac{2}{3} \log_2 \frac{2}{3} = \frac{2}{3} (-0.176) \\ = -0.176$$

$$E = - \sum_{i \in C} P_i \log_2 P_i$$

$$= - \left[(-0.584) \times \frac{1}{3} + \frac{2}{3} (-0.176) \right]$$

$$= -[-0.528 - 0.384]$$

$$= -[-0.917]$$

$$= 0.917$$

ii) * Gini Impurity

$$= 1 - \sum_{c \in C} P_c^2$$

Ex $C = \{\text{Yes, No, Maybe}\}$

$$P(\text{Yes}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{No}) = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{Maybe}) = \frac{1}{6}$$

$$GI = 1 - \left[\left(\frac{1}{3} \right)^2 + \left(\frac{1}{2} \right)^2 + \left(\frac{1}{6} \right)^2 \right]$$

$$= 1 - \frac{1}{9} + \frac{1}{4} + \frac{1}{36} = \frac{36}{36} - \frac{4}{36} + \frac{9}{36} + \frac{1}{36} = \frac{42}{36} = \frac{7}{6}$$

$$= 1 - \frac{4 + 9 + 1}{36}$$

$$= 1 - \left[\frac{14}{36} \right] = 1 - \frac{14}{36} = \frac{22}{36} = \frac{11}{18}$$

$$= \frac{11}{18}$$

$$E =$$

$$FIP = \frac{11}{18}$$

2/10/24

* Decision Tree

iii) * Information Gain

A - attribute

S - collection of dataset

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{value}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\begin{aligned}
 &= (0.64)(-0.63) + (-0.35)(-1.48) \\
 &= +0.403 + 0.51 \\
 &= +0.913 \rightarrow \approx 0.940 \\
 &\approx +0.940
 \end{aligned}$$

18/11/24

* Recent Trends in various data collection
and analysis techniques

* Unit 5

2) Data Visualization

art & science of communicating info
visually

* Types of Visualization

1. Histogram / Line chart

2 Bar chart

3 Piechart

4 Bar chart

5 Scatter plot

6 Heatmap

7 Map

8 Network diagram

9 Tree

10 Boxplot

* Technologies

D3.js

Plotly

Matplotlib

* Tools

- Tableau
- Power BI
- QlikView
- Google Data Studio
- Excel

* Techniques

- i) Colour Coding
- ii) Shading
- iii) Symbol size
- iv) Line thickness
- v) Animation
- vi) Interactive Elements

18/11/24

Page No.	
Date	

* Unit 3

Feature Engineering

→ Techniques of Feature Selection

- i)
- ii)
- iii)

* Unit 4

* Common Distance Measures

i) Euclidean Distance

$$\sqrt{x_1^2 - x_2^2}$$

ii) Manhattan Distance

$$= |P_1 - P_2|$$

iii) Minkowski

$$= ((x_1 - x_2)^x)^{1/x}$$

Q In Sklearn, which distance technique they are using?

Ans Euclidean

* Different Clustering ~~techniques~~ techniques

- 1.) Hierarchical
- 2.) Partitional
- 3.) Fuzzy
- 4.) Model-Based

* Cluster distance measures

- Single linkage

$$D(c_1, c_2) = \min_{\substack{x_i \in c_1 \\ x_j \in c_2}} D(x_i, x_j)$$

- Complete linkage

$$D(c_1, c_2) = \max D(x_i, x_j)$$

- Average linkage

$$D(c_1; c_2) = \frac{1}{|c_1|} \frac{1}{|c_2|} \sum \sum D(x_i, x_j)$$

- ~~Centroids~~ Centroids

$$D(c_1, c_2) = D\left(\left(\frac{1}{|c_1|} \sum_{x_i \in c_1} \bar{x}_i\right), \left(\frac{1}{|c_2|} \sum_{x_i \in c_2} \bar{x}_i\right)\right)$$

- Ward's Method

$$TD_{c_1 \cup c_2} = \sum_{x_i \in c_1 \cup c_2} D(x_i, M_{c_1 \cup c_2})^2$$

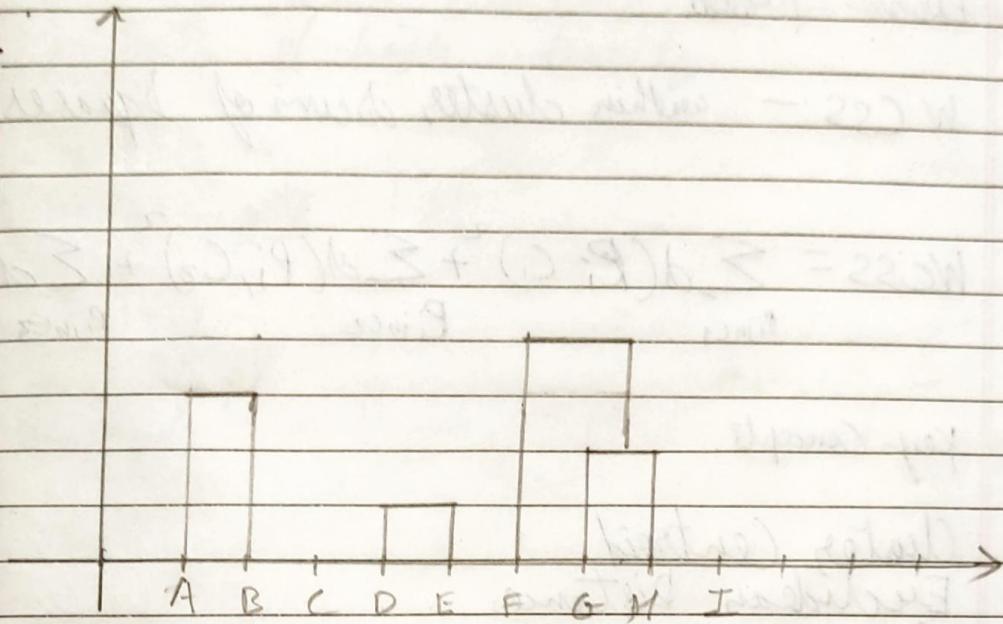
(1) P-L

(2) G-H

(3) A-B

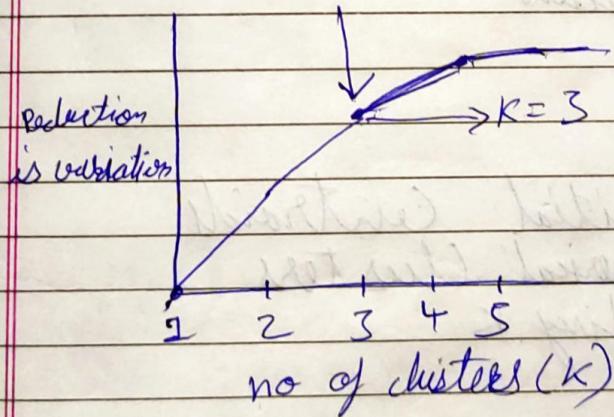
(4) G-H - F

★ Dendrogram



★ Kmeans - elbow clustering

elbow plot



stat Question with Tash Stat

YT

2/1/24

Page No.	
Date	

* Kmean's Clustering

Elbow Method

WCSS - within clusters sum of squares

$$WCSS = \sum_{P_i \in C_1} d(P_i, C_1)^2 + \sum_{P_i \in C_2} d(P_i, C_2)^2 + \sum_{P_i \in C_3} d(P_i, C_3)^2$$

→ Key Concepts

- Cluster Centroid
- Euclidean Distance
- Convergence

Strength

- Simple & Efficient
- Scalable

Weakness

- Sensitive to Initial Centroids
- Assumes Spherical Clusters
- Requires Specifying K

22/11/24

* DB Scan - density based spatial clustering
of applications with Noise

This algo defines clusters as continuous
regions of high density

Key concepts

- Core pt
- Border pt
- Noise pt

min = 4

ϵ - per cap

* Strength

- Handles clusters of Arbitrary Shape
- Noise Tolerance
- Does not require the no of clusters

* Weakness

- Sensitivity of Parameters
- Clustering Density Variations

25/11/24

Page No.

Date

* Choosing the right clustering approach:

- No of clusters
- cluster shape
- presence of noise & outliers
- computational complexity
- Interpretability

- * • Distance Matrix
- * • Feature Selection, Scaling

(A) Hierarchical Clustering

- Strength - not defining no of clusters
any x shape, visualize & interpret
- weakness - combining all pts into one cluster
less columns / small data set

(B) K means clustering

- S - large dataset
hyper parameters
- W - simple & efficient
defining no of clusters
non spherical
sensitive to initial centroid loc

③ DB Scan

- S - non spherical, ~~large~~ data sets, not define no of clusters, handle noise & outliers
- W - sensitive, for min pts & ϵ (distance) ^{radius}
density
computationally expensive for large dataset

④ Fuzzy Clustering

- S - can handle overlapping clustering
more realistic representation of data
- W - can be computationally expensive
interpretation of results can be challenging

Model-Based clustering (Gaussian Mixture Model)

- S - handle → model complex data set
overlapping clusters
estimate no of clusters

- W - sensitive to

* Clustering Tendency

Y?

- avoids futile clustering
- optimizes algo selection

how?

- Visual Inspection - scatter & density plots
- State Test - Hopkins test
- model-based approaches

* Specific application in Data Mining

Customer segmentation

Anomaly detection

Document clustering