

## Module No. 2

### Exploratory Data Analysis

Exploratory Data Analysis and the Data Science Process, Basic tools (plots, graphs and summary statistics) of EDA, Measuring similarity and dissimilarity.

#### **Exploratory Data Analysis:**

##### **Introduction:**

Exploratory Data Analysis (EDA) is an approach to analyzing and summarizing data sets with the primary goal of gaining insights into the data's structure, patterns, relationships, and distributions. It involves using statistical graphics, visualization techniques, and summary statistics to understand the main characteristics of the data and uncover hidden patterns or trends.

##### **Key aspects of Exploratory Data Analysis include:**

- **Data Cleaning:** Before diving into analysis, it's essential to clean and preprocess the data. This involves handling missing values, outliers, and any inconsistencies in the dataset.
- **Descriptive Statistics:** EDA often begins with computing and examining descriptive statistics such as mean, median, mode, range, and standard deviation. These statistics provide a basic summary of the main features of the data.
- **Visualization:** Visualizations play a crucial role in EDA. Plots like histograms, box plots, scatter plots, and correlation matrices help in visually understanding the distribution of data, relationships between variables, and potential patterns or trends.
- **Pattern Recognition:** EDA aims to identify any patterns or trends in the data. This could involve looking for clusters, trends over time, or any other recurring structures within the dataset.

- **Outlier Detection:** Identifying outliers is important in understanding the data distribution and ensuring that extreme values do not unduly influence the analysis.
- **Correlation Analysis:** EDA often involves exploring relationships between variables. Correlation analysis helps in understanding the strength and direction of relationships between pairs of variables.
- **Dimensionality Reduction:** In some cases, high-dimensional data may be simplified using techniques like Principal Component Analysis (PCA) to identify the most important features.

EDA is a crucial step in the data analysis process because it helps analysts and data scientists to formulate hypotheses, guide subsequent analyses, and make informed decisions about modeling and further exploration. It is often the first step in any data analysis project, providing a foundation for more advanced statistical and machine learning techniques.

A data scientist involves almost 70% of his work in doing the EDA of his dataset.

## Exploratory Data Analysis (EDA) Using Python Libraries

For the simplicity of the article, we will use a single dataset. We will use the employee data for this. It contains 8 columns namely – First Name, Gender, Start Date, Last Login, Salary, Bonus%, Senior Management, and Team. We can get the dataset - [Employees.csv](#)

Let's read the dataset using the Pandas [read\\_csv\(\)](#) function and print the 1st five rows. To print the first five rows we will use the [head\(\)](#) function.

```
import pandas as pd

import numpy as np

# read dataset using pandas

df = pd.read_csv('employees.csv')
```

```
df.head()
```

### Output:

	First Name	Gender	Start Date	Last Login Time	Salary	Bonus %	Senior Management	Team
0	Douglas	Male	8/6/1993	12:42 PM	97308	6.945	True	Marketing
1	Thomas	Male	3/31/1996	6:53 AM	61933	4.170	True	NaN
2	Maria	Female	4/23/1993	11:17 AM	130590	11.858	False	Finance
3	Jerry	Male	3/4/2005	1:00 PM	138705	9.340	True	Finance
4	Larry	Male	1/24/1998	4:47 PM	101004	1.389	True	Client Services

*First five rows of the dataframe*

Let's see the shape of the data using the shape.

```
df.shape
```

### Output:

(1000, 8)

This means that this dataset has 1000 rows and 8 columns.

Let's get a quick summary of the dataset using the pandas [describe\(\)](#) method. The describe() function applies basic statistical computations on the dataset like extreme values, count of data points standard deviation, etc. Any missing value or NaN value is automatically skipped. describe() function gives a good picture of the distribution of data.

### Example:

```
df.describe()
```

### Output:

	Salary	Bonus %
count	1000.000000	1000.000000
mean	90662.181000	10.207555
std	32923.693342	5.528481
min	35013.000000	1.015000
25%	62613.000000	5.401750
50%	90428.000000	9.838500
75%	118740.250000	14.838000
max	149908.000000	19.944000

*description of the dataframe*

Note we can also get the description of categorical columns of the dataset if we specify ***include = 'all'*** in the describe function.

Now, let's also see the columns and their data types. For this, we will use the [info\(\)](#) method.

```
# information about the dataset

df.info()
```

## Output:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   First Name            933 non-null    object
1   Gender                855 non-null    object
2   Start Date            1000 non-null   object
3   Last Login Time       1000 non-null   object
4   Salary                1000 non-null   int64
5   Bonus %               1000 non-null   float64
6   Senior Management     933 non-null    object
7   Team                  957 non-null    object
dtypes: float64(1), int64(1), object(6)
memory usage: 62.6+ KB
```

*Information about the dataset*

We can see the number of unique elements in our dataset. This will help us in deciding which type of encoding to choose for converting categorical columns into numerical columns.

```
df.nunique()
```

### Output:

First Name	200
Gender	2
Start Date	972
Last Login Time	720
Salary	995
Bonus %	971
Senior Management	2
Team	10
dtype: int64	

Till now we have got an idea about the dataset used. Now Let's see if our dataset contains any missing values or not.

## Handling Missing Values

You all must be wondering why a dataset will contain any missing values. It can occur when no information is provided for one or more items or for a whole unit. For Example, Suppose different users being surveyed may choose not to share their income, and some users may choose not to share their address in this way many datasets went missing. Missing Data is a very big problem in real-life scenarios. Missing Data can also refer to as NA(Not Available) values in pandas. There are several useful functions for detecting, removing, and replacing null values in Pandas DataFrame :

- [isnull\(\)](#)
- [notnull\(\)](#)
- [dropna\(\)](#)
- [fillna\(\)](#)
- [replace\(\)](#)
- [interpolate\(\)](#)

Now let's check if there are any missing values in our dataset or not.

### Example:

```
df.isnull().sum()
```

### Output:

```
First Name      67
Gender          145
Start Date       0
Last Login Time  0
Salary           0
Bonus %         0
Senior Management 67
Team            43
dtype: int64
```

We can see that every column has a different amount of missing values. Like Gender has 145 missing values and salary has 0. Now for handling these missing values there can be several cases like dropping the rows containing NaN or replacing NaN with either mean, median, mode, or some other value.

=====XXXXXX

## Types of Exploratory Data Analysis

There are three main types of EDA:

- Univariate
- Bivariate
- Multivariate

**In univariate analysis**, the output is a single variable and all data collected is for it. There is no cause-and-effect relationship at all. For example, data shows products produced each month for twelve months.

**In bivariate analysis**, the outcome is dependent on two variables, e.g., the age of an employee, while the relation with it is compared with two variables, i.e., his salary earned and expenses per month.

**In multivariate analysis**, the outcome is more than two, e.g., type of product and quantity sold against the product price, advertising expenses, and discounts offered. The analysis of data is done on variables that can be numerical or categorical. The result of the analysis can be represented in numerical values,

visualization, or graphical form. Accordingly, they could be further classified as non-graphical or graphical.

## 1. Univariate Non-Graphical

It is the simplest of all types of data analysis used in practice. As the name suggests, uni means only one variable is considered whose data (referred to as population) is compiled and studied. The main aim of univariate non-graphical EDA is to find out the details about the distribution of the population data and to know some specific parameters of statistics. The significant parameters which are estimated from a distribution point of view are as follows:

**Central Tendency:** This term refers to values located at the data's central position or middle zone. The three generally estimated parameters of central tendency are mean, median, and mode. Mean is the average of all values in data, while the mode is the value that occurs the maximum number of times. The Median is the middle value with equal observations to its left and right.

**Range:** The range is the difference between the maximum and minimum value in the data, thus indicating how much the data is away from the central value on the higher and lower side.

**Variance and Standard Deviation:** Two more useful parameters are standard deviation and variance. Variance is a measure of dispersion that indicates the spread of all data points in a data set. It is the measure of dispersion mostly used and is the mean squared difference between each data point and mean, while standard deviation is the square root value of it. **The larger the value of standard deviation, the farther the spread of data, while a low value indicates more values clustering near the mean.**

## 2. Univariate Graphical

The graphs in this section are based on Auto MPG dataset available on the UCI repository. Some common types of univariate graphics are:

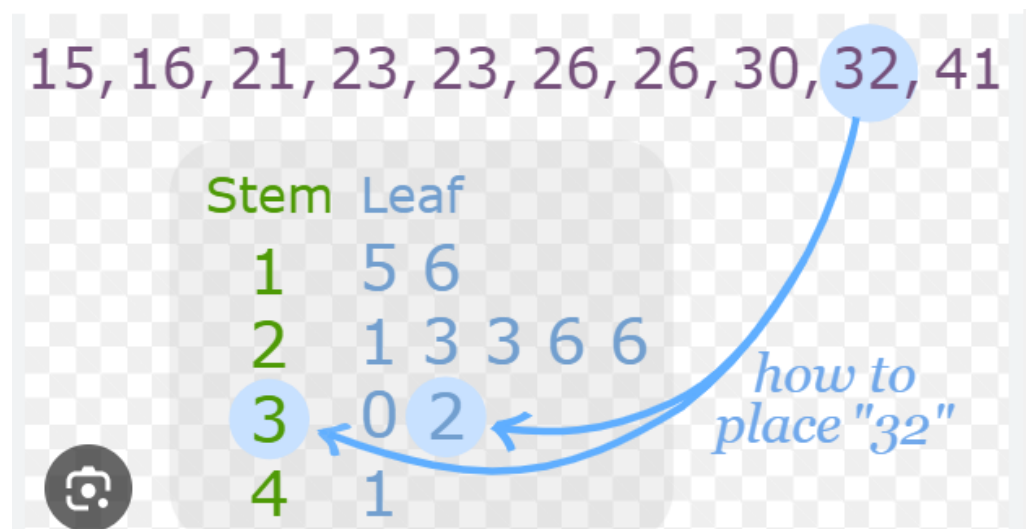
**Stem-and-leaf Plots:** This is a very simple but powerful EDA method used to display quantitative data but in a shortened format. It displays the values in the data set, keeping each observation intact but separating them as stem (the leading digits) and remaining or trailing digits as leaves. But histogram is mostly used in its place now.

A stem-and-leaf plot is a graphical representation used in statistics to display the distribution of a set of data. It is particularly useful for small to moderate-sized datasets. The plot separates each data point into a "stem" (the leading digit or digits) and a "leaf" (the trailing digit).

Here's how a stem-and-leaf plot is constructed:

- Separate the Data:
- Identify the leading digit (or digits) of each data point as the "stem."
- Identify the trailing digit of each data point as the "leaf."
- Order the Data:
- Arrange the data in ascending order.
- Create the Plot:
- Write the stems in a vertical column on the left side of the plot.
- Write the corresponding leaves to the right of their respective stems.

Example:



Histograms (Bar Charts): These plots are used to display both grouped or ungrouped data. On the x-axis, values of variables are plotted, while on the y-axis are the number of observations or frequencies. Histograms are very simple to quickly understand your data, which tell about values of data like central tendency, dispersion, outliers, etc. The simplest fundamental graph is a histogram, which is a



bar plot with each bar representing the frequency, i.e., the count or proportion (the ratio of count to the total count of occurrences) for various values.

There are many types of histograms, a few of which are listed below:

**Simple Bar Charts:** These are used to represent categorical variables with rectangular bars, where the different lengths correspond to the values of the variables.

**Multiple or Grouped charts:** Grouped bar charts are bar charts representing multiple sets of data items for comparison where a single color is used to denote one specific series in the dataset.

**Percentage Bar Charts:** These are bar graphs that depict the data in the form of percentages for each observation. The following image shows a percentage bar chart with dummy values.

**Box Plots:** These are used to display the distribution of quantitative value in the data. If the data set consists of categorical variables, the plots can show the comparison between them. Further, if outliers are present in the data, they can be easily identified. These graphs are very useful when comparisons are to be shown in percentages, like values in the 25 %, 50 %, and 75% range (quartiles).

### 3. Multivariate Non-Graphical

The multivariate non-graphical exploratory data analysis technique is usually used to show the connection between two or more variables with the help of either **cross-tabulation or statistics**.

For categorical data, an extension of tabulation called cross-tabulation is extremely useful. For two variables, cross-tabulation is preferred by making a two-way table with column headings that match the amount of one variable and row headings that match the amount of the opposite two variables, then filling the counts with all subjects that share an equivalent pair of levels.

For each categorical variable and one quantitative variable, we can generate statistical information for quantitative variables separately for every level of the specific variable. We then compare the statistics across the number of categorical variables.

### 4. Multivariate Graphical

Graphics are used in multivariate graphical data to show the relationships between two or more variables. Here the outcome depends on more than two variables, while the change-causing variables can also be multiple.

Some common types of multivariate graphics include:

#### A) Scatter Plot

The essential graphical EDA technique for two quantitative variables is the scatter plot, so one variable appears on the x-axis and the other on the y-axis and, therefore, the point for every case in your dataset. This can be used for bivariate analysis.

#### B) Multivariate Chart

A Multivariate chart is a type of control chart used to monitor two or more interrelated process variables. This is beneficial in situations such as process control, where engineers are likely to benefit from using multivariate charts. These charts allow monitoring multiple parameters together in a single chart. A notable advantage of using multivariate charts is that they help minimize the total number of control charts for organizational processes. Pair plots generated using the Seaborn library are a good example of multivariate charts as they help visualize the relationships between all numerical variables in the entire dataset at once.

#### C) Run Chart

A run chart is a data line chart drawn over time. In other words, a run chart visually illustrates the process performance or data values in a time sequence. Rather than summary statistics, seeing data across time yields a more accurate conclusion. A trend chart or time series plot is another name for a run chart. The plot below depicts dummy values of sales over a period of time.

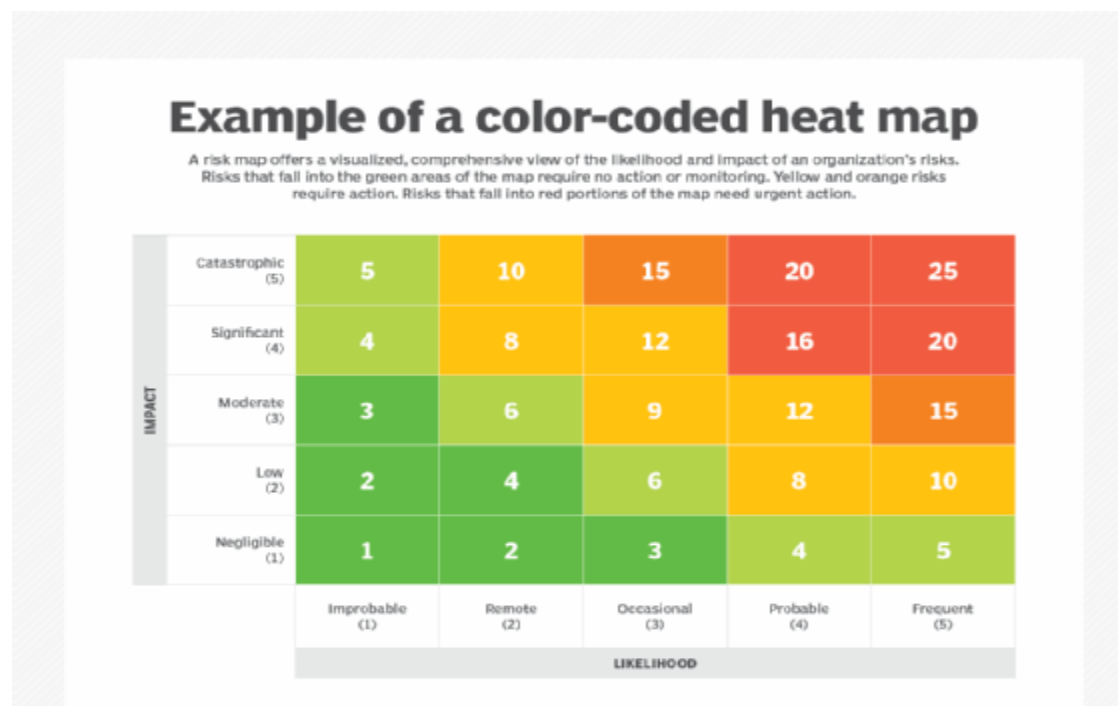
#### D) Bubble Chart

Bubble charts scatter plots that display multiple circles (bubbles) in a two-dimensional plot. These are used to assess the relationships between three or more numeric variables. In a bubble chart, every single dot corresponds to one data point, and the values of the variables for each point are indicated by different positions such as horizontal, vertical, dot size, and dot colors.

## E) Heat Map

A heat map is a colored graphical representation of multivariate data structured as **a matrix of columns and rows**. The heat map transforms the correlation matrix into color coding and represents these coefficients to visualize the strength of correlation among variables. It assists in finding the best features suitable for building accurate Machine Learning models.

Apart from the above, there is also the 'Classification or Clustering analysis' technique used in EDA. It is an unsupervised type of machine learning used for the classification of input data into specified categories or clusters exhibiting similar characteristics in various groups. This can be further used to draw important interpretations in EDA.



## Exploratory Data Analysis Tools

### 1. Python

Python is used for different tasks in EDA, such as finding missing values in data collection, data description, handling outliers, obtaining insights through charts, etc. The syntax for EDA libraries like Matplotlib, Pandas, Seaborn, NumPy, Altair,

and more in Python is fairly simple and easy to use for beginners. You can find many open-source packages in Python, such as D-Tale, AutoViz, PandasProfiling, etc., that can automate the entire exploratory data analysis process and save time.

## 2. R

R programming language is a regularly used option to make statistical observations and analyze data, i.e., perform detailed EDA by data scientists and statisticians. Like Python, R is also an **open-source programming** language suitable for statistical computing and graphics. Apart from the commonly used libraries like ggplot, Leaflet, and Lattice, there are several powerful R libraries for automated EDA, such as Data Explorer, SmartEDA, GGally, etc.

## 3. MATLAB

MATLAB is a well-known commercial tool among engineers since it has a very strong mathematical calculation ability. Due to this, it is possible to use MATLAB for EDA but it requires some basic knowledge of the MATLAB programming language.

## **Advantages of Using EDA**

Here are a few advantages of using Exploratory Data Analysis -

### **1. Gain Insights Into Underlying Trends and Patterns**

EDA assists data analysts in identifying crucial trends quickly through data visualizations using various graphs, such as box plots and histograms. Businesses also expect to make some unexpected discoveries in the data while performing EDA, which can help improve certain existing business strategies.

### **2. Improved Understanding of Variables**

Data analysts can significantly improve their comprehension of many factors related to the dataset. Using EDA, they can extract various information such as averages, means, minimum and maximum, and more such information is required for preprocessing the data appropriately.

### **3. Better Preprocess Data to Save Time**

EDA can assist data analysts in identifying significant mistakes, abnormalities, or missing values in the existing dataset. Handling the above entities is critical for any organization before beginning a full study as it ensures correct preprocessing of data and may help save a significant amount of time by avoiding mistakes later when applying machine learning models.

#### 4. Make Data-driven Decisions

The most significant advantage of employing EDA in an organization is that it helps businesses to improve their understanding of data. With EDA, they can use the available tools to extract critical insights and make conclusions, which assist in making decisions based on the insights from the EDA.

=====

Exploratory Data Analysis (EDA) often involves assessing the similarity or dissimilarity between data points or features. Various statistical and visual methods can be used for this purpose.

##### **Descriptive Statistics:**

Measure of Central Tendency: Mean, median, and mode can provide insights into the central location of the data.

Measure of Dispersion: Range, variance, and standard deviation can help understand how spread out the values are.

##### **Correlation:**

Pearson Correlation Coefficient: Measures the linear relationship between two continuous variables. It ranges from -1 to 1, where -1 indicates a perfect negative linear relationship, 1 indicates a perfect positive linear relationship, and 0 indicates no linear relationship.

##### **Distance Metrics:**

**Euclidean Distance:** Measures the straight-line distance between two points in space. It is commonly used when dealing with numerical data.

**Manhattan Distance (L1 Norm):** Sum of the absolute differences between corresponding coordinates. It is often used when dealing with categorical data.

**Cosine Similarity:** Measures the cosine of the angle between two non-zero vectors. It is commonly used in text mining and collaborative filtering.

### **Clustering Analysis:**

Techniques like k-means clustering or hierarchical clustering group similar data points together.

### **Visualization:**

Scatter Plots: Visualize the relationship between two variables.

Heatmaps: Visualize the correlation matrix of variables.

Pair Plots: Display scatter plots for pairs of variables in a dataset.

### **Principal Component Analysis (PCA):**

Reduces the dimensionality of the data while retaining most of its variability. It can be used to identify patterns and similarities in high-dimensional data.

### **Data Profiling:**

Summary Statistics: Provide a quick overview of the data, including mean, median, mode, and quartiles.

Value Counts: Count the occurrences of unique values in categorical variables.

### **Statistical Tests:**

T-tests, ANOVA, etc.: These tests can be used to compare means across different groups or conditions.

=====

**Why Preprocessing? Data Cleaning; Data Integration; Data Reduction: Attribute subset selection, Histograms, Clustering and Sampling; Data Transformation & Data Discretization: Normalization, Binning, Histogram Analysis and Concept hierarchy generation**

## **Why Preprocessing is required in data science:**

Preprocessing is a crucial step in the data science pipeline, and it involves cleaning and transforming raw data into a format that is suitable for analysis.

There are several reasons why preprocessing is required in data science:

### **Handling Missing Data:**

Real-world datasets often contain missing values, which can lead to biased or inaccurate results. Preprocessing involves techniques for handling missing data, such as imputation (replacing missing values with estimated ones) or removal of incomplete records.

### **Dealing with Outliers:**

Outliers, or data points significantly different from the rest of the dataset, can skew analysis and model performance. Preprocessing helps identify and handle outliers through methods like trimming, transformation, or imputation.

### **Data Cleaning:**

Raw data may contain errors, inconsistencies, or noise. Cleaning involves removing duplicates, correcting errors, and standardizing formats to ensure data quality.

### **Normalization and Scaling:**

Different features in a dataset may have different scales, which can affect the performance of certain machine learning algorithms. Normalization and scaling techniques ensure that all features contribute equally to the analysis and prevent dominance by features with larger scales.

### **Encoding Categorical Variables:**

Machine learning models often require numerical input, so categorical variables need to be converted into a numerical format through techniques like one-hot encoding or label encoding.

### **Feature Engineering:**

Creating new features or modifying existing ones can enhance the predictive power of models. Feature engineering involves transforming variables, creating interaction terms, or extracting relevant information from existing features.

### **Handling Imbalanced Datasets:**

In classification tasks, imbalanced datasets (where one class is significantly more prevalent than others) can lead to biased models. Preprocessing techniques, such as oversampling or undersampling, help balance class distributions.

### **Text Preprocessing:**

In natural language processing (NLP), text data often requires special preprocessing steps like tokenization, stemming, and removing stop words to extract meaningful information.

### **Dimensionality Reduction:**

High-dimensional datasets may suffer from the curse of dimensionality, leading to increased computational complexity and potential overfitting. Techniques like Principal Component Analysis (PCA) or feature selection help reduce dimensionality while retaining important information.

### **Data Standardization:**

Standardizing data ensures that different variables are on the same scale, making it easier to compare and analyze them. This is particularly important for algorithms that rely on distance metrics.

### **Handling Skewed Distributions:**

Some machine learning algorithms assume that the data is normally distributed. Preprocessing techniques, such as log transformation, can help address skewed distributions and improve model performance.



=====

### **A case study on an online e-commerce dataset:**

Creating a case study on an online e-commerce dataset involves describing a hypothetical scenario, the dataset used, and the analysis or insights gained. Following example of a case study on an online e-commerce dataset:

*Note: In the context of scientific research and data analysis, a hypothesis is a specific, testable proposition or prediction about the relationship between variables or the outcome of a scientific study. A hypothesis is a tentative explanation for an observed phenomenon or a statement that can be tested through empirical research.*

### **Case Study: Understanding Customer Behavior in an Online E-Commerce Platform**

#### **Background:**

Company XYZ is an online e-commerce platform that sells a variety of products, including electronics, clothing, and home goods. The company has collected data on customer transactions, website interactions, and product reviews over the past year. The objective is to gain insights into customer behavior and improve the overall customer experience.

#### **Dataset:**

The dataset includes the following key information:

##### **Customer Demographics:**

Age

Gender

Location

##### **Transaction Data:**

Date and time of purchase

Products purchased

Transaction amount

Website Interactions:

Page views

Time spent on the website

Click-through rates on product pages

### **Product Reviews:**

Ratings given by customers

Textual reviews

### **Analysis:**

#### **Customer Segmentation:**

Analyze customer demographics to identify key segments.

Explore purchasing patterns among different customer groups.

#### **Purchase Behavior:**

Examine transaction data to identify popular products.

Analyze the average transaction amount and frequency of purchases.

#### **Website Engagement:**

Investigate website interaction data to understand which pages are most visited.

Evaluate the effectiveness of product pages and identify potential improvements.

#### **Customer Satisfaction:**

Analyze product reviews to assess customer satisfaction.

Identify common themes or issues raised in customer feedback.

### **Sales Forecasting:**

Use historical transaction data to forecast future sales.

Identify trends and seasonality in purchasing behavior.

### **Recommendation System:**

Implement a recommendation system based on customer preferences and past purchases.

Evaluate the impact of the recommendation system on sales.

### **Insights and Recommendations:**

#### **Targeted Marketing:**

Tailor marketing campaigns to specific customer segments based on demographics and purchasing behavior.

#### **Website Optimization:**

Improve website navigation and the user interface based on insights from website engagement data.

#### **Product Recommendations:**

Enhance the recommendation system to provide more accurate and personalized suggestions to customers.

#### **Customer Retention:**

Implement strategies to address issues raised in customer reviews to improve overall satisfaction and retention.

#### **Seasonal Promotions:**

Plan targeted promotions and discounts based on identified trends and seasonality in sales.

## Conclusion:

- By leveraging the insights gained from the e-commerce dataset, Company XYZ can make informed decisions to enhance the customer experience, optimize its marketing strategies, and ultimately increase customer satisfaction and sales.
- This case study provides a framework for exploring an e-commerce dataset and deriving actionable insights to improve business outcomes.
- Depending on the specifics of the dataset and the company's goals, the analysis and recommendations can be customized accordingly.

=====XXXXXXXXXX=====

## Numerical:

The Euclidean distance between these two points can be calculated using the formula:

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

Suppose we have the following two data points:

Point A: (2, 4, 5)

Point B: (1, 7, 9)

Now, let's calculate the Euclidean distance:

$$d = \sqrt{(1 - 2)^2 + (7 - 4)^2 + (9 - 5)^2}$$

$$d = \sqrt{(-1)^2 + 3^2 + 4^2}$$

$$d = \sqrt{1 + 9 + 16}$$

$$d = \sqrt{26}$$



In real-world applications, Euclidean distance is frequently used in clustering, classification, and similarity analysis within the field of data science.

=====XXXXXX=====

## Data transformation

- Data transformation in data science involves modifying, converting, or restructuring raw data to make it more suitable for analysis, modeling, or visualization.
- The primary goal of data transformation is to enhance the quality and relevance of the data, making it more meaningful and valuable for the intended analytical tasks.
- There are various techniques and methods used for data transformation, depending on the characteristics of the data and the specific requirements of the analysis.

Here are some common aspects of data transformation:

Handling Missing Data:

Imputation: Filling in missing values with estimated or calculated values, such as mean, median, or interpolation, to ensure completeness in the dataset.

### Normalization and Scaling:

Scaling numerical features to a common scale. Common methods include Min-Max scaling (scaling values between 0 and 1) and Z-score normalization (scaling to have zero mean and unit variance). This ensures that all features contribute equally to the analysis, especially in machine learning models that rely on distance metrics.

### Handling Categorical Data:

One-Hot Encoding: Converting categorical variables into binary vectors, where each category becomes a binary feature.

Label Encoding: Assigning numerical labels to categories. It is suitable for ordinal categorical variables.

### Log Transformation:

Applying a logarithmic transformation to data, often used to stabilize variance and make the data more symmetric. This is particularly useful for data with skewed distributions.

### Box-Cox Transformation:

A family of power transformations that can be applied to stabilize variance and make the data more normal. It is suitable for data with heteroscedasticity and non-constant variance.

### Handling Outliers:

Truncation or winsorization: Capping extreme values to reduce the impact of outliers on analysis or modeling. Transformation (e.g., log transformation) can also mitigate the influence of outliers.

Dummy Variable Creation:

Creating dummy variables from categorical variables for use in statistical models, allowing the inclusion of categorical data in regression models.

Binning or Discretization:

Grouping continuous variables into intervals or bins, which can simplify the representation of the data. It is particularly useful for certain types of analyses, such as decision tree models.

### **Time Series Decomposition:**

Decomposing time series data into its trend, seasonality, and residual components, making it easier to analyze and model.

### **Aggregation and Grouping:**

Aggregating data by grouping and summarizing information based on specific attributes, such as calculating mean, sum, or other summary statistics for subsets of the data.

### **Feature Engineering:**

Creating new features or modifying existing ones to extract more valuable information for modeling. This might involve mathematical transformations, interaction terms, or creating derived variables.

=====XXXX=====

**Data reduction** in data science refers to the process of reducing the volume but producing the same or similar analytical results. This can involve reducing the dimensionality of the dataset, compressing the data, or summarizing it in a way that retains essential information while making it more manageable for analysis or storage. The primary goals of data reduction include improving efficiency, speeding up analysis, and mitigating the challenges associated with large and complex datasets.

There are several techniques and methods used for data reduction:

### **Dimensionality Reduction:**

**Principal Component Analysis (PCA):** PCA is a technique that transforms high-dimensional data into a lower-dimensional space while retaining the most important information. It identifies the principal components (linear combinations of the original features) that capture the maximum variance in the data.

**Feature Selection:**

Selecting a subset of the most relevant features is a form of data reduction. Feature selection methods aim to identify and retain the most informative variables while discarding less important ones.

**Data Aggregation:**

Aggregating data involves combining multiple data points into summary statistics or representative values. For example, converting daily data to monthly averages or summing values across different categories can reduce the dataset's size.

**Binning and Discretization:**

Binning involves grouping continuous data into intervals or bins, which can reduce the precision of the data but simplify its representation. Discretization is particularly useful in certain types of analyses and can be a form of data reduction.

**Sampling:**

Taking a subset or sample of the original data is a common form of data reduction. Random sampling or systematic sampling methods can be employed to create a smaller but representative sample.

**Data Compression:**

Data compression techniques aim to represent the data in a more compact form, reducing the number of bits needed to represent the information. This is often used in storage and transmission of data.

**Clustering:**



Clustering techniques group similar data points together, allowing for the representation of a dataset using fewer cluster representatives. This can reduce the size of the dataset while preserving essential patterns.

#### Summary Statistics:

Calculating summary statistics such as mean, median, standard deviation, or percentiles can provide a concise representation of the data, especially for large datasets.

- Data reduction is particularly important when dealing with big data, where the sheer volume of information can pose computational challenges.
- By reducing data complexity, data scientists can expedite analysis, save computational resources, and often achieve similar insights with a smaller, more manageable dataset.
- However, it's crucial to carefully choose the appropriate data reduction techniques based on the specific requirements of the analysis and the characteristics of the data.

=====XXXXXXXX=====

### Data discretization

**Data discretization** in data science refers to the process of converting continuous data into discrete intervals or bins. This transformation is useful in various scenarios, especially when dealing with numerical data, and it is applied to make the data more manageable, understandable, and suitable for certain types of analyses or modeling techniques. Discretization is often used in feature engineering and preprocessing stages to handle specific challenges associated with continuous data.

Here are some common methods of data discretization:

#### Equal-Width Binning:

Divide the range of values into equal-width intervals. This method is simple, but it may not be suitable for datasets with unevenly distributed values.

**Example:**

Input=[5,10,50,72,92,104,215], Bins:3

## Equal Width Binning

Range Of Each Bin :

[min + W]

[min + 2W]

.....

[min + nW]

n : Number Of Bins

$$W = (\text{max}-\text{min})/\text{Number Of Bins}$$

## Equal Width Binning ⓘ

**Input:** [5,10,50,72,92,104,215] , **Bins:** 3

$$W = (215-5)/3 = 70$$

Bin 1 Range = [5+70] = 75

Bin 2 Range = [5+2(70)] = 145

Bin 3 Range = [5+3(70)] = 215

B1=[5,10,50,72]

B2=[92,104]

B3=[215]

=====XXX=====

### Equal-Frequency Binning:

Divide the data into intervals such that each interval contains approximately the same number of data points. This approach can be more robust to unevenly distributed data but may result in uneven ranges.

Example: Bin=3, Total elements presents =6.

## Equal Frequency Binning

**Input:** [5,10,50,72,92,104]

Number of Values/Number of Bins:

**$6/3=2$  Numbers per Bin**

# Equal Frequency Binning

**Input:** [5,10,50,72,92,104]

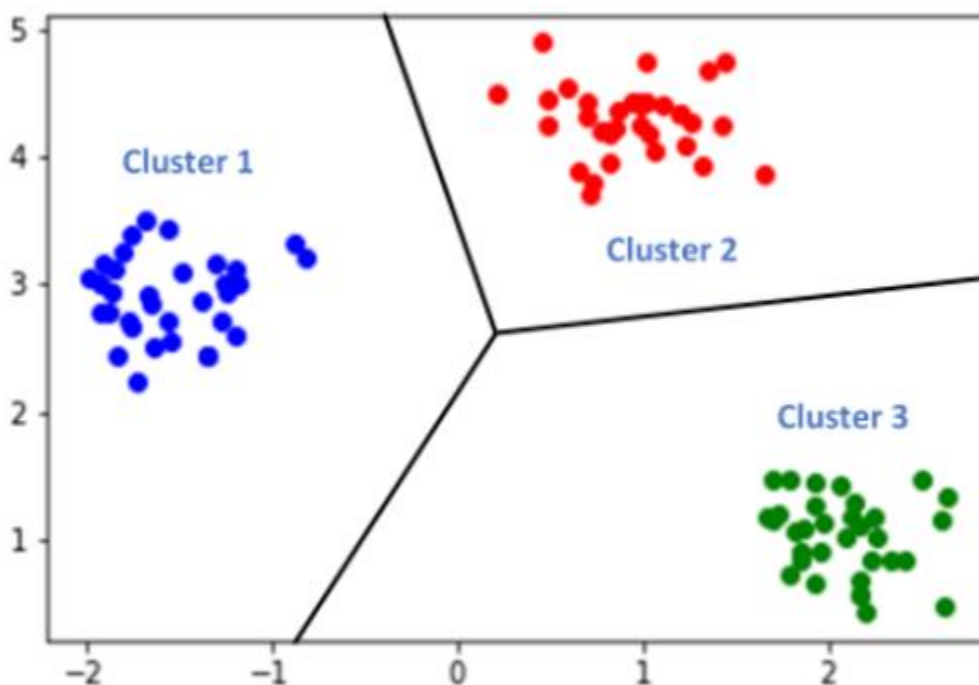
**Bin 1 :** [5,10]

**Bin 2 :** [50,72]

**Bin 3 :** [92,104]

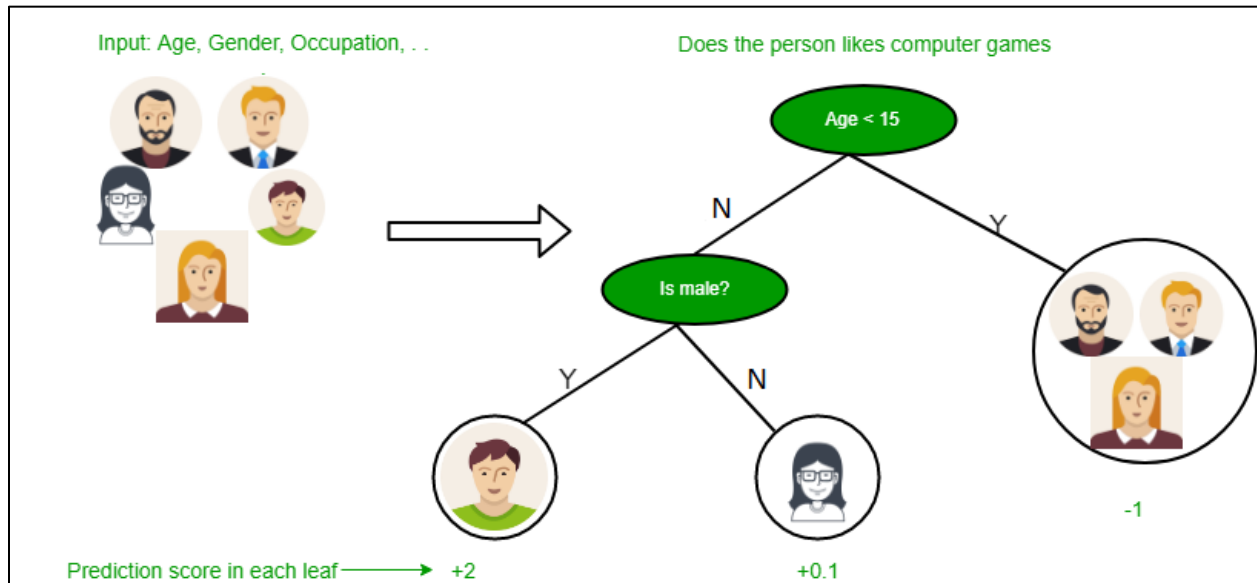
## Clustering-Based Discretization:

Use clustering algorithms to group similar data points into clusters and assign a discrete label to each cluster. The clusters effectively become the bins for the discretized data.



## Decision Tree-Based Discretization:

Utilize decision tree algorithms to identify optimal split points in the continuous data, creating intervals based on these splits. This approach is often part of the decision tree modeling process.



## Entropy-Based Discretization:

Apply information theory concepts, such as entropy, to find the most informative split points for discretization. This method aims to maximize the information gain in the discretized categories.

## Entropy-based Discretization

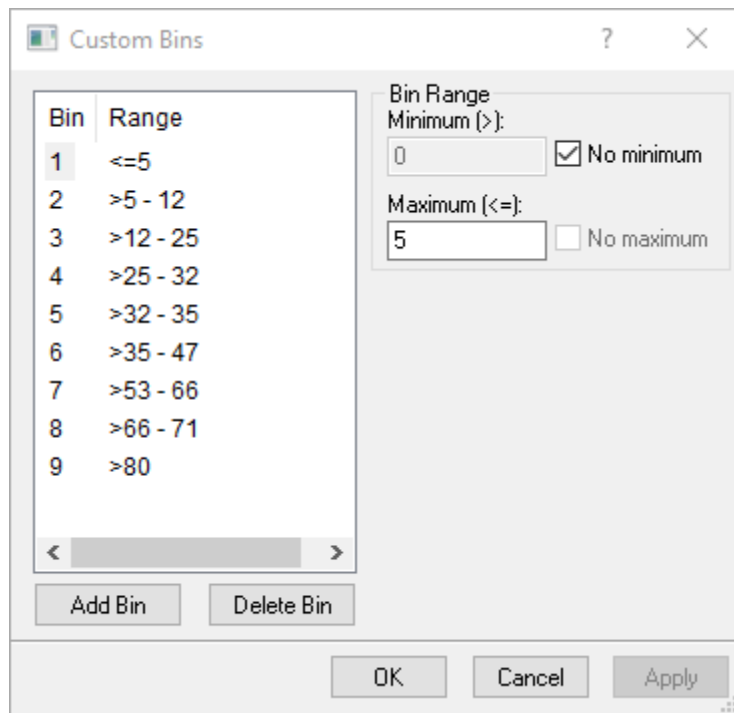
- A widely-used supervised discretization method
- Entropy is a measure of impurity
  - Higher entropy implies data points are from a large number of classes (heterogeneous)
  - Lower entropy implies most of the data points are from the same class

$$Entropy = -\sum_j p_j \log p_j$$

Where  $p_j$  is the proportion of data points belonging to class  $j$

### Custom Binning:

Define custom intervals based on domain knowledge or specific requirements. This approach allows for the creation of bins that align with the significance of values in the context of the problem.



Data discretization is often employed when working with machine learning algorithms that assume categorical or ordinal input. It can also be beneficial for simplifying complex datasets, reducing noise, and improving the interpretability of models. However, it's important to choose an appropriate discretization method based on the characteristics of the data and the goals of the analysis.

=====XXXXXXXXXX=====

## Numerical

Jaccard distance is commonly used to measure the dissimilarity between two sets. It is defined as the size of the difference between two sets divided by the size of their union. The formula for Jaccard distance (J) is given by

The formula for Jaccard distance (J) is given by:

$$J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Let's consider an example where we have two sets, A and B:

$$A = \{1, 2, 3, 4\}$$

$$B = \{3, 4, 5, 6\}$$

To calculate the Jaccard distance:

1. Find the intersection of sets A and B (elements common to both sets):  $A \cap B = \{3, 4\}$
2. Find the union of sets A and B (all unique elements from both sets):  $A \cup B = \{1, 2, 3, 4, 5, 6\}$
3. Calculate the Jaccard distance using the formula:

$$J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

$$J(A, B) = 1 - \frac{2}{6}$$

$$J(A, B) = 1 - \frac{1}{3}$$

$$J(A, B) = \frac{2}{3}$$

So, the Jaccard distance between sets A and B is  $\frac{2}{3}$  or approximately 0.6667 when rounded to four decimal places.

In practical terms, Jaccard distance is often used in text analysis, recommendation systems, and other applications where the similarity between two sets needs to be measured. The result ranges from 0 (identical sets) to 1 (completely dissimilar sets).

=====**xENDx**=====