



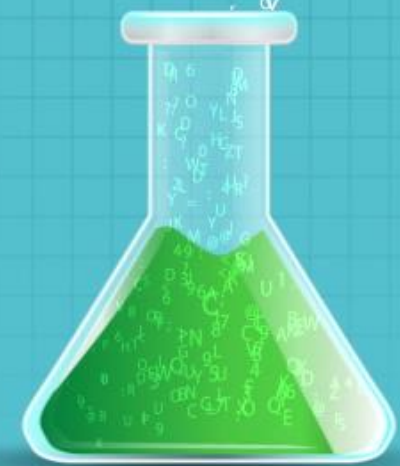
Data Science with Python

Lesson 2 – Data Analytics Overview

DATA
SCIENCE

What's In It For Me

- Data Analytics process and its steps
- Skills and tools required for Data Analysis
- Challenges of the Data Analytics Process
- Exploratory Data Analysis technique
- Data visualization techniques
- Hypothesis testing to analyze data



Why Data Analytics

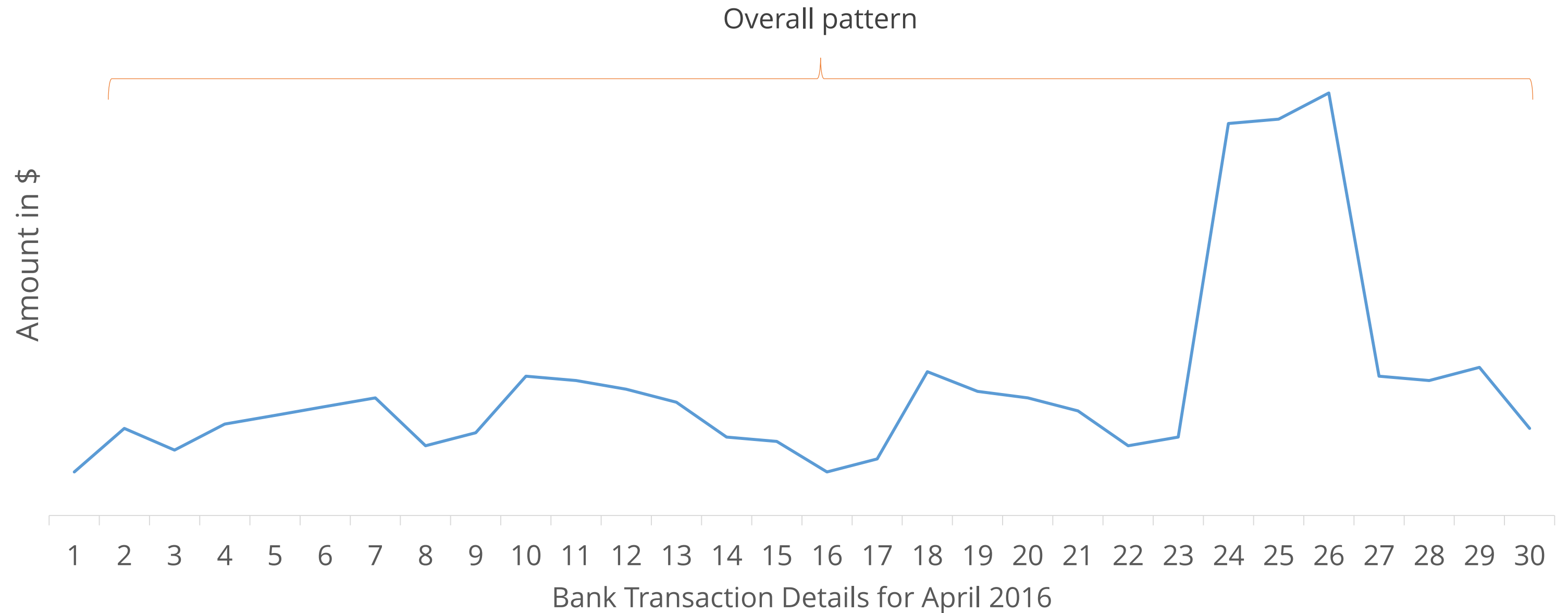
Data by itself is just an information source. But unless you can understand it, you will not be able to use it effectively.

Date	Description	Deposit	Withdrawal	Balance
Apr 1	ATM Post Debit		100	\$200,000
Apr 2	Paypal Tranfer 231054	200		\$202,000
Apr 3	Simplilearn course fee		150	\$200,500
Apr 4	Starluck Café		210	\$198,400
Apr 5	Walcart TX		230	\$196,100
Apr 6	ebuy swiss watch 239		250	\$193,600
Apr 7	Caterpallor black boots men		270	\$190,900
Apr 8	Halo blue shirt 831		160	\$189,300

Information source;
overall patterns not
clearly visible

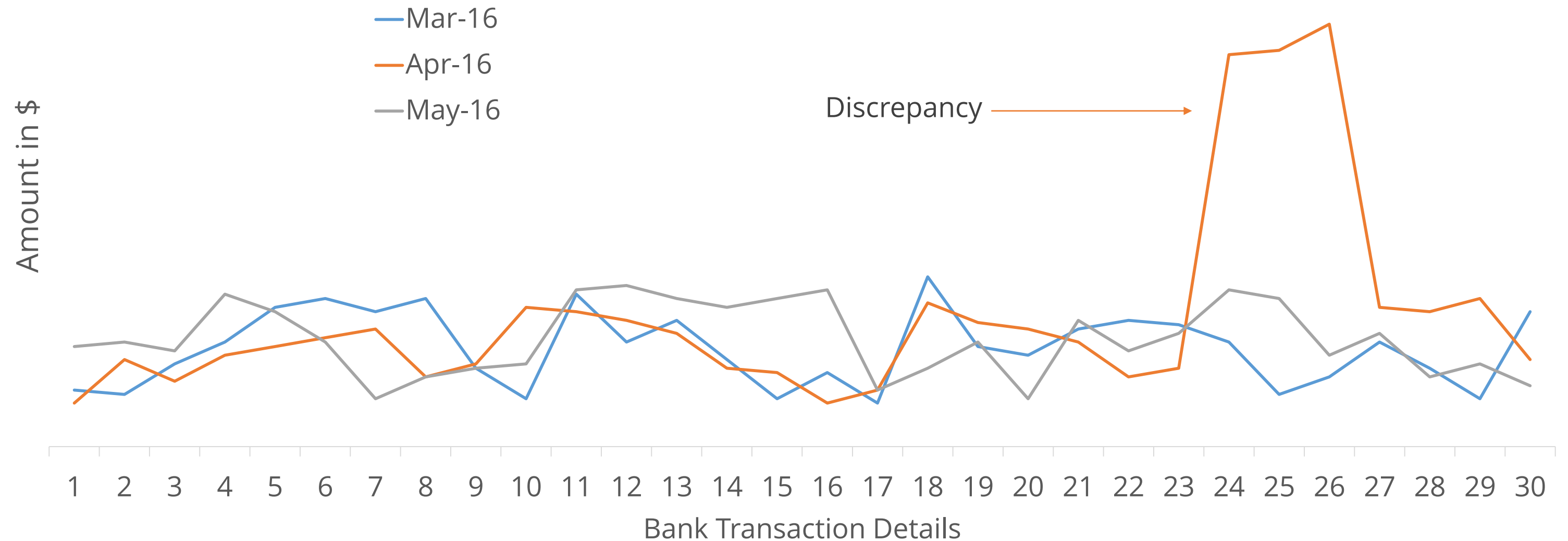
Why Data Analytics (contd.)

When the transaction details are presented as a line chart, the deposit and withdrawal patterns become apparent.



Why Data Analytics (contd.)

When the transaction details are presented as a line chart, the deposit and withdrawal patterns become apparent. It helps view and analyze general trends and discrepancies.



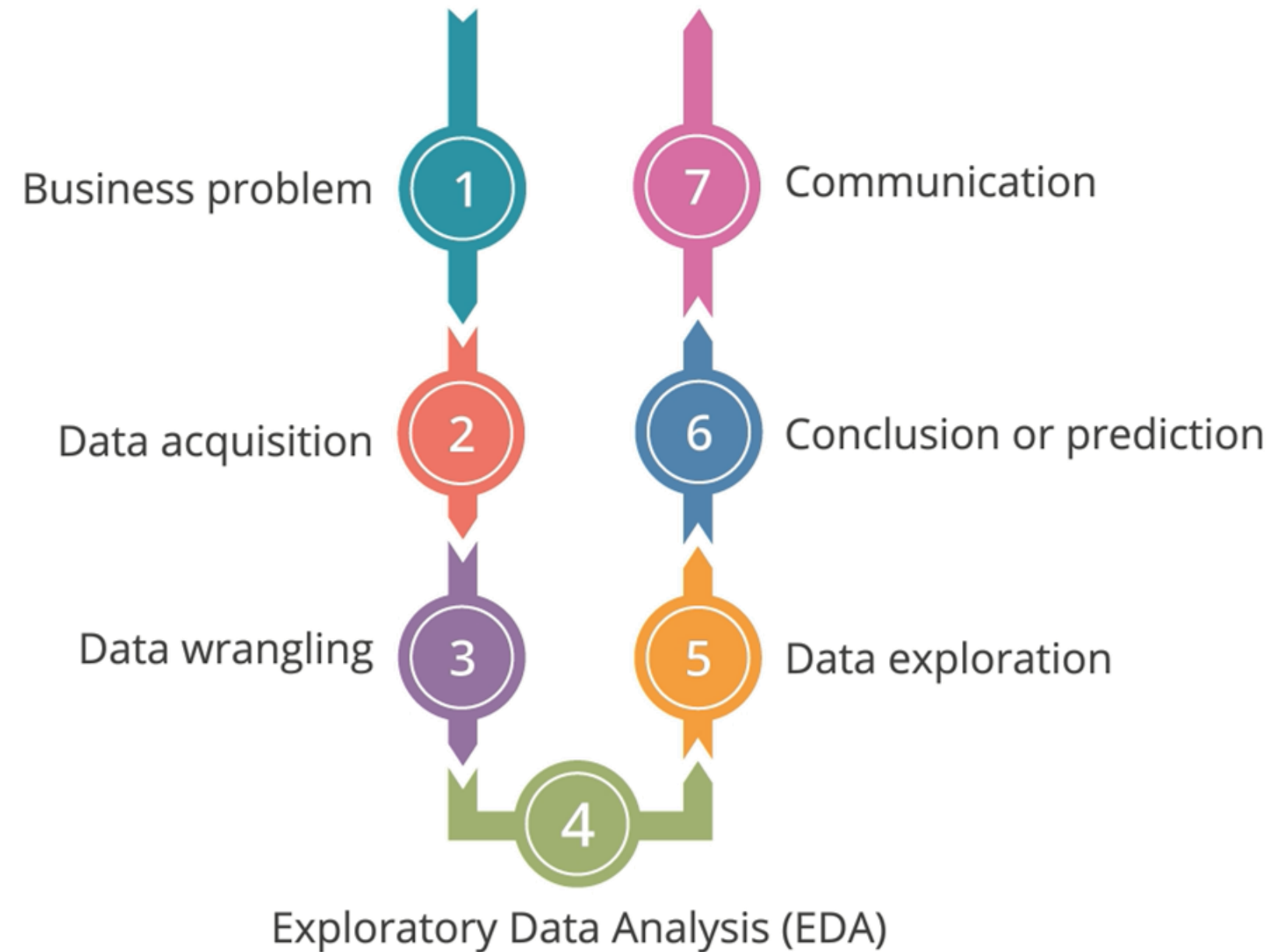
Introduction to Data Analytics

Data Analytics is a combination of processes to extract information from datasets.



Introduction to Data Analytics

Data Analytics is a combination of processes to extract information from datasets.



Business Problem

The process of analytics begins with questions or business problems of stakeholders.



Data Acquisition

Collect data from various sources for analysis to answer the question raised in step 1.



Data Scientist Expertise:

- File handling
- File formats
- Web scraping



Twitter, Facebook, LinkedIn, and other social media and information sites provide streaming APIs.



Server logs can be extracted from enterprise system servers to analyze and optimize application performance.

Data Wrangling and Exploration

Data wrangling is the most important phase of the data analytic process.



Data cleansing



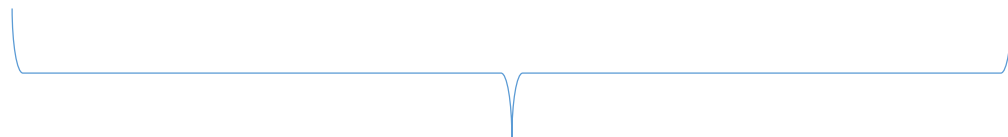
Data manipulation



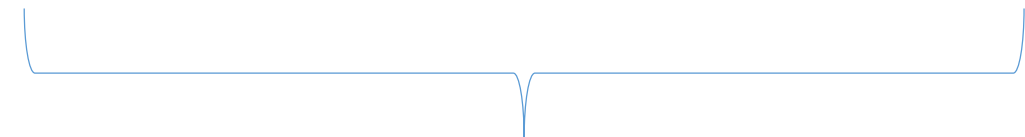
Data discovery



Data pattern



Data Wrangling



Data Exploration

Data Wrangling—Challenges

This phase includes data cleansing, data manipulation, data aggregation, data split, and reshaping of data.



Causes of challenges in the data wrangling phase:

- Unexpected data format
- Erroneous data
- Voluminous data to be manipulated
- Classifying data into linear or clustered
- Determining relationship between observation, feature, and response



Data wrangling is the most challenging phase and takes up 70% of the data scientist's time.

Data Exploration—Model Selection

This phase includes data cleansing, data manipulation, data aggregation, data split, and reshaping of data.



Model selection

- Based on the overall data analysis process
- Should be accurate to avoid iterations
- Depends on pattern identification and algorithms
- Depends on hypothesis building and testing
- Leads to building mathematical statistical functions

Exploratory Data Analysis (EDA)

Let's take a look at the exploratory data analysis phase.



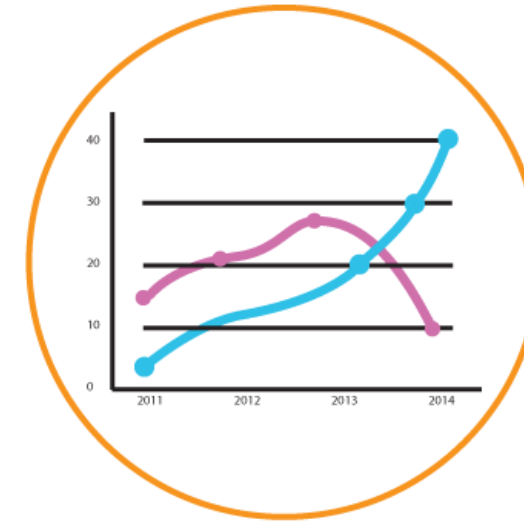
APPROACH

EDA approach studies the data to recommend suitable models that best fit the data.



FOCUS

The focus is on data; its structure, outliers, and models suggested by the data.



ASSUMPTIONS

EDA techniques make minimal or no assumptions. They present and show all the underlying data without any data loss.



EDA TECHNIQUES

Quantitative: Provides numeric outputs for the inputted data
Graphical: Uses statistical functions for graphical output

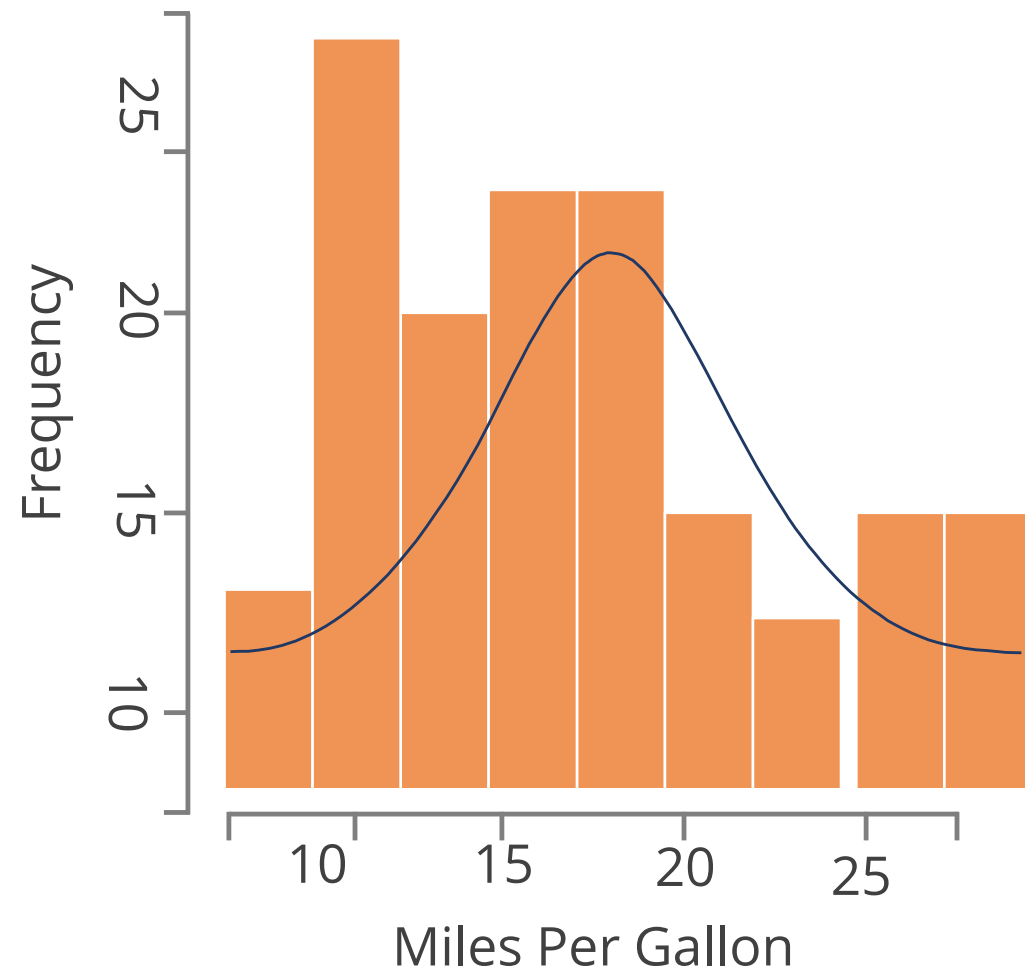
EDA— Quantitative Technique

EDA – Quantitative technique has two goals, measurement of central tendency and spread of data.

	Measurement of Central Tendency
Mean	Mean is the point which indicates how centralized the data points are. <ul style="list-style-type: none">• Suitable for symmetric distributions
Median	Median is the exact middle value. <ul style="list-style-type: none">• Suitable for skewed distributions and for catching outliers in the dataset
Mode	Mode is the most common value in the data (frequency).
	Measurement of Spread
Variance	Variance is approximately the mean of the squares of the deviations.
Standard Deviation	Standard deviation is the square root of the variance.
Inter-quartile Range	Inter-quartile range is the distance between the 75 th and 25 th percentile. It's essentially the middle 50% of the data.

EDA – Graphical Technique

Histograms and Scatter Plots are two popular graphical techniques to depict data.



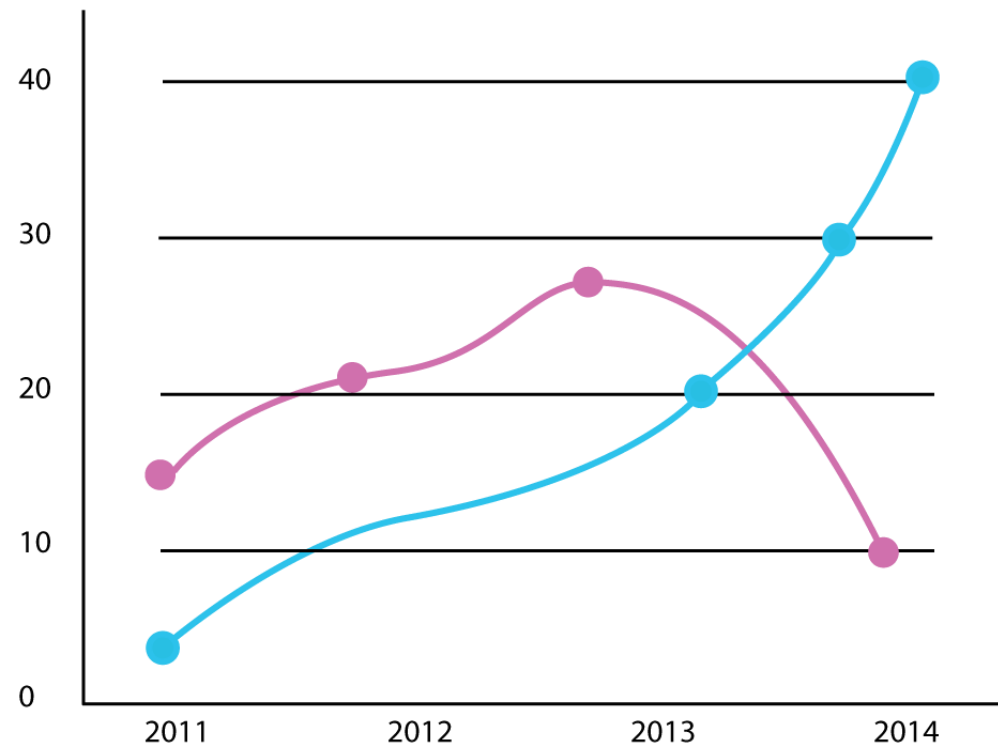
Histogram graphically summarizes the distribution of a univariate dataset.

It shows:

- the center or location of data (mean, median, or mode)
- the spread of data
- the skewness of data
- the presence of outliers
- the presence of multiple modes in the data

EDA – Graphical Technique

Histograms and Scatter Plots are two popular graphical techniques to depict data.



A Scatter plot represents relationships between two variables. It can answer these questions visually:

- Are variables X and Y related?
- Are variables X and Y linearly related?
- Are variables X and Y non-linearly related?
- Does change in variation of Y depend on X?
- Are there outliers?



Knowledge Check

KNOWLEDGE
CHECK

What is the goal of data acquisition?

Select all that apply.

- a. Collect data from various data sources
- b. Answer business questions through graphics
- c. Collect web server logs
- d. Scrape the web through web APIs



KNOWLEDGE
CHECK

What is the goal of data acquisition?
Select all that apply.

- a. Collect data from various data sources
- b. Answer business questions through graphics
- c. Collect web server logs
- d. Scrape web through web APIs



The correct answer is **a, c, d**

Explanation: Data acquisition is a process to collect data from various data sources **such as** RDBMS, No SQL databases, web server logs and also **scrape** the web through web APIs.

KNOWLEDGE
CHECK

What is the Exploratory Data Analysis technique?

Select all that apply.

- a. Analysis of data using quantitative techniques
- b. Conducted only on a small subset of data
- c. Analysis of data using graphical techniques
- d. Suggests admissible models that best fit the data



KNOWLEDGE
CHECK

What is the Exploratory Data Analysis technique?
Select all that apply.

- a. Analysis of data using quantitative techniques
- b. Conducted only on small subset of data
- c. Analysis of data using graphical techniques
- d. Suggests models that best fit the data



The correct answer is **a, c, d**.

Explanation: Most EDA techniques are graphical in nature with a few quantitative techniques and also suggest models that best fit the data. They use almost the entire data with minimum and no assumptions.

Conclusion or Predictions

This step involves reaching a conclusion and making predictions based on the data analysis.

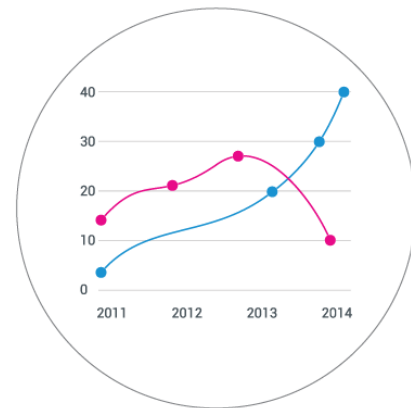


This phase:

- Involves heavy use of mathematical and statistical functions
- Requires model selection, training, and testing to help in forecasting
- Is called “machine learning” as data analysis is fully or semi-automated with minimal or no human intervention

Meaning of Hypothesis

Hypothesis is used to establish the relationship between dependent and independent variables.



Data Exploration Stage

Hypothesis building begins in the data exploration stage but becomes more mature in the conclusion or prediction phase.



Conclusion and Prediction

Key Considerations of Hypothesis Building

Testable explanations of a problem or observation

Used in quantitative and qualitative analyses to provide research solutions

Involves two variables, one dependent on another

Independent variable manipulated by the researcher

Dependent variable changes when the independent variable changes

Hypothesis Building Using Feature Engineering

Domain knowledge leads to hypothesis building using feature engineering.



Feature engineering involves domain expertise to:

- Make sense of data
- Construct new features from raw data automatically
- Construct new features from raw data manually

Hypothesis Building Using a Model

There are three phases to hypothesis building which are model building, model evaluation, and model deployment.

Phase 1: Model Building

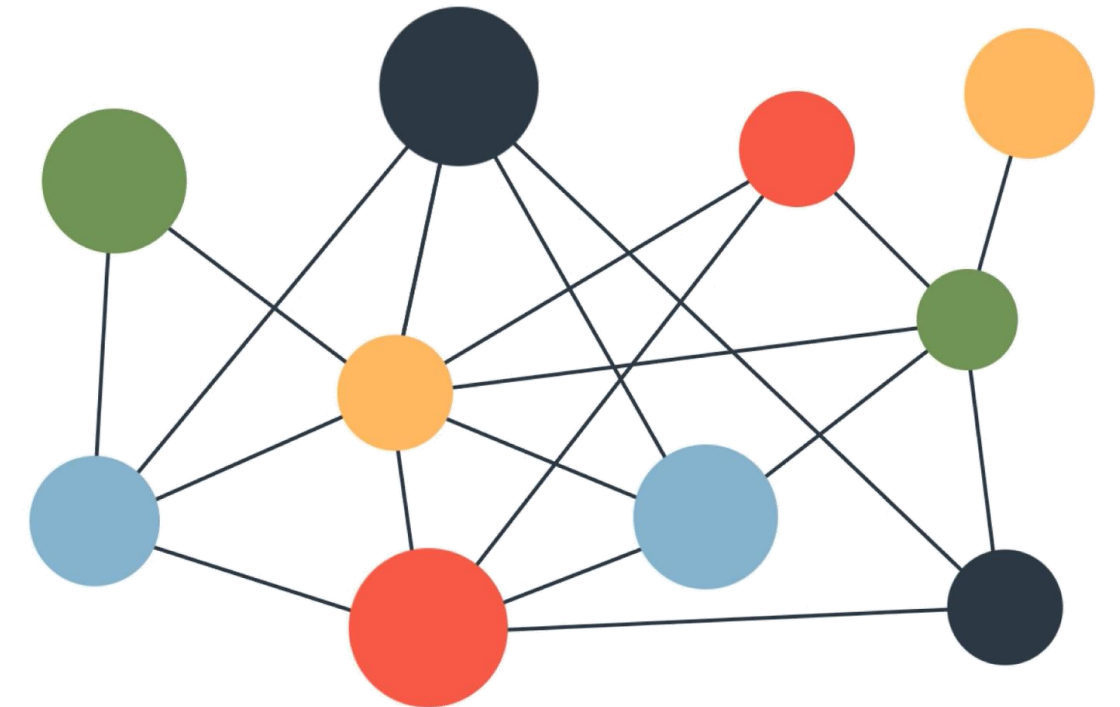
- Identify best input variables
- Evaluate the model's capacity to forecast with these variables

Phase 2: Model Evaluation

- Train and test the model for accuracy
- Optimize model accuracy, performance, and comparisons with other models

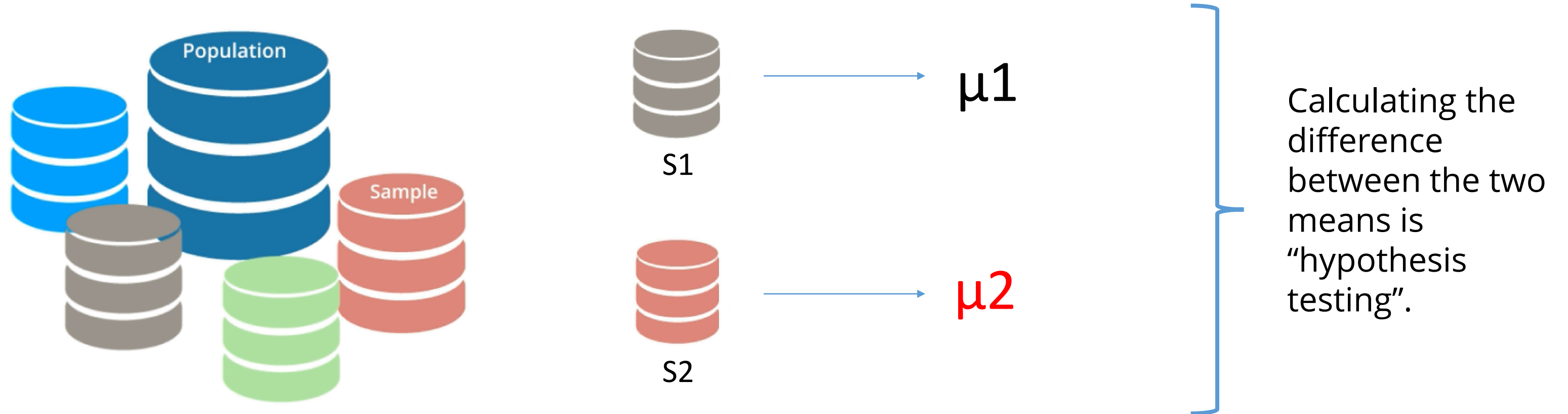
Phase 3: Model Deployment

- Use the model for prediction
- Use the model to compare actual outcome with expectations



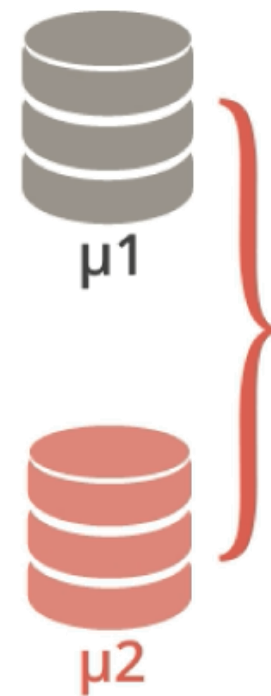
Hypothesis Testing

Draw two samples from the population and calculate the difference between their means.



Hypothesis Testing

Draw two samples from the population and calculate the difference between their means.



$$\mu_1 \neq \mu_2$$

$$\mu_1 = \mu_2$$

Alternative Hypothesis

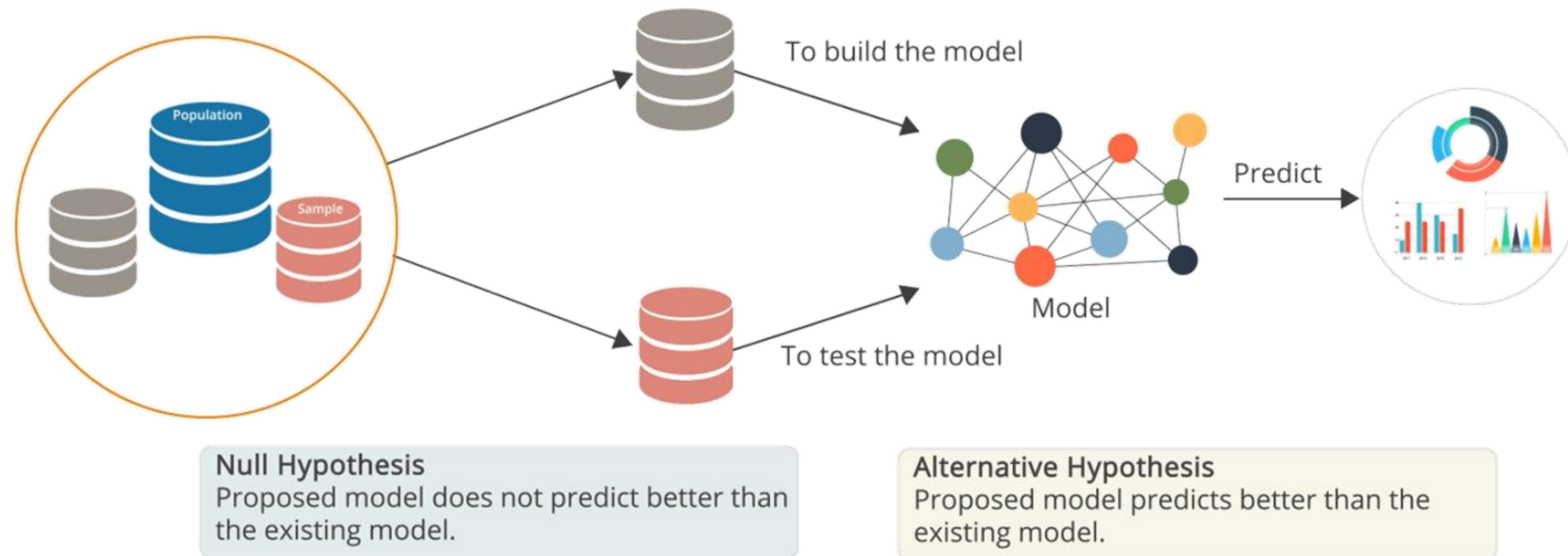
- Proposed model outcome is accurate and matches the data.
- There is a difference between the means of S1 and S2.

Null Hypothesis

- Opposite of the alternative hypothesis.
- There is no difference between the means of S1 and S2.

Hypothesis Testing Process

Choosing the training and test dataset and evaluating them with the null and alternative hypothesis.



Usually the training dataset is between 60% and 80% of the big dataset and the test dataset is between 20% and 40% of the big dataset.

Communication

Data analysis process and results are presented to stakeholders.



Forms of Data analysis presentation:

- Visual graphs
- Plotting maps
- Reports
- Whitepaper reports
- PowerPoint presentations

Data Visualization

Data visualization techniques are used for effective communication of data.



Benefits of data visualization:

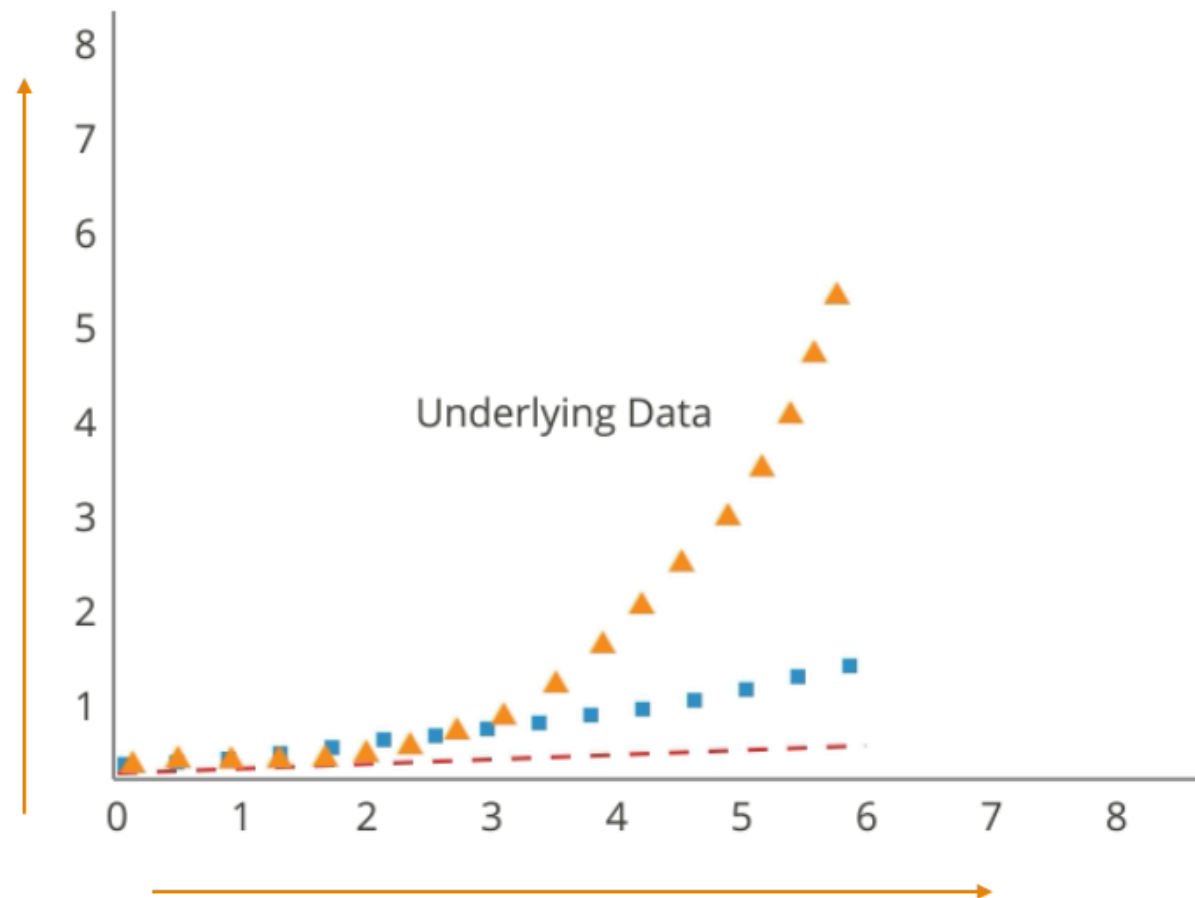
- Simplifies quantitative information through visuals
- Shows the relationship between data points and variables
- Identifies patterns
- Establishes trends

Examples of data visualization:

- Presenting information about new and existing customers on the website and their behavior when they access the website
- Representing web traffic pattern for the website, for example, more activity on the website in the morning than in the evening

Plotting

Plotting is a data visualization technique used to represent underlying data through graphics.

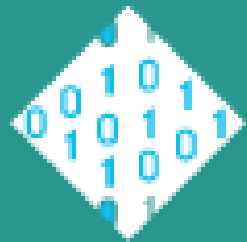


Features of plotting:

- Plotting is like telling a story about data using different colors, shapes, and sizes.
- Plotting shows the relationship between variables.
- Example:
 - Change in value of Y results in change in value of X.
 - X is independent of y.

Data Types for Plotting

There are various data types used for plotting.



Numerical Data

There are two types of numerical data:

Discrete Data – Distinct or counted values

Example: Number of employees in a company or number of students in a class

Continuous Data – Values within a range that can be measured

Example: Height can be measured in feet or inches and weight can be measured in pounds or kilograms



Categorical Data

There are two types of categorical data:

Cluster or group – Grouped values

Example: Students can be divided into different groups based on height – Tall, Medium, and Short

Ordinal data – Grouped values as per ranks

Example: A ranking system; a five-point scale with ranks like “Agree,” “Strongly agree,” and “Disagree”



Time Series

Data measured in time blocks such date, month, year, and time (hours, minutes, and seconds)

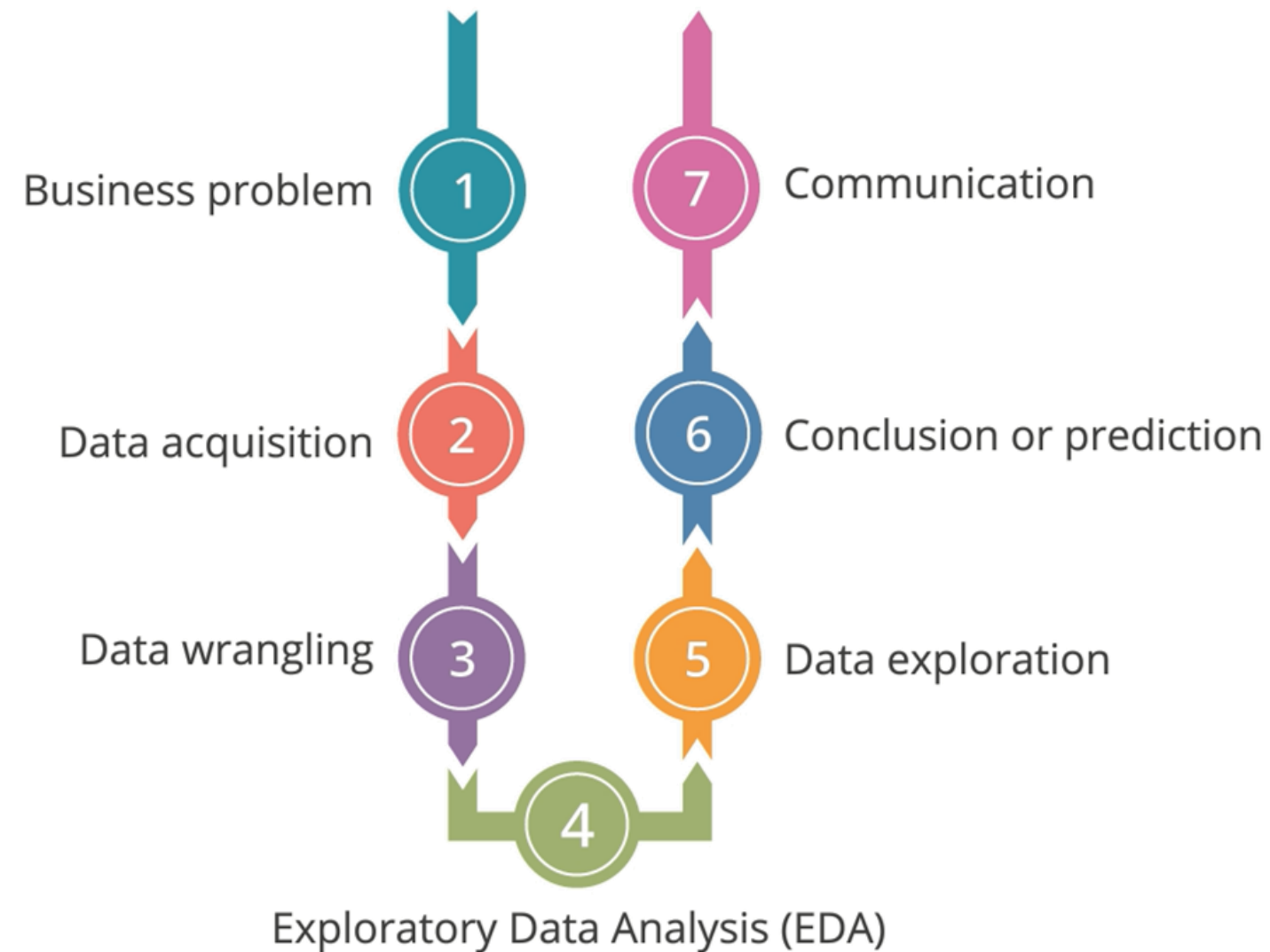
Types of Plot

Different data types can be visualized using various plotting techniques.



Data Analytics – An Iterative Process

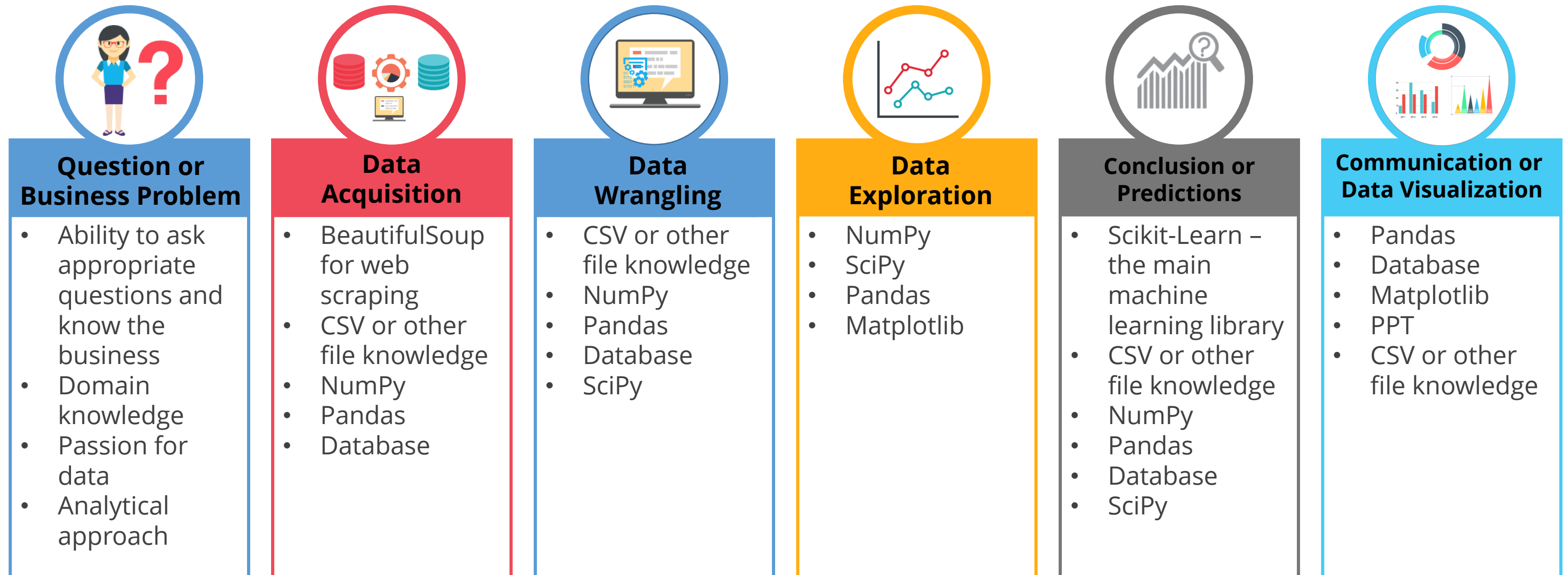
Data analytics is an iterative process involving tracing back the steps, often to ensure that you are on the right track.



Process Result: Question is answered or business problem is solved.

Data Analytics – Skills and Tools

Skills and tools required for each step of the data analysis process.





Knowledge Check

KNOWLEDGE
CHECK

Which plotting technique is used for continuous data?
Select all that apply.

- a. Regression plot
- b. Line chart
- c. Histogram
- d. Heat map



KNOWLEDGE
CHECK

Which plotting technique is used for continuous data?
Select all that apply.

- a. Regression plot
- b. Line chart
- c. Histogram
- d. Heat map



The correct answer is **b, c**.

Explanation: Line charts and histograms are used to plot continuous data.



QUIZ 1

Which Python library is the main machine learning library?

- a. Pandas
- b. Matplotlib
- c. Scikit-learn
- d. NumPy



QUIZ

1

Which Python library is the main machine learning library?

- a. Pandas
- b. Matplotlib
- c. Scikit-learn
- d. NumPy



The correct answer is **c**.

Explanation: SciKit-learn is the main machine library in Python.

QUIZ 2

Which of the following includes data transformation, merging, aggregation, group by operation, and reshaping?

- a. Data Acquisition
- b. Data Visualization
- c. Data Wrangling
- d. Machine learning



QUIZ 2

Which of the following includes data transformation, merging, aggregation, group by operation, and reshaping?

- a. Data Acquisition
- b. Data Visualization
- c. Data Wrangling
- d. Machine learning



The correct answer is **c**.

Explanation: Data wrangling includes data transformation, merging, aggregation, group by operation, and reshaping.

QUIZ 3

Which measure of central tendency is used to catch outliers in the data?

- a. Mean
- b. Median
- c. Mode
- d. Variance



QUIZ 3

Which measure of central tendency is used to catch outliers in the data?

- a. Mean
- b. Median
- c. Mode
- d. Variance



The correct answer is **b**.

Explanation: Median is the exact middle value and most suitable to catch outliers.

QUIZ 4

In hypothesis testing, the proposed model is built on:

- a. the entire dataset.
- b. the test dataset.
- c. a small subset.
- d. the training dataset.



QUIZ 4

In hypothesis testing, the proposed model is built on:

- a. the entire dataset.
- b. the test dataset.
- c. a small subset.
- d. the training dataset.



The correct answer is **d**.

Explanation: The proposed model is built on the training dataset in hypothesis testing.

QUIZ 5

Beautiful soup library is used for ____.

- a. data wrangling
- b. web scraping
- c. plotting
- d. machine learning



QUIZ 5

Beautiful soup library is used for ____.

- a. data wrangling
- b. web scraping
- c. plotting
- d. machine learning.

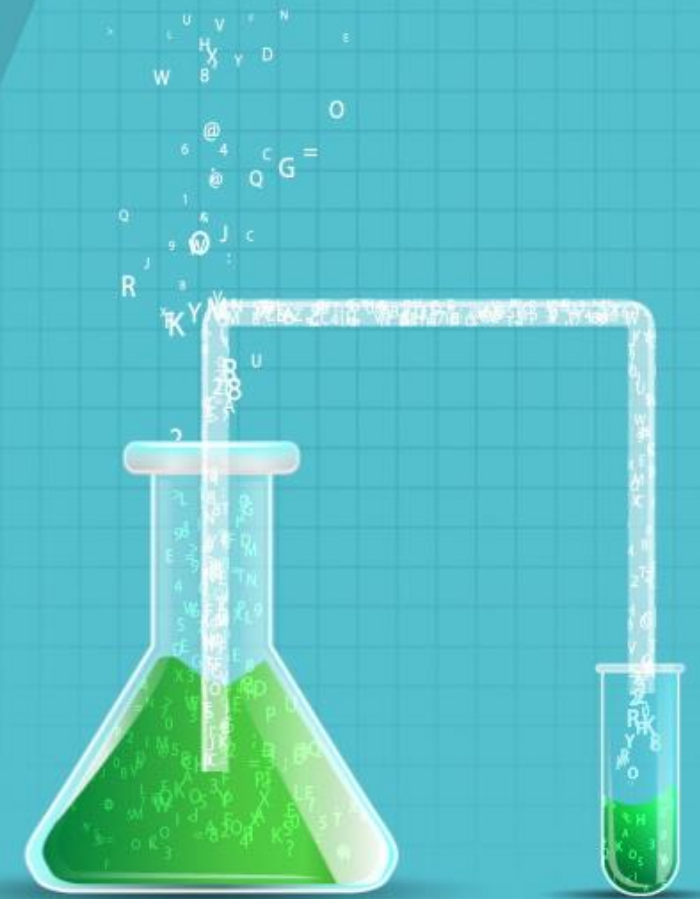


The correct answer is **b**.

Explanation: BeautifulSoup is used for web scraping and mainly used in the data acquisition phase.

Key Takeaways

- Data analytics is used to solve business problems.
- Data analysis requires a number of skills and tools.
- Data wrangling, data exploration, and model selection processes are challenging.
- EDA includes quantitative and graphical techniques.
- Data visualization helps show data characteristics and patterns effectively.
- Hypothesis testing establishes the relationship between dependent and independent variables in data analytics.



This concludes “Data Analytics”

The next lesson is “Statistical Analysis and Business Applications”