

(بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ)

(Data Science) Courses

(1)

Introduction

Definition :-

A Data Science is an interdisciplinary field that includes

- Statistics
- Artificial Intelligence
- Programming
- Domain Expertise

Question :- Differentiate

B/w Machine learning, Deep Learning and Data Science?

VS

ML vs Data Science

- | | | |
|----|---|---|
| ML | Focuses on Studying models and algorithms | → Works one steps from ML. |
| | data | → Solves Real-world Problems. |
| | Just we interested in Model in ML / DL | → Follows an End-to-End Process from problem formulation to presentation. |

Decision tree

What is Data Science?

End-To-End Data Science Process

Data Science cover the entire Journey from problem formulation insight delivery.

* Problem formulation

Defining the question

* Model Development

Building Predictive models

* Insight Communication

Sharing Result

④ you can Represent all insight

Assignment to Data Science Management.

• Data Science is End-to-End discipline To deliver insight from with the help of Stat and ML Models.

* Data Acquisition

Gathering Relevant Data

* Model deployment

Implementing in Production

Step 1:- (Role of domain Expertise)

* Problem formulation :-

Ex:- House prices Predicting Application

④ Input features ④ Expected output

Location, Size, Age house price

of house. (Selecting appropriate algorithms)

This insight of house To Prediction Prices.

Step 2:- * Data Acquisition (y/jpo)

Collect Relevant data for problems data Sources

* APIs

→ Programmatic

interface To Access data

* Data-base

→ Structured collection

of information

* Web-Scraping

→ Extracting

data from web-sites

* Public datasets:- Open Source Collection

Repository:- Kaggle, many university.

Step 3:- Data Preparation

- Understanding and cleaning data ^(EDA) Preprocessing.

④ Cleaning data

- handling missing values and Removing outliers

⑤ Transforming data

- Scaling, encoding categorical variables

⑥ Feature Engineering:-

- Creating new meaningful variables.

Step 4:- Data Analysis / Model Development:

- ① Selecting machine learning models "Evaluation" and "Fine Tuning" models Performance metrics

- Collection Ratio → True-Positive Accuracy
- Finding all Prediction(positive)

Step 5:- Model Development

(3)

Trained models accessible for real world use. **Methods:-**

- ⑥ Web-application || Embedded Systems.
- User friendly hardware implementation interface

⑦ APIs :- Service Integration.

Step 6:- Communicating insight:-

- Sharing these insight to (non-technical) Stakeholders (**CEO, Management**), all decisions maker.

* data visualization

- graphical representation of your insight

- Business Reports, formal - documentation

* Storytelling

- Crafting with narratives of data

Question:- Which types of skill are required for Data Science?

Essential Skills:-

(4) :-

1 Statistics

- Understanding

data distribution

4 Data Visualization

- Creating graphs

5 Communication Skills

- Finding clarity

Exploring.

2 Machine learning

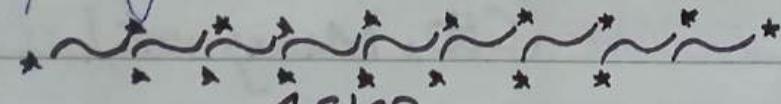
- Predicted models

3 Python - Progr

amming

- Implementing

Solution



1st step:-

(EDA) :- An other method use
in data understanding.

Exploring the data types:

structured data / unstructured Data

if structured types then we
have seen which types of data."

Attributes: Normal values, ordinary
types.

→ Data Types

→ Distribution of values

• Numerical, category • how data is spread

binary

→ Missing values and outliers :-

- data quality issues.

→ Central tendency Measure :-

- Mean, Mode, Median.

→ Dispersion Measure :-

- Variance, Stand deviation

→ Proximity Measure

Corelation similarly.

2nd Step:- Data Preprocessing:

1. Handling Missing values

- imputation or removal

2. Recovering Inconsistency

- fixing errors

3. Data Transformation

(0 to Range 0 to 10)

- Scaling normalization

(A house which has price is
"30" Rs. This is very inconsistency)

4. Data Reduction

- Feature Selection

Impact of Data Processing:-

(5)

- ⑥ Clean Structured data ensures effects model training & it helps:-

→ Improve model Accuracy :-

Better Predictions

→ Reduce Computational time :-

Faster Processing

→ Eliminate Misleading information :-
(Give More reliable Result.)

→ ~~Like ML, Statistics~~

are "developed." Use different Techniques in ML/AI Tools

(1) Regression :-

Predicted Continuous value

(2) Classification

Categorizing data

(4) Statistical Modeling :-

(3) Clustering :-

Identifying groups

Understanding Relationships.

(5) Self-Supervised Learning

Ex:- like house-price prediction (Score)

Using Regression Algorithms models.

Next Stage:- Step 3

Data Analysis :-

- ① The Core Stage when models

* Algorithms

It is Just Technique, Algorithm to see

like:- Linear Regression, DL interpretable model data and learn data

Polynomial Regression non-linear relationships

Random Forest Regressor Decision Tree

* Models

Algorithm to see

data and learn data make models.

Data Knowledge.

Evaluating Models:-

- It helps measure model performance using evaluation metrics

MSE	RMSE
Mean Squared Error	Root Mean Square Error
Average Squared diff b/w Predicted and actual values	Square root of MSE in same units as target

R^2 :- (R-Squared Value)

Proportion of variance explained
by models.

Model Development:-

- It is very important part of data analysis. So User of Stakeholder can used and share experience and experiment.

انٹرنیٹ سرورز میں اسکے اور سوال کو مرکزی۔

- Web Application
 - ⑥ Browser-based interface
 - Mobile Application
 - ⑥ On Smart phones and Tablets
 - APIs
 - ⑥ Application Programming interface
 - Embedded System
 - ⑥ hardware devices

Communication & Visualization

- ⑥ Data driven insight must be represented in graph. Used Tools

Matplotlib	Seaborn	Power BI/Tableau
python plotting library	Static data visual	Business Intelligence Tools

Types:-, Bar charts → Line charts
→ Scatter plots → Pie charts

Lectures No # 3+4

7

Topic:- Data Science Processes
Understanding
of Data - Types of
....Attributes.... (EDA)

- Analyzing data distribution
- Detecting outliers
- Measure similarity b/w data
- Objects. (Relation).

① Data - Types:-

→ "Structure Data" which
is easy to understand
(like, Spreadsheet, Tables)

→ Un-structured Data:-

Speech Signals, images, text
Written Text.

IMPORTANT:-

- Understanding different types
of Attributes in data

What is an Attributes?

① Column : an Excel sheet
representing a Property of data
Alternate Names

Attributes: from "database"

Feature: from "Machine Learning"

Variable: from "Math"

Example:-

① A data of Person:

Name, Height, weight, hair color.

Data Object:-

A **Row** in a data sets representing an individual entity

Example:-

In a Person datasets, each row represents a different person
 (Irfan, Sheeraz, Mujtaba)

Simple Terms-

- Simple in Statistic
- Observation in data analysis Research
- Record in database
- Instance in "Machine Learning"
- Data Points in visualization

Types of Attributes

Four major types of

Attributes-

- Nominal Attributes → Binary Attribute
- Ordinal Attributes → Numeric ...

* Nominal Variable

- nominal Attributes
- categorical data that serve as **Labels or Names** without any quantitative value or natural **"order"**

Example:- Colors of hair

black, white ...) blood groups (A, B, AB, O) or

only country's names

Inherent order is not appear.

* Possible Encoding
 can Represent Using Symbols or Code (black = 6
 white = 1)

only one operation of math is apply Frequency

Meaning operation most frequently hair colors.

★ Binary Attributes:-

- It is a Special types of nominal attributes. **only**
- Two values.

Examples:-

- Gender : Male / Female
- Covid : Positive (1) / Negative (0)
- Smoking status : Smoker or non-smoker

★ Sub-Types:-

Symmetric Binary Attributes

Both values are equally important

Example:- Gender (Male / Female)

Asymmetric Binary Attributes

Example:- Smoker (1 = Smoker, 0 = Non-smoker)

One value is more important

★ Interchangeable words (class, category, labels)

★ Ordinal Attributes:-

- ① Attributes whose values possess a meaningful order or Ranking between "categories" allowing for comparative relationships.

IMPORTANT Characteristics

- ② While we can determine if one value is Greater Than, other, the magnitude of differences b/w values cannot be precisely quantified

Example:- Academic Grade A+ (ellent)

→ A (very Good), B (Good)

Ranks : Assistant Professor, Associate

→ Professor, full Professor.

Represent :- Code, Simple not numerical

values Two features most frequent values is (modes). Between values (Median), Extract

from numerical (Not perform arithmetic operation)

is a way to convert numeric attribute to ordinal attribute.

* Discretization (convert) & ordinal variables :-

- Convert a numeric attribute into an ordinal attribute.

Example:- Temperatures Range.
Low ($1-15^{\circ}\text{C}$) → Medium ($16-30^{\circ}\text{C}$) →
High ($31-45^{\circ}\text{C}$)

* Numeric Attributes-

Quantitative attributes that are measurable. ~~arithmetical operations~~

* Interval Scale Attr. Operations

Key characteristics:-

- No absolute "zero" are meaningful diff b/w values
- Measured in equal size scale.

Operations that

- Comparison of order

Continuous value: floating point values: 1.23
Discrete value: whole numbers: 2, 4, 6

(10)

Two Types:-

Interval Scale Attributes

Ratio Scale Attributes

Examples:-

→ Date No zero "date (X) and Ratio are Temperature ($^{\circ}\text{C}$, $^{\circ}\text{F}$) not meaningful can +ve, 0 or -ve (Years) ~~Ex 20 $^{\circ}\text{C}$ is not twice as high as 10 $^{\circ}\text{C}$~~

* Ratio Scale Attributes

- has meaning absolute zero points
- Support all mathematical operation.
- Ratios b/w values are meaningful and interpretable.

Examples:-

- Length: 10 meters is twice as long as 5m
- Weight: 80kg is 25% heavier than 64kg
- Time: 4 hours is half duration of 8h

(Lectures NO#4+5)

(11)

Topic:-

Understanding Statistical description In Data Science

We have study different attributes of data in previous class and now this class lectures

We have study how "data is descripted" in one columns or values, so we use this techniques.

* Data Distribution :-

- How data is spread across different values.

* Data Set & Understanding :-
Methods To gain insight from your data.

* Basic Statistical Description :-
Tools To Summarize and interpret data patterns.

* Why Statistical Description Matters

* Data Shape

Understanding overall distribution pattern

of your datasets (Average value)

* Central Values

Identify the typical or representative values in your data

* Data Spread

Measure how widely values are distributed from Central Points.

(Mean, Median, mode)

* Noise and bias

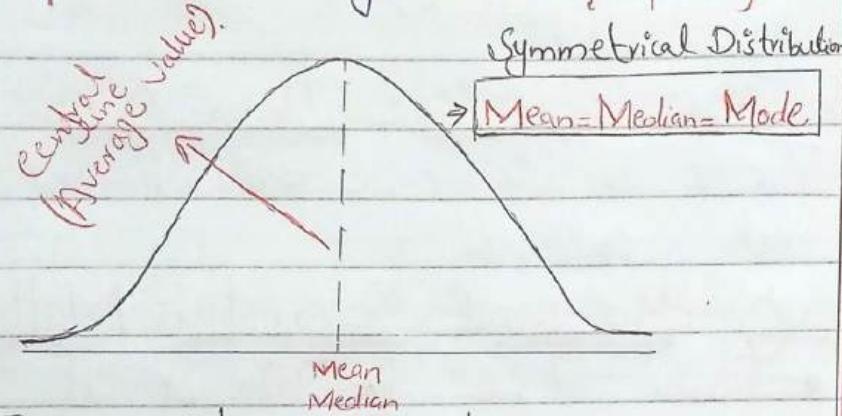
Detection :-
Spot Random Variation and Systematic errors in your dataset

(Types of data Shape)

(12)

* Symmetric data Distribution:-

Definition:- A form of distributed data where **Left** and **Right** parts are **Symmetric [Equal]**.



* Key Characteristic :-

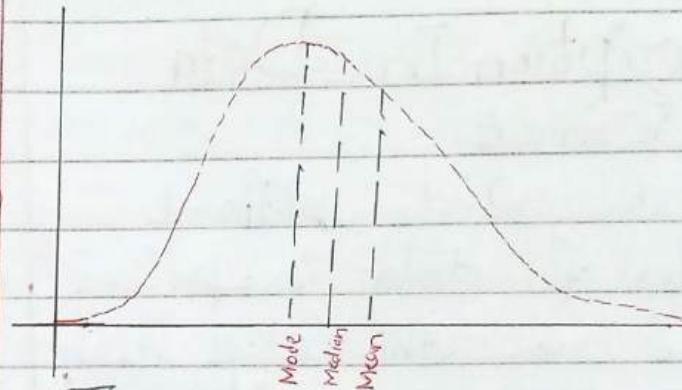
Data is Symmetric around the values.

* Contrast :-

The alternative form is called **skewed distribution**.

* Skewed distribution

Definition:- A data distribution that is not symmetric in shape.



* Types:-

Two Types Positive & negative skewed distribution.

Positive VS Negative

Neither type is symmetric around the central value.

(9)

★ Binary Attributes:-

- It is a Special Types of nominal attributes. only
- Two values.

Examples:-

- Gender : Male / Female
- Covid : Positive(1) / Negative(0)
- Smoking status : Smoker or non-smoker

★ Sub-Types :-

→ Symmetric Binary Attributes

Both values are equally important

Example:- Gender (Male / Female)

→ Asymmetric Binary Attributes

Example:- Smoker (1 = Smoker, 0 = Non-smoker)

One value is more important

★ Interchangeable words (class, category, Lable)

★ Ordinal Attributes:-

- ① Attributes whose values possess a meaningful order or ranking between "categories" allowing for comparative relationships.

IMPORTANT Characteristics

- ② While we can determine if one value is Greater Than, other, the magnitude of differences b/w values cannot be precisely quantified.

Example:- Academic Grade A+ (ellent)

A (very Good), B (Good),

Ranks :- Assistant Professor, Associate

Professor, full Professor.

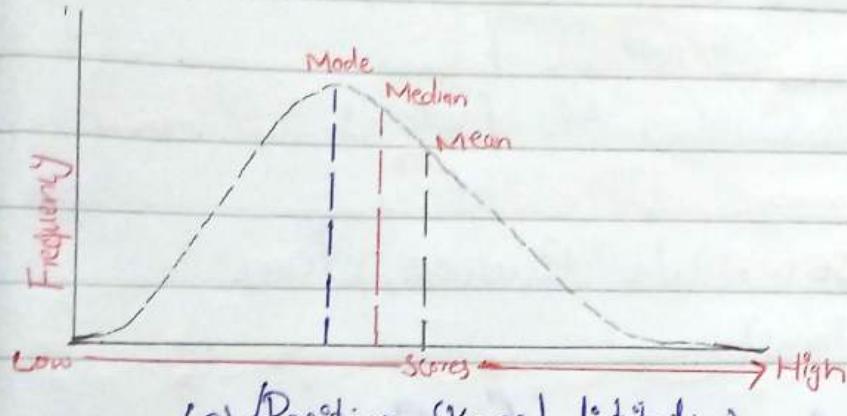
Represent :- Code, Simple not numerical

values Two features most frequent values

is (modes). Between values (Median), Extract from numerical (Not perform arithmetic operation)

Positive Skew (+)

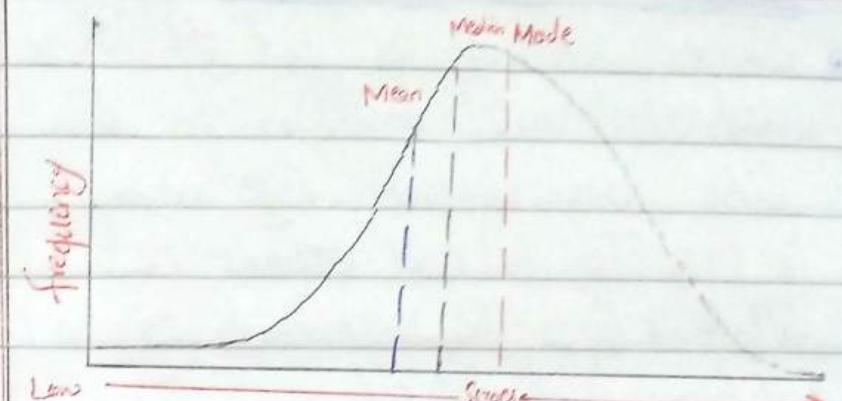
- Data Spread at Start, curve goes down slowly. Tail on Right side.



(a) Positive Skewed distribution

Negative Skewed :- (-)

- Curve gradually increases; then abruptly closes Tail on Left Side
- value (x^{restart}) > value (x^{left})



• This is not (b) (Negative Skewed distribution)

around the Symmetric central values."

* Basic Statistical distribution Measures

- Measure of Central Tendency
- Measure of Dispersion
- Measure of Approximety.

They all give good idea of data distribute.

→ Noise data :-

Random value

Example - (Online Quiz)

Take Online Quiz -
you complete an Online due to noise
Quiz get full Number but some
one distribute you like play song loudly

Skew/Biase data :-

Means data is leaning in one side (एक ओर लेने की)

particular side or data Biase ho

Ex:- student weight Range (54 to 63")

Add 8 inches in measurement all

All values are skewed Towards higher number (8). effect all data points

Outliers in DataSet:-

- Some values are out of Range which is not mention in data like one person Age column is (40 To 70) and others is fallal is (5 To 20) which is outliers in datasets.

Cleans:-

- Remove outliers and errors To ensure your analysis reflects actual patterns not Random variation.

Measures :-

- The are Three Main Types of Statistical description.
- (1) Central Tendency (2) Dispersion
- (3) Similarity

Central Tendency

Measures of Average or middle

values (Mean, Mode, Median, Trimmed Mean)

Similarity (Measure of

Relationships b/w features. (columns)
numeric attributes

Use:- Data Reduce, clean

Transforms:-

Measure Of Central

Tendency:-

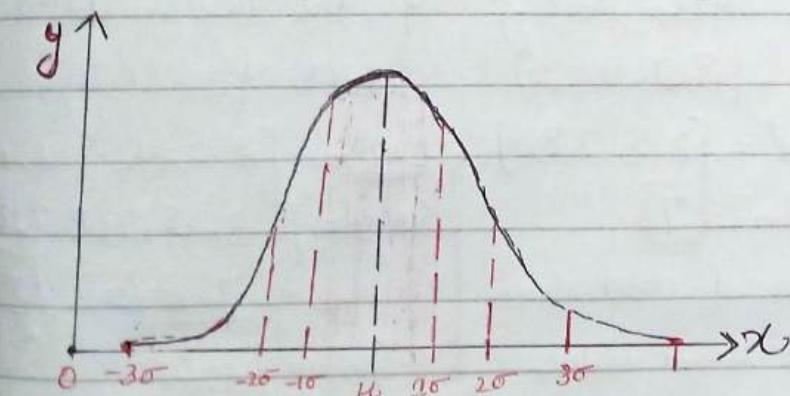
A Measure of central Tendency represent a Single value that describes the centre of dataset.

Variable, features, attributes which is divided into Two Parts (50% above, 50% down values)

* Data Division
It effectively divides the dataset into two equal halves

* Typical value
Indicates the most common or average value in the distribution.

Exp :- In a class where most students are similar height that common height represent the central tendency.



Arithmetic Mean)

Means (Simple Average).

Sum of all values divided by observation.

Formula:-

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

∴ $\sum x$ = Sum of all values
N = Num of values

Key Points :-

- ① Simple To Calculate
- highly sensitive to extreme values
- Used when data has no extreme value.

Example:- Salary of Employees.

Pos.

Crons (₹ 1 p.)

- Easy to calculate
- Very sensitive to extreme values.
- Cost Effective: Efficient computation power in large data
- When we used in data no extreme values is present.
- Comprehensive all calculation in every data points

(16)

Trimmed Mean

A Refined Average of Means.

is a modified mean where the highest and lowest values are Removed before Calculation.

Remove Extreme values. (4 person Salary)

Average and 5 highest Salary is Removed.

How Much You Trimmed Values
As a Rule, do not Trim more than 2% of the Total data Points.

Median : (Middle Value)

The median is the value that lies exactly in the Middle of Sorted dataset. (Ascending order)

Usually it is present in data. if you known in variable you can Able to 50% lower and 50% higher value differenate b/w them.

Calculate:-

Case 1 :-

when odd numbers

Select the middle value after sorting

Exp. :- [8, 5, 9, 3, 7]

Sorted = [3, 5, 7, 8, 9]

Median = 7

Numbers of observation is '5' Single value

Case 2 :-

when Even Numbers

Take the Average of the Two middles values after sorting.

Exp. :- [8, 5, 11, 9, 3, 7]

Sorted = [3, 5, 7, 8, 9, 11]

$$= \frac{7+8}{2} = \frac{15}{2}$$

Median = 7.5

Advantage	Disadvantage:-	(Advantage)	(disadvantage,)
→ not sensitive to extreme values.	→ Costly To calculate for large data	→ It apply all types of data	→ not exist in some dataset.
→ Use in data processing and cleaning	→ Limited data usage. only middle values.	→ not affected by extreme values.	→ Less useful for continuous data May not represent centre tendency well

① Mode (Most Frequent values)

- The mode is the value that appears most frequently in a datasets. It can be used with numerical and categorical values.

Example:-

dataset = [1, 2, 2, 5, 7, 8, 8, 9, 10]

$$\text{Mode} = \{2, 8\}$$

- Multiple values repeats and may be no mode in datasets.

↳ Multi Model data.

② Midrange: (Computation Purpose)

- The midRange is calculated as.

$$\frac{(\text{Minimum value} + \text{Maximum})}{2}$$

Provides a very quick and rough estimate of the centre. Not suitable for detailed analysis but for initial approximation

③ Data Shape:- [distribution]

Symmetric

$$\begin{array}{c|c|c} \text{Mean} = \text{Median} = \text{Mode} & \text{Mode} < \text{Median} & \text{Mean} < \text{Median} \\ \hline & & \text{Mean} < \text{Mode} \\ & & < \text{Mode} \end{array}$$

Positive Skew

Negative Skew

1)
2)

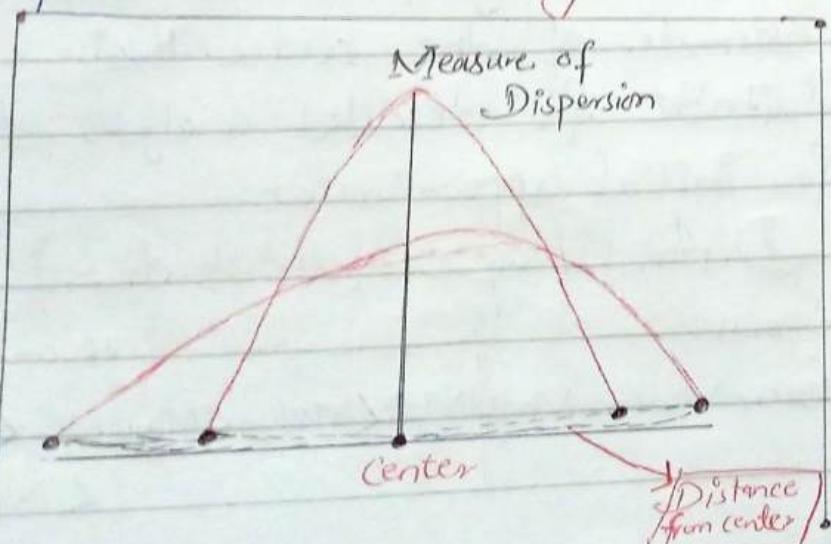
Lecture No # 6

Topic:-

Measure of Dispersion (BII 5)

Define:-

- ① It refers to how spread out the data is. It shows how far values lie from the central value (Mean or Median).
- ② Understanding dispersion helps us interpret the variability in data.



- Why is Dispersion Important
- Beyond Central Tendency
- It alone does not show data variability
- Added Context
- Dispersion adds context to the central value.

- ★ Choosing the Measure of Dispersion:
- Median based Measures
- If we use median we use quartiles and interquartile Range
- Mean based Measure
- If we use Mean, we use standard deviation and variance:

- Selection Criteria
- Choose depend upon the Measure of central Tendency used

Assuming:-

* (Dispersion from the Median) :-

- When the median is the central measure.

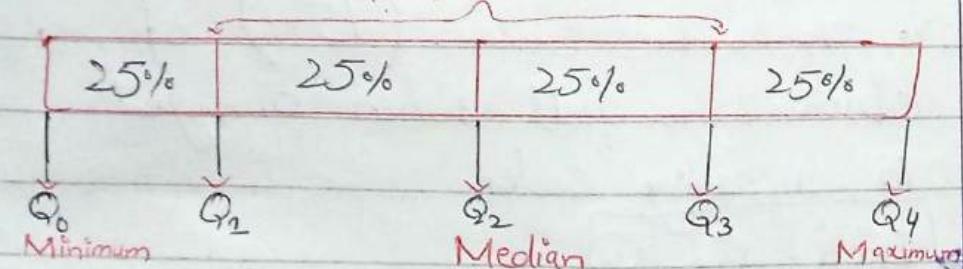
→ Quantiles

We use quantiles
(e.g. quartiles, percentiles)

→ Interquartile Range

We focus on
interquartile Range
(IQR)

(IQR)



* Understanding Quantiles :-

(1) Definitions -

Are the values that divide a dataset into equal-sized groups or portion.

(3) Median

When $q=2$, we get 1 quantile value → This is the Median, dividing data into 2 equal halves.

(2) Division :-

When we divide data into q equal parts, we get $(q-1)$ quantile values as boundaries b/w these parts.

(4) Quartiles

When $q=4$, we get 3 quantiles value → These are Quartiles dividing data into 4 equal quarters.

(20)

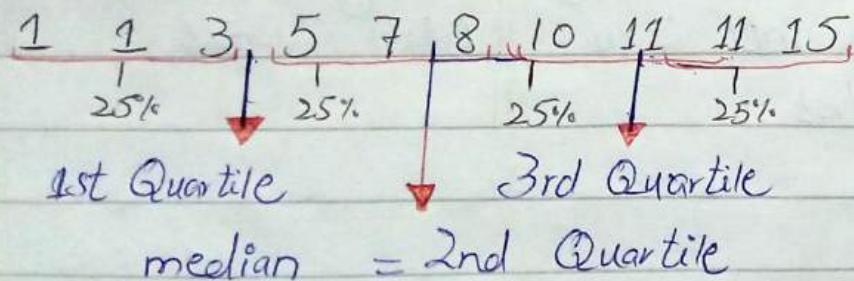
* Median as a Quantile

→ 2-Quantiles

The median is 20 especially, representing the middle value that divides a sorted dataset

→ Data Spread :-

① The median provides crucial insight into how values are distributed on either side of the central point revealing the dataset's balance.



* Quartiles Explained :-

→ Q1 first (Quartile)

Q1 Represent the 25th % the value below which 25% of observation fall when arranged in ascending order

Q3: (Third Quartile)

Q3:- Represent 0.75th

• the value below (75%) of observations fall when arranged in ascending order.

$$IQR = Q_3 - Q_1$$

(IQR):-

⑥ Measure data dispersion by focusing on the "middle values, eliminating the influence of outliers."

i) Definition & formula ii) Data Concentration

IQR represent the middle 50% of data points.

$IQR = Q_3 - Q_1$

Q_3 is the 75th% and Q_1 is 25th%

IQR measure the spread of the middle 50% of observations.

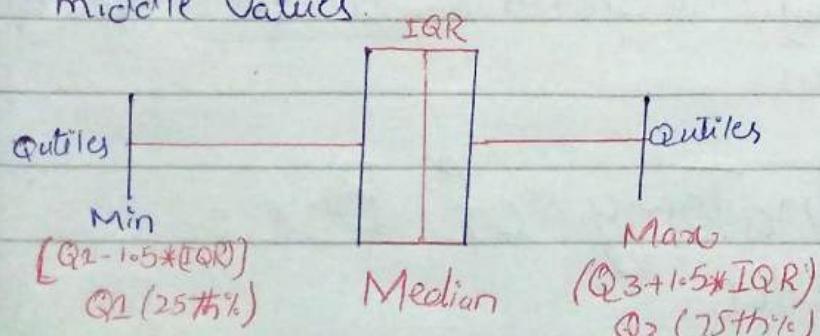
for extreme values.

(Small IQR)

iii) Interpreting IQR values :

⑥ middle values Large IQR widely spread

middle values.



*5-Number Summary:-

⑥ Give a complete picture of distribution and spread.

→ Minimum	→ Q_1 (first Quartile)	→ Median
The Smallest value in the dataset	25th% of the data	50th% percentile Middle value
→ Q_3 (3rd Quartile)	→ Maximum	
75th% Percentile	The largest value of the data	in the dataset.

Note:-

If you know the any Numeric feature we have use it To find.

⑥ Central Tendency of data
data Spread through IQR
Overall distribution of Shape.
Simply used IQR quartiles.

* Understanding data Shape:

① The data is skewed if more spread on one side of median.

i) Positively Skewed

$$Q_3 - Q_2 > Q_2 - Q_1$$

ii) Symmetric

Data is symmetric

$$Q_3 - Q_2 \approx Q_2 - Q_1$$

iii) Negative Skewed

$$Q_2 - Q_1 > Q_3 - Q_2$$

* Outliers :- What Are They?

① Are the extreme values outside the normal Range. They can distort analysis and give misleading insights.

* Identifying Outliers:-

→ Calculate IQR

$$IQR = Q_3 - Q_1$$

→ Lower Bound

Lower Bound

$$Q_1 - 1.5 \times IQR$$

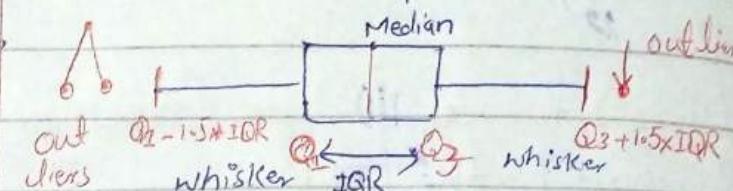
→ Upper Bound

Upper bound

$$Q_3 + 1.5 \times IQR$$

→ Identify

② Values beyond these are outliers



* Box & Whisker Plot

for visualization.

★ Complete Visualization

visual Representation of

Min, Q_1 , Median, Q_3

Max] Also Shows

outliers

★ Compositional

Multiple

distribution.

for Comparing.

* Anatomy of Box Plot

The Box

box show IQR
(Q₁.to Q₃)

The whiskers

Extent To Min
and Max (non-
outliers) values

Median Line

Line inside is
the 'Median' (Q₂)

Outliers

Dots show
outliers.

Variance:

It measure the Average
Squared deviation from Mean.

$$\rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$



if Measure is Means:-

Using this Standard Deviation

Square Root of Variance.

is the Square root of variance

$$\rightarrow \sigma = \sqrt{\sigma^2}$$

↓
effected
observation

Deviation Measure

Indicates how much data
deviates from the mean

→ The Box

box shows IQR
(Q₁ to Q₃)

→ The whiskers

Extent to Min and Max (non-outliers) values

→ Median Line

Line inside is the 'Median' (Q₂)

→ Outliers.

Dots show outliers

★ if Measure is Mean:-

① Using this Standard Deviation

② Square Root of Variance

is the Square root of variance

$$\rightarrow \sigma = \sqrt{\sigma^2}$$

(σ^{effective})
- ith observation

③ Deviation Measure

indicates how much data deviates from the mean

★ Variance:

It measures the Average Squared deviation from Mean.

$$\rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

~~*~*~*

Lectures No # 7 + 8

Topic:- Proximity & Similarity

★ Introduction:-

① In data Science, understanding the Relationship b/w different features

or Attributes / columns of a dataset is crucial for analysis and model development.

for
next

* Proximity and Similarity:-

- Measures are Techniques used To quantify how close or similar Two or More Attributes are.

* Importance of Similarity in DS

→ Improved data understanding

Feature Relationship

It identifies which features are related or independent.

Preprocessing Guidance

This information guides data preprocessing & modeling.

Numeric Attribute focused Techniques.

It is especially imp for Numeric attributes.

Similarity Range:-

(1) : vector are perfectly aligned

(0) : vector are orthogonal (no similarity)

(-1) : vectors are perfectly opposite.

* Representing Attributes as Vectors

Vector Representation

We represent each attribute in vectors

Mathematical foundation

The representation

covers the rich Theory of Linear algebra

Similarity calculation

Using vectors, we can easily calculate similarity

* Different Proximity Measures:-

Dot Product

We will begin with the dot product (•)

Covariance

Next, we discuss it with the dot product (•)

Cosine Similarity

(-∞, +∞)

Cosine Similarity

We move to the cosine

similarity

(-1, 1)

Correlation

Finally we introduce correlation

for evaluating strength and direction.

Rang (-1, 1)

* (.) Dot Product :-

→ It Takes Two vector as input and give output as number if $\boxed{0}$ numbers There is minimum difference Similarity.
if $\boxed{\text{num} > 0}$ then more chance in Similarity.

”push“ کا جو اسٹریکٹ push کے نتائج میں دیکھ سکتے ہیں اسے push کا نتیجہ کہا جاتا ہے۔

* Pros & Cons

Simple Calculation
Simple to calculate
and computationally
efficient

Magnitude dependency

* Cosine Similarity :-

→ Direction focus

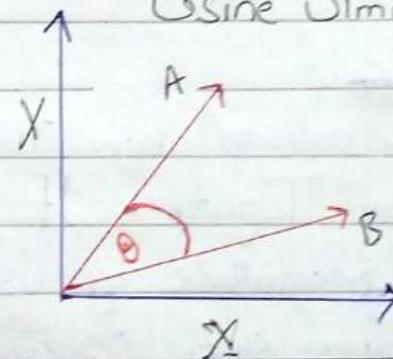
Focuses Solely on the direction of vectors

→ Magnitude direction
Resolve the magnitude
issue of dot product

→ Perfect alignment

$(+1)$ implies perfect alignment and (-1) implies opposite alignment.

↑ Cosine Similarity



$$\text{Sim}(A, B) = \text{Cos}(\theta) = \frac{A \cdot B}{|A| |B|}$$

Benefits :-

- its Magnitude b/w (-1, +1).
- if Cosine Similarity (-1) The ^{both} attributes is different.
- if Cosine Similarity (+1) Same direction attributes of both sides. (Trends Same) hy.
- if Cosine Similarity is '0' both is independent (no Relation).

Covariance :-

- Covariance Measures how Two Variable Vary Together. It Consider the Magnitude by (Subtracting the mean from each ^{vector} observation.)
- (+ve) Covariance indicates that Both attributes Trends To move Together.

→ (-ve) indicates that They move in opposite direction

Formula:-

* Covariance formula b/w Random Variable X, and Y.

$$\rightarrow \text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

where

- n = Observation Num
 - x_i = i-th value of variable X
 - y_i = i-th value of Variable Y
 - \bar{x} = Mean of X vector
 - \bar{y} = Mean of Y vector
- Range = $(-\infty, +\infty)$.

Correlation :-

- Direction of variable provided
But not strength is not provided.

Correlation :- (Mostly used) Data Reduction

- Perfect Relationship \rightarrow Bound Range.
 (± 1) indicates perfect relationship
- Standardization
- Derived by dividing Covariance by Standard deviation
- Relationship Measure \rightarrow Quantifies both Strength and direction

$$\rightarrow \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

" X by variables and \div by Standard Variation and although is variance"

Summary :-

The Summary of Proximity.

Measures	Characteristics
Dot Product	Simple operation, Magnitude dependent.
Cosine Similarity	Normalize vectors, focus
Covariance	Incorporates magnitude and direction but unbound
Correlation	Standardized version of Covariance b/w $(-1, +1)$

Topic:- Data Quality &

Pre-Processing

Data Preparation Components:

- Exploratory Data Analysis
- Data Preprocessing

سے دیجیٹل دادا کی کام کے لئے ایک بڑی کمپنی

④ Importance of Data Quality in Preprocessing.

Essential Step before modeling.

- Good Models
- Successful Machine Learning
- Essential foundation before modeling begins.

Good data

Quality preprocessing

Good data → Good

Models

Poor data → Poor Prediction

(Garbage In, Garbage Out)

* Data Quality:-?

→ Accuracy of values
Ensuring data correctly represents Reality.

Minimal Misery values

Complete inform for analysis.

→ Timelines and Relevance: Update inform from current needs.

Consistency in Attributes

Uniform formats and standards across dataset

not duplication

④ Real-World Data challenges

→ Missing values.

Incomplete Record

requiring handling

Inconsistent format

Varying Standard across records.

* Kaggle:-

- Using dataset and get it websites
- all issues in actual data.

* Preprocessing of Data:-

④ To improve the quality of data is called Preprocessing of data which give insight of business help.

Clean	Integrate	Reduce	Transform
→ Remove errors & inconsistencies	→ Combine multiple Resources	→ Reduce volume while Maintaining quality	→ Convert To suitable format for machine

* Techniques for Data Reduction

Several Methods use.

→ Feature aggregation

→ Combining multiple attributes into meaningful summary.

→ Feature Selection:

→ Identifying and retaining only the most relevant based variables.

→ (PCA: Principle Component Analysis) :-

→ Creating new uncorrelated variables that capture Max variance in dataset.

Dimensionality Reduction :-

→ Using T-SNE or UMAP To visualize and compress high-dimension into lower dimension.

* Rows Sampling:-

There is no duplication in our datasets. (20 Millions) Rows.

→ Original dataset → Sampling Process
→ Reduce dataset.

* Data Transformations:-

→ Scaling features To Similar range.
→ Making data uniform for algorithms
→ Prevent bias b/w larger scale feature Scales b/w (0,1)
→ Helps in faster and better convergence.

Before Transform	After Transform
→ features on diff Scales	→ Uniform features Scales
→ Algorithms bias Towards large Scale	→ Equal important in Algorithms