

Data Visualization:

is the representation of data through use of common graphics, such as Pie chart, bar chart, Histogram, scatter plot, and even animations. These visual information communicate complex data relationships and data-driven insights in a way that is easy to understand.

Importance of Data Visualization:

- i) Improved Understanding: Data visualization helps in comprehending complex data by allowing individuals to identify patterns, trends, and correlations that may not be apparent in raw data.
- ii) Enhanced Decision making: Data visualization enables decision-makers to gain insights quickly and make informed decisions. By presenting data in a visually appealing way, it becomes easier to identify key areas of focus, spot outliers, and understand the impact of different factors.
- iii) Effective Communication: Visual representations of data are often more accessible and engaging than raw numbers or text. Data visualization facilitates effective communication and storytelling by presenting information in a format that can be easily understood and remembered by the audience.

Purpose of Data Visualization:

- i) Data exploration: Visualization helps in understanding the structure and characteristics of data sets, enabling analysts to gain insight and formulate hypotheses.
- ii) Analyze Data: Visualization techniques allow analysts to identify trends, correlations, and patterns in data, aiding in data-driven decision-making and problem solving.
- iii) Communicate Findings: Visualization makes it easy to communicate complex data and analysis results to a broader audience. They enable effective story telling, supporting presentations, reports and dashboards.

Benefits of Data visualization Tools:-

- Solves data inefficiencies and absorb vast amount of data presented in visual formats.
- Increase the speed of decision making.
- Identify errors and inaccuracy in data quickly.
- Access real time information and assist in management function.
- It promotes story telling and conveys the right message to the audience.
- Explore business insights and achieve business growth.

Basic principles of data visualizations:-

- Truffe's law: This law states, that the human eye is drawn to the brightest and most contrasting elements in a visual. This means that it's important to use colours and fonts that will stand out and grab the viewer's attention.
- Gestalt's law: This theory of visual perception states that the human mind groups together similar objects or features. This can be used to create visual patterns and make data easier to understand.

Preattentive Processing:

is a type of visual processing that occurs before we consciously pay attention to something. This type of processing allows us to quickly identify objects and patterns in our environment. There are mainly four preattentive processing:

- Form: includes shape, size and orientation of an object.
- Color: includes hue, saturation, brightness of an object.
- Spatial positioning: includes location of an object in relation to other objects.
- Motion: includes movement of an object.

Data visualization is a powerful tool that can be used to communicate information in a clear and concise way. By following these one can create effective & engaging data visualizations.

Exploratory Data Analysis (EDA) :-

Exploratory Data Analysis (EDA) is an approach used in data science and statistics to analyze and summarize the main characteristics of a dataset. It involves a variety of techniques and visualizations to understand the data, identify patterns, detect anomalies, and generate hypotheses for further investigation.

The primary goal of EDA is to gain insights into the data and develop the understanding of its structure, distribution, relationships b/w variables. It helps researchers and data scientists to formulate appropriate questions, determine the quality and reliability of the data, and select appropriate statistical methods for further analysis.

Steps involved in EDA :-

- Data Collection (collecting required data from various sources)
- Finding all variables and understanding them (identify imp variables and understand their impact).
- Cleaning the Dataset (Remove null values and irrelevant information)
- Identify correlated variables (Find relationship b/w variables)
- Choosing the right statistical methods (use appropriate statistical tools for analysis)
- Visualizing and Analyzing Results: (Observe and interpret the findings)

Type of EDA:-
It is based on analyzing individual variables independently and uses

i) Univariate Analysis: focus on analyzing individual variables independently and uses visualizations such as histograms, box plots, bar charts, etc...

ii) Bivariate Analysis: focus on analyzing relationship b/w two variables. and uses visualizations such as scatter plots, line plots, correlation analysis

iii) Multivariate Analysis: analyzing relationship b/w multiple variables simultaneously and uses visualizations like heatmaps, pie charts and clustered analysis.

EDA user

- Understand the data better
- Identify patterns and trends
- spot anomalies
- Testing hypothesis
- Make better decisions

Machine learning Algorithms:

ML algorithms are computational models that enable computers to learn patterns and make predictions or decisions without being explicitly programmed. These algorithms learn and generalize patterns from input data to perform specific tasks, such as classification, regression, clustering or recommendation.

Supervised Learning: Labeled data is given and trained in presence of both input & output.

Eg:- Classification: image, text, voice recognition; Regression: Fraud detection.

Unsupervised :- unlabeled data, find hidden patterns by itself. Eg: Clustering: Segmentation, anomaly detection, association mining, Dimensionality reduction.

Semi-supervised: mixture of both sup & unsup learning. Data scientists may feed an algorithm mostly labeled training data, but the model is free to explore the data on its own and develop its own understanding. Eg (Machine translation), Fraud detection, labelling.

Reinforcement Learning: feedback from env by agent through actions.

Eg (Robotics, Video games, Resource management) Train with priors.

(Actions and their rewards, information gain) Address challenges of how to

Different ML algorithms:

i) Linear Regression: Supervised learning alg used for regression tasks. It models the relationship b/w a dependent variable and one (or) more independent variable.

fitting a linear eq to the observed data. The goal is to find best fit.

fitting line that minimizes the diff b/w the predicted and actual values.

From sklearn.model_selection import train_test_split

From sklearn.linear_model import LinearRegression, LogisticRegression.

From sklearn.tree import DecisionTreeClassifier.

From sklearn.ensemble import RandomForestClassifier.

From sklearn.svm import SVC, SVR

From sklearn.neighbors import KNeighborsClassifier

From sklearn.naive_bayes import GaussianNB

From sklearn.metrics import accuracy_score, mean_squared_error

import numpy as np.

sample input data

```

x = np.array([1,2],[2,3],[3,4],[4,5],[5,6])
y = np.array([10,20,30,40,50])

```

load and split the dataset:

```

x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.2,random_state=42)

```

linear_regression = Linear_Regression()

linear_regression.fit(x_train, y_train)

linear_Prediction = linear_regression.predict(x_test)

model evaluation

```

mse = mean_squared_error(y_test, linear_Prediction)
accuracy = accuracy_score(y_test, linear_Prediction)

```

```

print("Linear Regression MSE : ", mse)
print("Linear Regression Accuracy : ", accuracy)

```

i) Logistic Regression: supervised learning algo used for classification tasks. models relationship by fitting logistic function to data. Predicts probability of an instance belonging to a particular class.

ii) Decision Trees: unsup algo handles both reg & classi tasks. Tree like structure where internal nodes represents features, branch → decisions : leaf node → outcomes.

iii) Random Forest: combines multiple decision trees to improve predicted accuracy. creates forest of Decision trees. handles high-dim data and avoid overfitting.

iv) Support vector Machine: sup algo handles both reg & classi tasks. finds optimal hyperplane that maximizes the margin b/w diff classes, enabling effective separation.

v) K-Nearest Neighbors: sup for reg & classi tasks, classify new instances based on the majority vote of their NN. "K" determines no. of neighbors considered for decision making.

vi) K-mean clustering: unsup algo for classi, regg, divides dataset into K clusters based on the similarity of instances. Adjust centroids until convergence.

ML uses: predictive analysis, NLP, CV, Robotics, etc...

a) Underfitting and Overfitting in Machine Learning:

* Underfitting:-

Underfitting occurs when a machine learning model does not learn the training data well enough. This can happen for a number of reasons, including:

- The model is too simple and cannot capture the underlying patterns in the data.
- The model is not trained on enough data.
- The data is noisy and corrupted.

Underfitting can be identified by looking at the model's performance on the training data and the test data. If the model performs poorly on both training and testing data, then it is likely that the model is underfitting.

* Overfitting:-

Overfitting occurs when a machine learning model learns the training data too well. This can happen for a number of reasons, including:

- The model is too complex and learns the noise in the data as well as the underlying patterns.
- The model is trained on too much data.
- The data is not representative of the real world.

Overfitting can be identified by looking at the model's performance on training and testing data. If the model performs well on the training data but poorly on the test data, then it is likely that the model is overfitting.

* Preventing Underfitting and Overfitting:

There are a number of things that can be done to prevent underfitting and overfitting, including:

- Selecting the right model: The complexity of the model should be chosen based on the amount of training data and the complexity of the problem.

Datafication:-

- Process of turning aspects of the physical world into data that can be analyzed, done through using sensors, cameras, etc..
- Key part of development of IoT that makes it useful, smart & useful.
- Benefits, such as helping us to understand the world around us better and improve decision making, e.g.: Traffic tracking, monitoring our health, used to discriminate against people.
- Having risks, such as invading our privacy and being used to discriminate against people.

Big data:

- E.g.: Smart phones; fitness trackers; self-driving cars.
- Datafication rapidly growing → major impact on our lives.
- Primary characteristics: Volume → Big data sets are extremely large (Petabytes insize).
- Velocity → Big data sets are generated & processed at high speed.
- Variety → Structured (organized in regular way e.g.: databases) unstructured (not organized in any way e.g. text, img, video).
- Examples: Social media data, sensor data, log data, etc.

Hypothesis testing:

- is an act in statistics whereby an analyst tests an assumption regarding a population parameter. Statistical method to determine whether there is significant difference between two datasets.
- The goal of hypothesis testing is to determine whether there is enough evidence to reject the null hypothesis. If null hyp. is rejected then alternative hyp. is accepted.
- Type I error: Rejecting the null hypothesis when it is actually true.
 - Significance level: (a) Probability of making type I error, often set at 0.05 or 0.01.

Type II error: Failing to reject the null hypothesis when it is actually false.

Power of a test (β) The probability of Type II error calculated by $1 - \beta$.

→ best way to reduce type I & II error is to increase the sample size which helps to detect diff b/w two sets of data.

Clustering:-

unsupervised ml algorithm that groups data points into clusters based on similarity. The goal of clustering is to find grp of data pts similar to each other, and diff from data points in other clusters.

Types of clustering:

i) K-mean clustering:

is a simple and popular clustering alg that groups data pts into fixed no/of clusters. The alg work by iteratively assigning data pts to clusters with the nearest mean, until the convergence.

steps:- i) Choose no/of clusters → ii) Initialize cluster centroids → iii) Assign data pts to ~~cluster~~ cluster with nearest centroid → iv) update the mean of centroids in a cluster → v) repeat until convergence. (of centroids, produce suboptimal results)

→ k-mean is good choice for datasets with fixed no/of clusters, sensitive in initial choice

ii) Hierarchical clustering:

is a type of clustering alg that builds a tree-like structure of clusters where each cluster can be a parent (or) child of other clusters. The algorithm starts by treating each data points as its own cluster. Then it repeatedly merge the similar clusters until there is only one cluster left.

Clusters until there is only one cluster left.

→ It can be agglomerative (bottom-up) or divisive (top-down)

→ It can be agglomerative (bottom-up) or divisive (top-down)

→ does not requires to specify the no/of clusters in advance.

→ does not requires to specify the no/of clusters in advance.

iii) DBSCAN: (Density based spatial clustering of Applications with Noise) identifies and groups data points that are densely packed, while considering pts that are not densely packed as noise. doesn't requires to specify the no/of clusters in advance.

→ Capable of identifying clusters of diff shapes & size, computationally expensive, tough to choose esp, min pts.

→ core points, border pts, noise min pts : min pts within neighbourhood to make it as core pts

Eg of clustering: Customer segmentation, Data exploration, NLP, image segmentation, anomaly detection,

Python code for data cleaning

```
import pandas as pd
```

```
data = pd.read_csv('data.csv') # Load the data set
```

```
missing_values = data.isnull().sum() # Check for missing values.
```

```
Point("Missing values!", missing_values)
```

```
data = data.dropna() # Remove missing values & duplicates
```

```
data = data.drop_duplicates() # Remove duplicates
```

```
data['Date'] = pd.to_datetime(data['Date']) # Convert data type.
```

```
data['Value'] = pd.to_numeric(data['dateValue'])
```

```
data['Name'] = data['Name'].str.strip() # Remove trailing (or) leading white
```

```
data['Name'] = data['Name'].str.lower() # Convert to lower case.
```

```
data = data[(data['Value'] > lower_threshold) & (data['Value'] < upper_threshold)]
```

```
data.to_csv('cleaned_data.csv', index=False) # Export cleaned data into new CSV
```

```
(CSV file contains all states (USA), brands taken after cleaning (removal of  
columns having null values present))
```

Student T-test: It is also known as independent sample t-test that is used to check if the difference between two independent groups is statistically significant. Student T-test is a statistical hypothesis test that is used to compare the mean of two independent group of data. The t-test calculates the probability that the difference between the two means is due to chance.

to the standard error of the difference.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{S_{\bar{x}}}$$

$$\text{where } S_{\bar{x}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

where \bar{x}_1 = mean of first set of values,

\bar{x}_2 = mean of second set of values,

s_1 = standard deviation of first set of values

" second " "

s_2 = standard deviation of second set of values

n_1 = Total no. of values in first set,

" second " "

n_2 = Total no. of values in second set,

i) Data preprocessing :-

- It is the process of transforming raw data into a format that can be easily analyzed by ML algorithms. It involves several steps, including data cleaning, integration, transformation and reduction.

ii) Data cleaning :-

- process of identifying and removing (or) correcting errors, inconsistencies and the missing values from data.

Techniques used for data cleaning includes imputation, deduplication, outlier removal and handling missing values.

- Example: filling in missing values with estimated values when there is missing data in the 'age' column of dataset, this is called imputation.

iii) Data integration :-

- Process of combining data from multiple sources and resolving inconsistencies or conflicts.

Techniques used includes data fusion, data alignment and data merging,

- Example: merging of customer data based on common identifier, such as customer ID and then align the information.



iv) Data transformation :-

- Process of converting data into a format suitable for machine learning algorithms

Techniques used includes normalization, feature scaling, encoding & data discretization

- Example: Converting categorical variable like colour into numerical value by encoding enables the ML algorithm to use it directly.

v) Data Reduction:-

- Process of reducing the amount of data while retaining as much relevant info as possible

Techniques used includes feature selection, dimensionality reduction and data sampling,

- Example: Selecting a subset of features that are most relevant to the analysis by dimensionality reduction to increase computation speed of the model.

Statistics is the study of the collection, analysis, interpretation, presentation and organization of data, in almost every field of life.

e.g.: mean of marks obtained by 50 students in a class.

Basics of statistics include the measure of central tendency (mean, median, mode) and the measure of dispersion (variance and standard deviation).

Types of Statistics:-

i) Descriptive statistics :-

These statistics are used to summarize and describe a data set. Histograms, pie charts, bar and scatter plots are common ways to summarize data and present it in table and graphs.

Here, the summarization is done using:-

Measures of central tendency : mean, median, mode

Measures of dispersion : range, variance, standard deviation.

Measures of shape : skewness, kurtosis.

For instance, if we have a data set containing the weights of 20 people,

we could use descriptive statistics to summarize and describe the data, such as calculating the mean weight, the range of weights, or the

skewness of the weight distribution.

ii) Inferential statistics

These statistics are used to draw conclusions or make predictions about a population based on a sample of data.

Here, the prediction is done using (t-tests, ANOVA, chi-square test & linear regression).

Hypothesis testing : testing whether a sample mean is significantly different from population mean.

Confidence intervals : estimating the range of values within which a population parameter is likely to fall.

Regression analysis : analyzing the relationship between variables and predicting values based on that relationship.

For example, we might use inferential statistics to test whether a new drug reduces cholesterol levels in a population; we would take a sample of people and randomly assign them to a treatment group or a control group, and then use statistical tests to determine whether there is a significant difference in cholesterol levels between the two groups.

Probability :-

Probability is a mathematical concept that predicts how likely events are to occur. The probability values are expressed between 0 & 1. The definition of probability is the degree to which something is likely to occur. These are mainly two types of probability distributions:

i) Discrete probability :-

These distributions are associated with random variables that can only take on a finite or countable number of values, such as the outcome of a coin toss or the number of defective items in a batch of products.

ii) Continuous probability :-

These distributions on other hand are associated with random variables that can take on any value within a range, such as the height or weight of individuals in a population.

Examples of common probability distribution include the normal distribution (also known as Gaussian distribution), the binomial distribution, the poisson distribution, and the exponential distribution.

Probability distributions are widely used in fields such as statistics, physics, engineering, finance, and many other areas of science and industry to model and analyze a wide range of phenomena.

SVD (singular value decomposition) :-

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}$$

rotation stretching rotation

$$A^T A = (V \cdot \Sigma^T \cdot U^T)(U \cdot \Sigma \cdot V^T)$$

$$= V \cdot \Sigma^T \cdot \Sigma \cdot V^T$$

$$A^T A = (V \cdot \Sigma^T \cdot U^T)(U \cdot \Sigma \cdot V^T)$$

$$= V \cdot \Sigma^T \cdot \Sigma \cdot V^T$$

$$A^T A = V \cdot \Sigma^T \cdot \Sigma \cdot V^T$$

rotation cd

$$\text{eg: } A = \begin{pmatrix} 3 & 1 & 1 \\ 1 & 3 & 0 \\ 1 & 1 & 3 \end{pmatrix}_{3 \times 3}$$

$$A^T = \begin{pmatrix} 3 & -1 & 1 \\ 1 & 3 & 1 \\ 1 & 1 & 3 \end{pmatrix}_{3 \times 3}$$

$$A \cdot A^T = \begin{pmatrix} 3 & 1 & 1 \\ 1 & 3 & 0 \\ 1 & 1 & 3 \end{pmatrix} \cdot \begin{pmatrix} 3 & -1 & 1 \\ 1 & 3 & 1 \\ 1 & 1 & 3 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}_{3 \times 3}$$

$$A \cdot A^T = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}_{3 \times 3}$$

$$\Rightarrow |A - \lambda I| = 0 \Rightarrow (10 - \lambda)^2 \cdot 1^2 = 0 \Rightarrow \lambda = 10$$

$$\Rightarrow \begin{vmatrix} 10 - \lambda & 1 & 1 \\ 1 & 10 - \lambda & 0 \\ 1 & 0 & 10 - \lambda \end{vmatrix} = 0$$

$$\begin{vmatrix} 10 - \lambda & 1 & x_1 \\ 1 & 10 - \lambda & x_2 \\ 1 & 0 & 10 - \lambda \end{vmatrix} = 0 \Rightarrow \begin{vmatrix} 10 - \lambda & x_1 \\ 1 & x_2 \end{vmatrix} = 0$$

$$\begin{vmatrix} 10 - \lambda & 1 & x_1 \\ 1 & 10 - \lambda & x_2 \\ 1 & 0 & 10 - \lambda \end{vmatrix} = 0 \Rightarrow \begin{vmatrix} 10 - \lambda & x_1 \\ 1 & x_2 \end{vmatrix} = 0$$

$$\text{we need to write in descending order so,}$$

$$\text{orthogonal, for above}$$

$$\text{orthogonal} \rightarrow \begin{pmatrix} 0 & 0 & \sqrt{6} \\ 0 & \sqrt{6} & 0 \end{pmatrix}$$

here U & V are orthogonal i.e. $U \cdot U^T = I$
 $V \cdot V^T = I$,
& Σ is a rectangular matrix with singular values

rotation : $\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$ unitary transformation

stretching : $\begin{pmatrix} \alpha & 0 \\ 0 & \alpha \end{pmatrix}$

$$A^T A = V \cdot \Sigma^T \cdot \Sigma \cdot V^T = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{pmatrix}$$

now for V : $A^T \cdot A$.

$$\rightarrow \begin{vmatrix} 3 & -1 & 1 \\ 1 & 3 & -1 \\ -1 & 1 & 3 \end{vmatrix} \begin{vmatrix} 3 & 1 & 1 \\ 1 & 3 & -1 \\ -1 & 1 & 3 \end{vmatrix} \Rightarrow \begin{vmatrix} 100 & 0 & 200 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{vmatrix} \rightarrow \lambda^3 - S_1 \lambda^2 + S_2 \lambda - S_3 = 0$$

$S_1 = 22$ (sum of diagonal elements)

$$S_2 = 100 + (+4) + 16 = 120$$

$$S_3 = 10(4) - 0 + 2(-20) = 0$$

position of $\begin{pmatrix} 0 & 0 & 2 \\ 0 & 0 & 4 \\ 0 & 0 & 2 \end{pmatrix}$: position

$$\rightarrow \lambda^3 - 22\lambda^2 + 120\lambda = 0$$

$$\rightarrow \lambda(\lambda^2 - 22\lambda + 120) = 0$$

$$\rightarrow \lambda^2 - 12\lambda - 10\lambda + 120 = 0$$

$$\therefore \lambda = 0; \lambda = 10; \lambda = 12 \rightarrow \lambda_1 = 12; \lambda_2 = 10; \lambda_3 = 0$$

by cofactor rule!

$$\therefore \text{for } \lambda = 12 \therefore \begin{vmatrix} 0 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = 0 \Rightarrow \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} + T_{A \cdot A}^{-1} \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix}$$

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\text{for } \lambda = 10 \therefore \begin{vmatrix} 0 & 0 & 2 \\ 0 & 0 & 4 \\ 2 & 4 & 2 \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = 0 \Rightarrow \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{matrix} -x_2 \\ -x_1 \\ 0 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} + T_{A \cdot A}^{-1} \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix}$$

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\text{for } \lambda = 0 \therefore \begin{vmatrix} 0 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{vmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = 0 \Rightarrow \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} = \begin{matrix} -x_2 \\ -x_1 \\ 0 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} + T_{A \cdot A}^{-1} \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} = \begin{matrix} 0 \\ 0 \\ 0 \end{matrix}$$

$$V = \begin{vmatrix} 1 & 2 & 1 \\ 2 & -1 & 2 \\ 1 & 0 & -5 \end{vmatrix} \quad \text{eigenvalues} \quad V^T = \begin{vmatrix} 1 & 2 & 1 \\ 2 & -1 & 0 \\ 1 & 0 & -5 \end{vmatrix}$$

$$\text{orthogonal } V = \begin{pmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{5}{\sqrt{30}} \end{pmatrix} \quad \text{eigenvalues} \quad \text{here } \lambda_1, \lambda_2, \lambda_3 \text{ are } 12, 10 \text{ & } 0 \text{ as}$$

$$\Sigma = \text{diagonal matrix with same dimensionality as } A \Rightarrow \begin{pmatrix} \sqrt{\lambda_1} & 0 & 0 \\ 0 & \sqrt{\lambda_2} & 0 \\ 0 & 0 & \sqrt{\lambda_3} \end{pmatrix}$$

$$\therefore A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{5}{\sqrt{30}} \end{pmatrix}$$

$$\begin{pmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{pmatrix} \leftarrow \text{descending order}$$

a) PCA for

x	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1	1.8
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9	n = 10

x	y	$x_i - \bar{x}$ Ⓐ	$y_i - \bar{y}$ Ⓑ	AB	A^2	$F \cdot B^2$
2.5	2.4	0.69	0.49	0.3381	0.4761	0.2401
0.5	0.7	-1.31	-1.21	1.5851	1.7161	1.4641
2.2	2.9	0.39	0.99	0.3861	0.1521	0.9801
1.9	2.2	0.09	0.29	0.0261	0.0081	0.0841
3.1	3.0	1.29	1.09	1.4061	1.6641	1.1881
2.3	2.7	0.49	0.79	0.3871	0.2401	0.6241
2.0	1.6	0.19	-0.31	0.4589	0.0361	0.0961
1.0	1.1	-0.81	-0.81	0.6561	0.6561	0.6561
1.5	1.6	-0.31	-0.31	0.0961	0.0961	0.0961
1.1	0.9	-0.71	-1.01	0.7171	0.5041	1.0201
				5.539	5.549	5.449
$\bar{x} = 1.81 \bar{y} = 1.91$						

Covariance matrix = $\begin{bmatrix} \text{Cov}(x, x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Cov}(y, y) \end{bmatrix}$ where $\text{Cov}(x, y) = \text{Cov}(y, x)$

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

\Rightarrow Covariance matrix = $\begin{bmatrix} 5.549 & 5.539 \\ 5.539 & 5.449 \end{bmatrix} = \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix}$

Now, we need to find eigen values & eigen vectors of the above covariance matrix.

$$\text{Q. } |\lambda I - A| = 0$$

$$\Rightarrow \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} (0.6166 - \lambda) & 0.6154 \\ 0.6154 & (0.7166 - \lambda) \end{bmatrix} = 0 \quad \text{--- eq(1)}$$

By finding det for above matrix, we will get 0 an equation, which can be solved to get the eigen values:

$$\Rightarrow (0.6166 - \lambda)(0.7166 - \lambda) - (0.6154)^2 = 0$$

$$\Rightarrow 0.441 - 0.7166\lambda - 0.6166\lambda + \lambda^2 - 0.378 = 0$$

$$\Rightarrow \lambda^2 - 1.3332\lambda + 0.063 = 0$$

now, we will find roots for above equation :- $\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

$$\Rightarrow \lambda_1 = 1.283 \quad \& \quad \lambda_2 = 0.049$$

(i) now, putting λ_1 in eq. (1) for eigen vectors

$$\begin{bmatrix} (0.6166 - 1.283) & 0.6154 \\ 0.6154 & (0.7166 - 1.283) \end{bmatrix} \begin{pmatrix} x, y \end{pmatrix}_{\text{vectors}} = \begin{pmatrix} 0, 0 \end{pmatrix}_{\text{vectors}}$$

(or)

eigen vectors can be found by using :- $CV = \lambda V$

$$\Rightarrow \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 1.283 \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}$$

$$\Rightarrow 0.6166x_1 + 0.6154y_1 = 1.283x_1 \Rightarrow x_1 = 0.923y_1$$

$$\Rightarrow 0.6154x_1 + 0.7166y_1 = 1.283y_1$$

~~$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 0.923 \\ 1 \end{bmatrix}$$~~

$$\therefore \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} \frac{x_1}{\sqrt{x_1^2+y_1^2}} \\ \frac{y_1}{\sqrt{x_1^2+y_1^2}} \end{bmatrix} = \begin{bmatrix} \frac{0.923}{\sqrt{1.851}} \\ \frac{1}{\sqrt{1.851}} \end{bmatrix} \cdot \begin{bmatrix} \frac{0.923}{1.36} \\ \frac{1}{1.36} \end{bmatrix} = \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}$$

for $\lambda = 0.049 \Rightarrow CV = \lambda V$

$$\begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 0.049 \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}$$

$$x_1 = -\frac{0.6154}{0.5676} y_1$$

$$\Rightarrow 0.6166x_1 + 0.6154y_1 = 0.049x_1 \Rightarrow$$

$$0.6154x_1 + 0.7166y_1 = 0.049y_1$$

$$\Rightarrow x_1 = -1.084 y_1$$

$$\Rightarrow \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} \frac{-1.084}{\sqrt{2.175}} \\ \frac{1}{\sqrt{2.175}} \end{bmatrix} = \begin{bmatrix} \frac{-1.084}{1.474} \\ \frac{1}{1.474} \end{bmatrix} = \begin{bmatrix} -0.735 \\ 0.678 \end{bmatrix}$$

SVD: $A_{m \times n} \rightarrow \text{rank}(A)$

$\Sigma_{m \times n}; U_{m \times m}; V_{n \times n}$

$\Rightarrow A = U \Sigma V^T$ & find eigen value for that matrix
 from given "A" find $A^T \cdot A$ & find eigen values [like $\sigma_1 = \sqrt{\lambda_1}; \sigma_2 = \sqrt{\lambda_2}, \dots$].
 and also find σ values and place beside each other
 & from eigen values find eigen vectors and place beside each other
 by normalizing them to get "V".
 \Rightarrow find Σ , as it is a diagonal matrix of values " σ " & make its rank equal to A i.e. $\text{Rank}(\Sigma) = \text{Rank}(A)$.

$$\Rightarrow U = [U_1, U_2]$$

$$\text{now } U_1 = \frac{1}{\sigma_1} AV_1 \text{ & } U_2 = \frac{1}{\sigma_2} AV_2$$

$$A = \frac{U}{\sigma_1} \Sigma \frac{V^T}{\text{diagonal matrix}}$$

$U = \text{normalized Eigen vector matrix } A A^T$

$V^T = \text{normalized Eigen vector matrix } A^T A$