

FDS & ML

6/8/24

Data Science Process:-

- 1) Define the problem Statement.
 - 2) Data acquisition.
 - 3) Data preparation - Apply EDA [Exploratory Data Analysis]
 - 4) Model planning
Data finalization
 - 5) Visual Communication. {Visualization}
 - 6) Model Deployment
- To Instance Based Future Selection
Model-based.

- Titanic data

Know your data (Commands)

df - dataframe

- Data size - df.shape {row & column. of data}
- Data look like - df.head()
- Detail of data - df.info()
- Calculate null value - df.isnull().sum()
- Look mathematically - df.describe()
- Checking duplicates - df.duplicated().sum()
- Checking correlation b/w columns - df.corr()

Type of Data

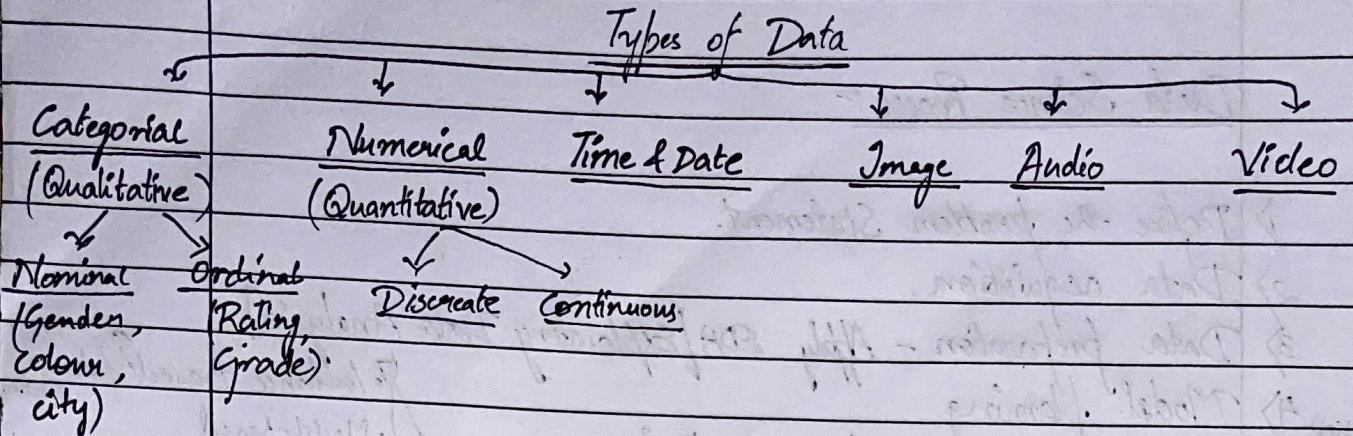
→ Categorical
→ Numerical.

- 1) Categorical Data → EDA univariate {Exploratory Data Analysis}
- 2) Numerical Data.

Categorical Data

- ↳ Pie chart. \Rightarrow df['column name'].Value_counts().plot(kind='pie')
- ↳ Count chart \Rightarrow sns.countplot(df['column name'])

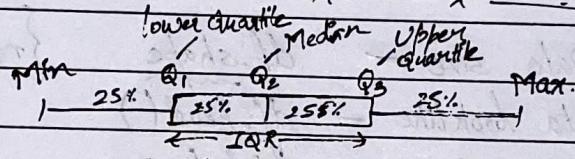
Types of Data



- PassengerId - Nominal
- pClass - Ordinal
- SibSp - Discrete
- Fare - continuous
- Survival - nominal
- name - nominal
- Parch - discrete
- Cabin - nominal
- sex - nominal
- Ticket - nominal
- embarked - nominal
- age - continuous

Numerical univariate

→ Box plot ⇒



→ Histogram ⇒

Inter Quartile Range

plot
sns.hist(df['Age'])

→ Distplot = histogram + PDA

(Probability density....)

sns.distplot(df['Age'])

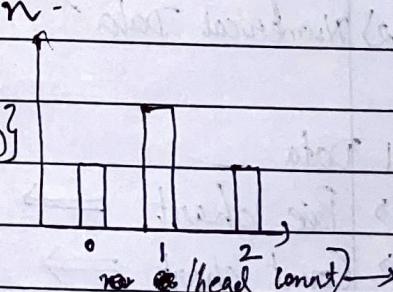
Probability - It is the measure of likelihood of an event occurring.

1) Random Variable

2) Probability distribution - description of possible value of each random variable.

Case: Prob of getting Head

Sample space: { (H,H) (H,T) (T,H) (T,T) }



3) Conditional probability - If first event occurs, then prob of second event occurring

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

4) Baye's theorem.

Statistics : Science of collecting, analysing, interpreting, presenting & organizing data.

Types:-

1) Descriptive statistic - Set of data is defined.

2) Inferential statistic - Sample space is not defined

Hypothesis - statistic also helped in hypothesis.

• Correlation - relationship b/w 2 entities.

• Regression - Model defining relationship b/w dependent variable and one or more independent variable.

Bivariate / Multivariate

→ heat map

sns. heatmap

pd. crosstab

sns. distplot(titanic['Survived'], hist=False)

→ Scatterplot

sns. scatterplot(x="Survived", y="Age")

Commands :- (Titanic dataset)

1) Age > 29

df[df['Age'] > 29]

2) Drop a column

df1 = df.drop(columns=['Age'])

df2 = df.drop(3) # index no./column no.

3) Delete rows which have null value

df.dropna(inplace=True)

4) Change datatype of column

df['Amount'] = df['Amount'].astype('int')

5) Update value of a column

df[2, 'Age'] = 23 → column value

6) To find null values.

df.isnull().sum()

inplace=True → will drop column from the current dataframe & store in it.

Tendency (Measures of the center)

- The mean (the average value)
- The median (the mid point value)
- The mode (the most common value)

Spread (Measures of Variability)

min & max.

variance

Moment

$$\text{Variance} \rightarrow \sum (x - \bar{x})^2 / (N-1)$$

$$\text{Standard Deviation} \rightarrow \sqrt{\text{Variance}}$$

$$\text{Mean}(x) = \frac{\sum_{i=1}^N x_i}{N}$$

m, 3, 8, 4

Moment

momentum

Mean

Central tendency

~~standard deviation~~

Variance
(Volatility)

Dispersion

$$= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

Skewness

Symmetry (+ve or -ve)

$$= \frac{1}{N} \sum_{i=1}^N \frac{(x_i - \bar{x})^3}{\sigma^3}$$

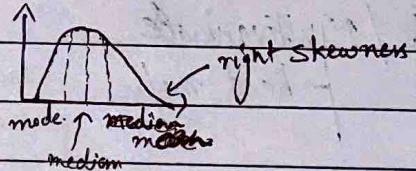
Kurtosis

Shape (Tall or flat)

$$= \frac{1}{N} \sum_{i=1}^N \frac{(x_i - \bar{x})^4}{\sigma^4}$$

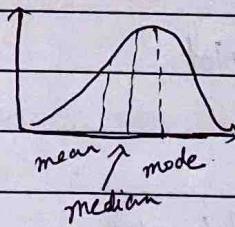
Types of Skewness

1) Right skewness
(+ve skew)



mode < median < mean

2) Left skewness
(-ve skew)



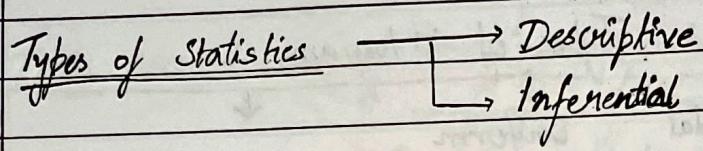
mean < median < mode

Types of Kurtosis

1) Mesokurtosis: An excess kurtosis of 0. Normal distributions are mesokurtotic.

2) Platykurtosis: Neg. excess kurtosis. Few outliers

3) Leptokurtosis: Positive excess kurtosis. Many outliers.



- (how the categorical value is related to the result)
- Analysis of Variance (ANOVA) - involves comparing means across multiple groups to determine if there are any significant differences. For example, comparing the mean height of individuals from different regions. Useful with categorical data.
 - ↳ within
 - ↳ between
 - Regression Analysis - Involves modelling the relationship b/w a dependent variable & one or more independent variables. For ex - predicting the sales of a product based on advertising expenditure.
 - Confidence intervals - involves estimating the range of values that a population parameter could take based on a sample of data. ex. estimating the population mean height
 - Chi-square Test - Test independence or association (correlation) b/w 2 diff categorical variables ex. Testing gender & occupation are indep. or not
 - Sampling techniques - Involves ensuring that the sample of data is representative of the population. ex. random sampling to select individuals from a population
 - Bayesian Statistics - Alternative approach to statistical inference that involves updating beliefs about the probability of an event based on new evidence. for ex:- updating the probability of a disease given a positive test results.
 - * Weighted mean - sum / weight of product.
 - * Trimmed mean

Measure of Dispersion -

Coefficient of variation = $(\text{Standard Deviation} / \text{mean}) \times 100\%$

Types of histogram

Symmetric

Bimodal

Uniform

Types of Bivariate.

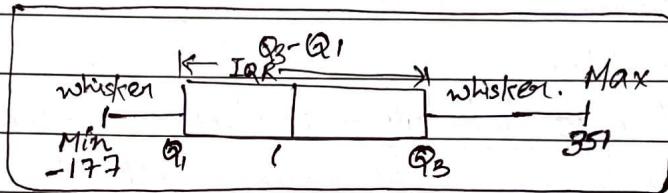
Contingency Table / Cross tab

- 1) Categorical - categorical
- 2) Numerical - numerical → Scatter plot
- 3) Categorical - Numerical → Contingence.

* Distribution Type

(Odd Index) Box Plot

Sorted	1	6
data	2	10
in ascending order	3	11
	4	31
	5	35
	6	50
	7	65
	8	75
	9	100
	10	105
	11	201
	12	700
	13	1000



$$\text{Min} = Q_1 - 1.5 \times \text{IQR} = -177$$

$$\text{Max} = Q_3 + 1.5 \times \text{IQR} = 357$$

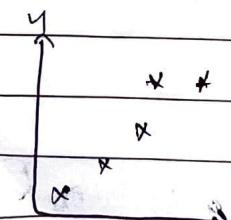
Outliers: 700, 1000.

$$\begin{aligned} Q_1 &= 21 \\ Q_2 &= 65 \\ \text{IQR} &= 153 - 21 \\ &= 132 \end{aligned}$$

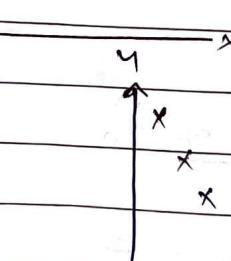
Even Index

1	20	Even Index	166
2	30		166 $\frac{153}{132}$
3	33	$Q_1 = 34$	Min = -215
4	35		
5	37		
6	40	$42 \rightarrow Q_2$	Max = 449
7	46		
8	50		
9	100	$Q_3 = 200$	Outlier = 500, 700
10	300		
11	800		
12	700		

Correlation:

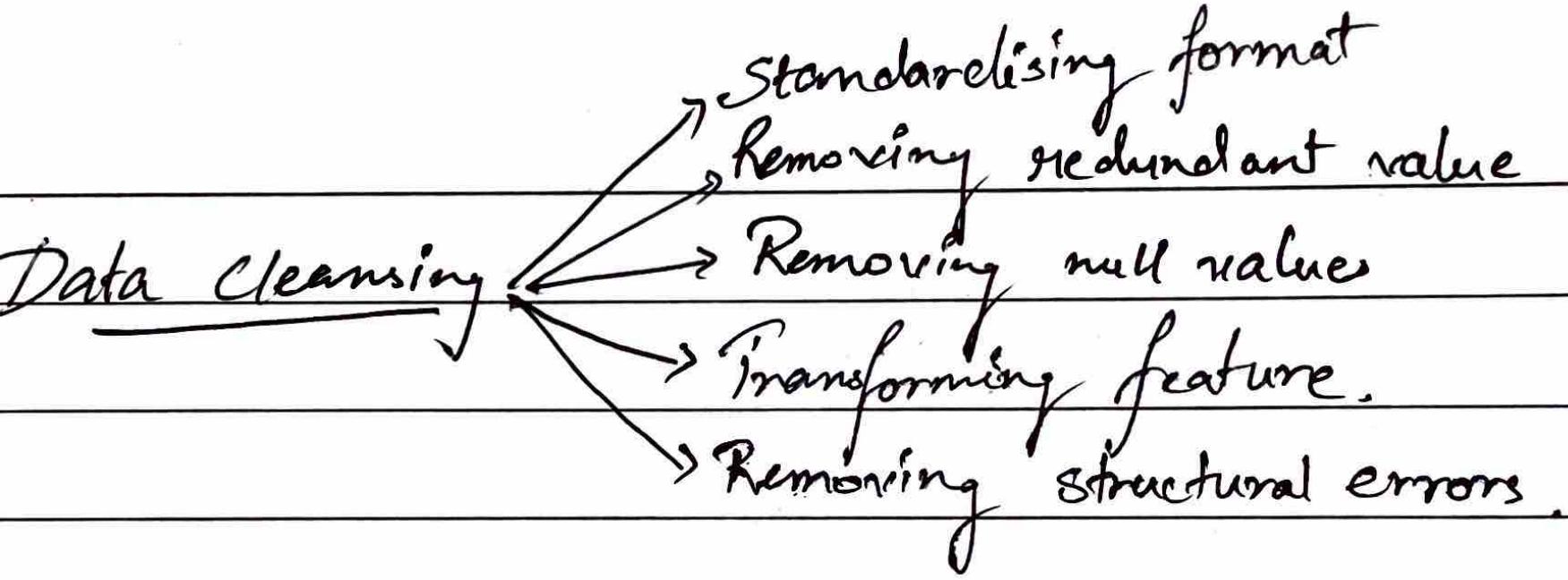


Positively correlated.



Negatively correlated

$$\frac{300}{267} = 1.13$$



Data Cleansing

- Handling Noisy Data →
- Binning → Clustering
 - Regression → Combined computer & human inspection.

Lossless & Lossy Compression

~~TA - Symbols~~

- Data science process
- univariate analysis.
- Types of distribution
- Types of data.
- Bivariate analysis.
- (uniform, normal etc)
- Basic command of python "know your data"
- Momentum(0,2,3,4)
- Importance of Data Preprocessing
- Steps in Data Preprocessing
- Python code
- Confusion matrix

Confusion Matrix (Mã mìn ko confusion hoi)

		Actual		Prediction:	Actual
		1	0	1 → Positive	Same → True
y	ŷ	1	TP	FN	0 → Negative
		0	FP	TN	Difference → False

$$P = TP + FN$$

$$N = FP + TN$$

$$\text{Accuracy} = \frac{TP + TN}{P + N}$$

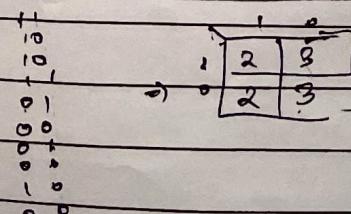
$$\boxed{\text{TPR} = \text{Recall}}$$

$\frac{TP}{N}$
 $N=6$

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

$$\begin{aligned} TP &= 3 \\ TN &= 3 \\ FP &= 3 \\ FN &= 1 \end{aligned}$$

		Prediction	
		0	1
Actual	0	1	2
	1	0	4
2	0	1	2



— / /

Define a function for following conversion.

x	x.
112233	11-22-33
45566	44-55-66

def con(v):

return x[:2]+ '-' + x[2:4]+ '-' + x[4:]

def -con(v):

return x[:2] + x[3:5] + x[6:]

11-22-33	112233
44-55-66	445566

Remove =>

all char
except number

11,223\$	11223
455667	455667