# EXPLOLATORY DATA ANALYSIS

# Content :

- Introduction to EDA
- Importance of EDA
- Data types
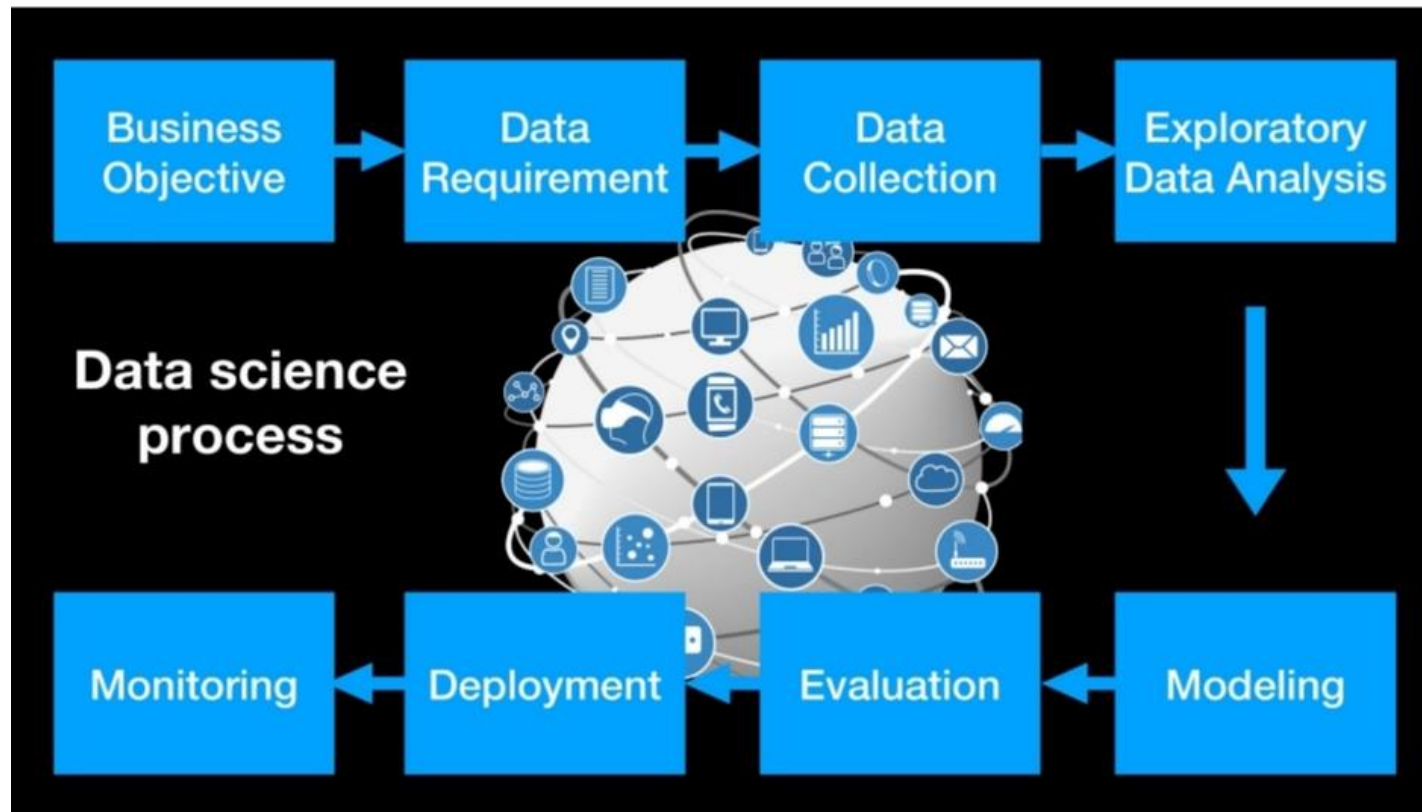- Python Packages for EDA
- Lists of Graphs
- Practical EDA

# 1. INTRODUCTION TO EDA

- Exploratory Data Analysis refers to the critical process of performing initial investigations on data so as to discover patterns, to spot anomalies, to test hypothesis and to check assumptions with the help of summary statistics and graphical representations .

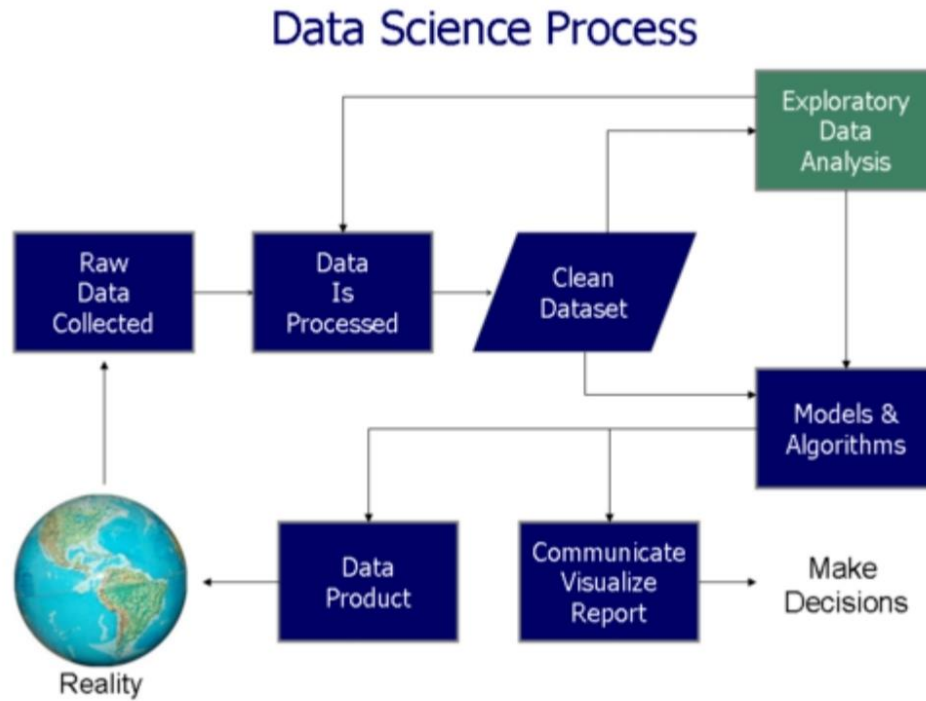- It is a good practice to understand the data first and try to gather as many insights from it.

# 2. IMPORTANCE OF EDA

- Identifying the most important variables/features in your dataset.

- Testing a hypothesis or checking assumptions related to the dataset.

- To check the quality of data for further processing and cleaning.

- Deliver data driven insights to business stakeholders.

- Verify expected relationships actually exists in the data .
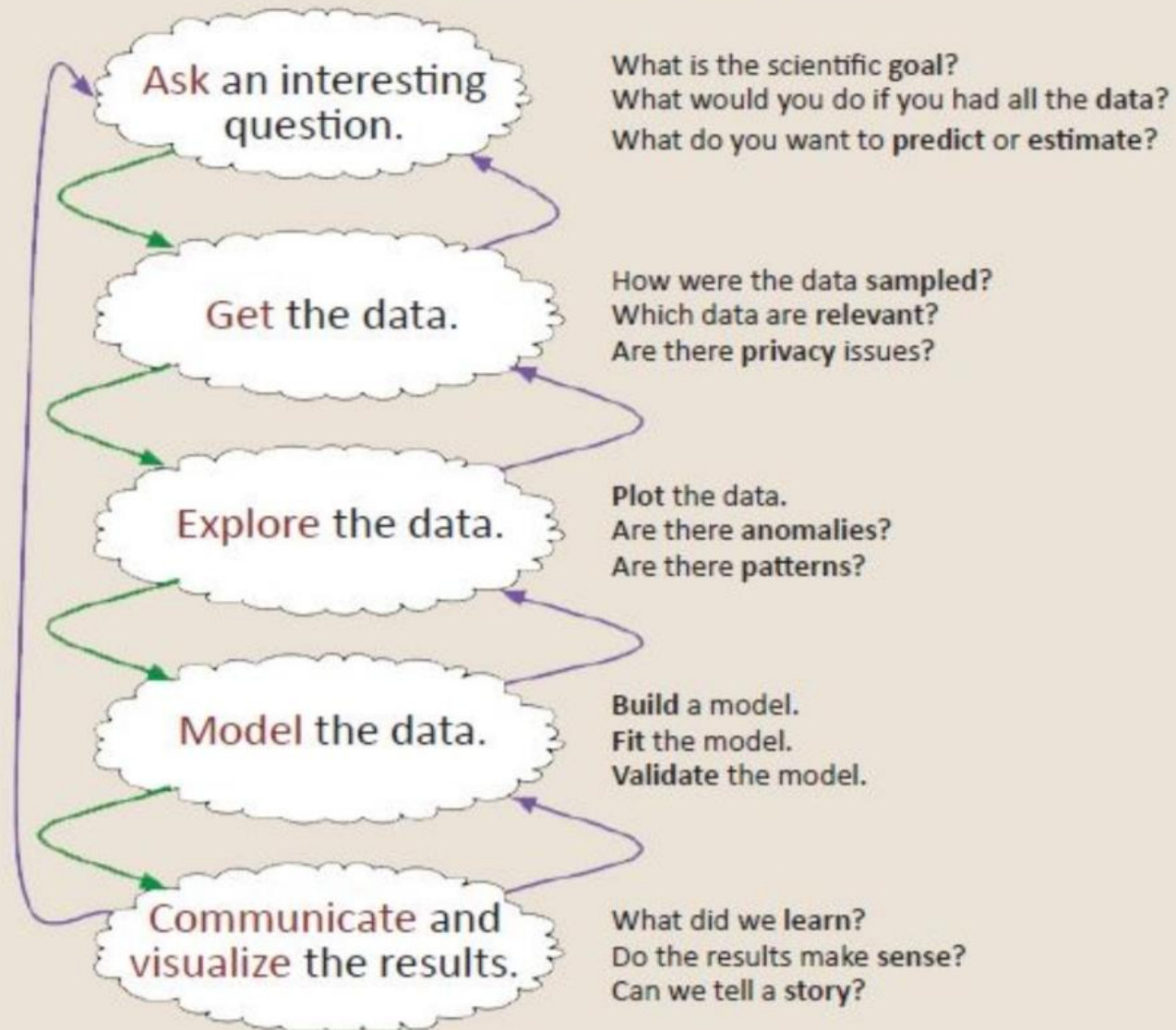
- To find unexpected structure or insights in the data.

# Data Science Modeling Process

# Data Science Process

# The Data Science Process

**Ask an interesting question.**

What is the scientific goal?
What would you do if you had all the **data**?
What do you want to **predict** or **estimate**?

**Get the data.**

How were the data **sampled**?
Which data are **relevant**?
Are there **privacy** issues?

**Explore the data.**

**Plot** the data.
Are there **anomalies**?
Are there **patterns**?

**Model the data.**

**Build** a model.
**Fit** the model.
**Validate** the model.

**Communicate and visualize the results.**

What did we **learn**?
Do the results make **sense**?
Can we tell a **story**?

Joe Blitzstein and Hanspeter Pfister, created for the Harvard data science course http://www.cs109.org/.

# Two categories of Data

- Structured Data types

  example :  CSV file , Excel file , Database file

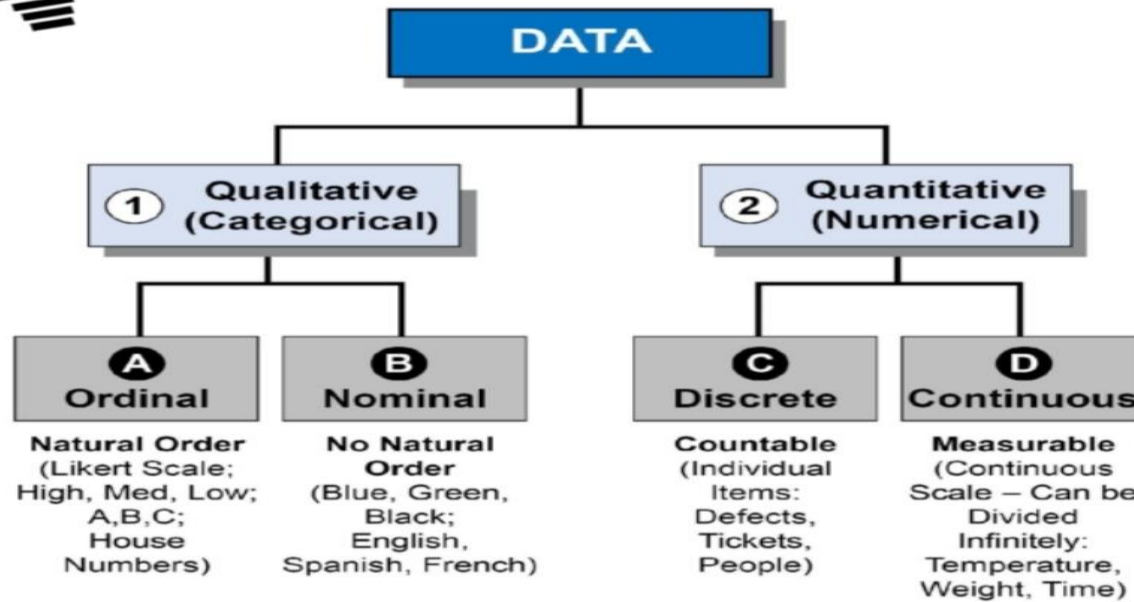- Unstructured Data types

  example :   Images , videos , audio

# Data Types



**DID YOU KNOW?**
#1 – TYPES OF DATA

**DATA**

① **Qualitative (Categorical)**

② **Quantitative (Numerical)**

**A** **Ordinal**
Natural Order (Likert Scale; High, Med, Low; A,B,C; House Numbers)

**B** **Nominal**
No Natural Order (Blue, Green, Black; English, Spanish, French)

**C** **Discrete**
Countable (Individual Items: Defects, Tickets, People)

**D** **Continuous**
Measurable (Continuous Scale – Can be Divided Infinitely: Temperature, Weight, Time)

Notes:
1. Nominal and Ordinal data can be treated as discrete when the frequencies within each group are counted. Bar, Line, Pareto, and Pie charts can be created with discrete data. Histograms, Line graphs, and Scatter diagrams can be created with continuous data.

2. The efficient problem-solver is always aware of the type of data to be analyzed.

# Structured Data Types

**Categorial – This is any data that is not a number .**

▶  Ordinal – have a set of order  eg.  Rating happiness on a scale of 1-10 .

▶ Binary – have only two values   eg.  Male or female

▶ Nominal – no set of order   eg .   Countries

**Numerical -  Data inform of numbers**

▶ Continious  -  numbers that don't have a logical end to them eg. Heights

▶ Discrete   -   have a logical end to them    eg.  Days in the Month

# Python Packages for EDA

# 1. Bar Chart

## 2. Pie Chart

Animals



- Chickens
- Cats
- Dogs
- Goats
- Cows
- Donkeys

# 3.Histogram

# 4.Scatter Plot

# 5. Heatmap



Sales per employee per weekday
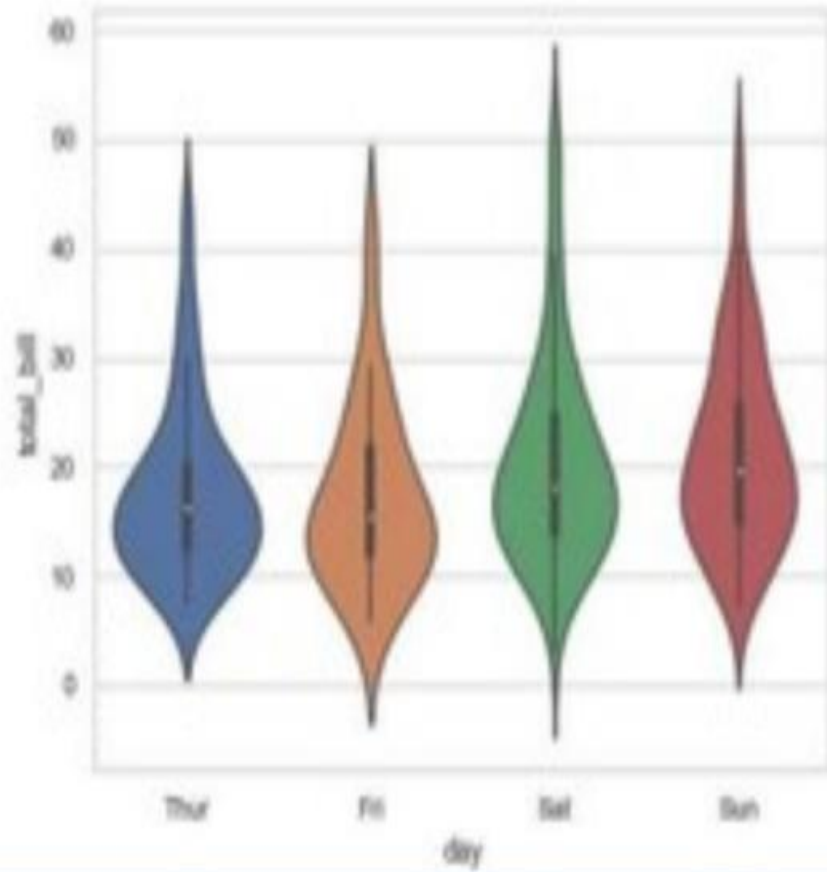
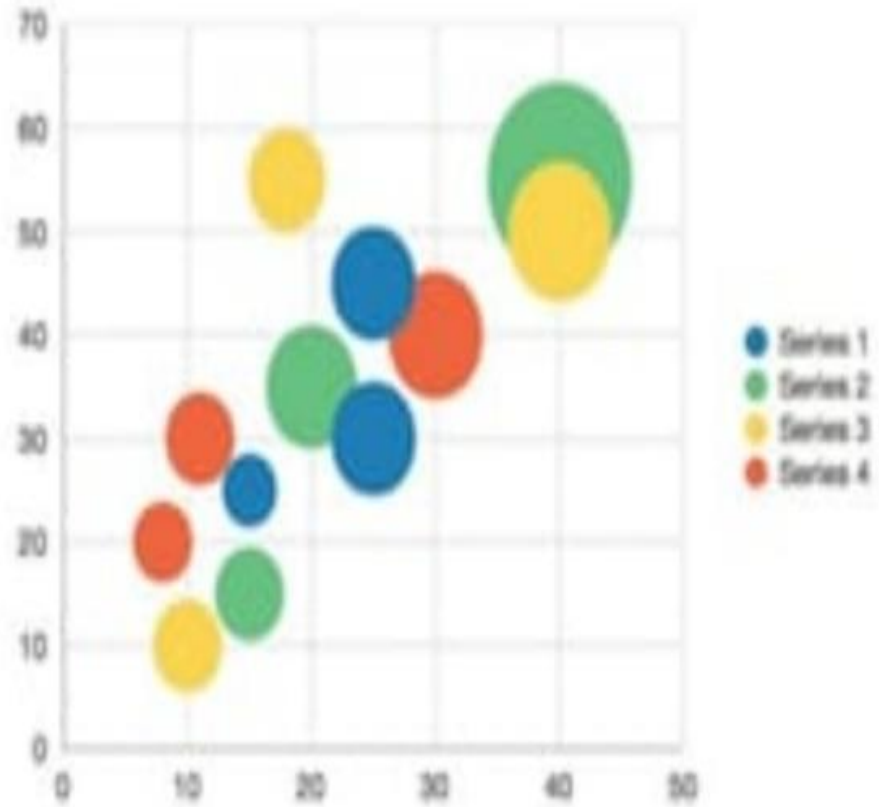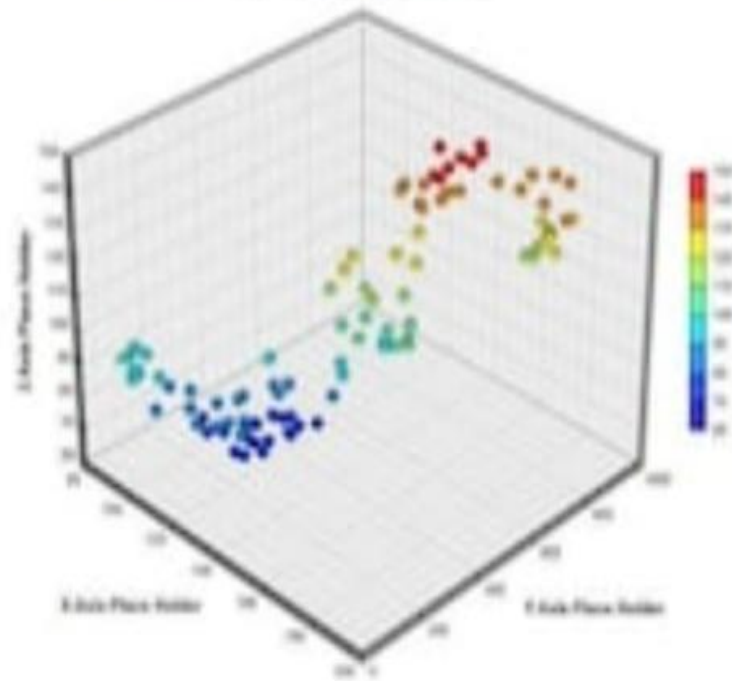# 6. Box Plot

# 7. Line Plot

Product Trends by Month

# 8. Violin Plot

# 9.Bubble Plot

**Bubble Chart**
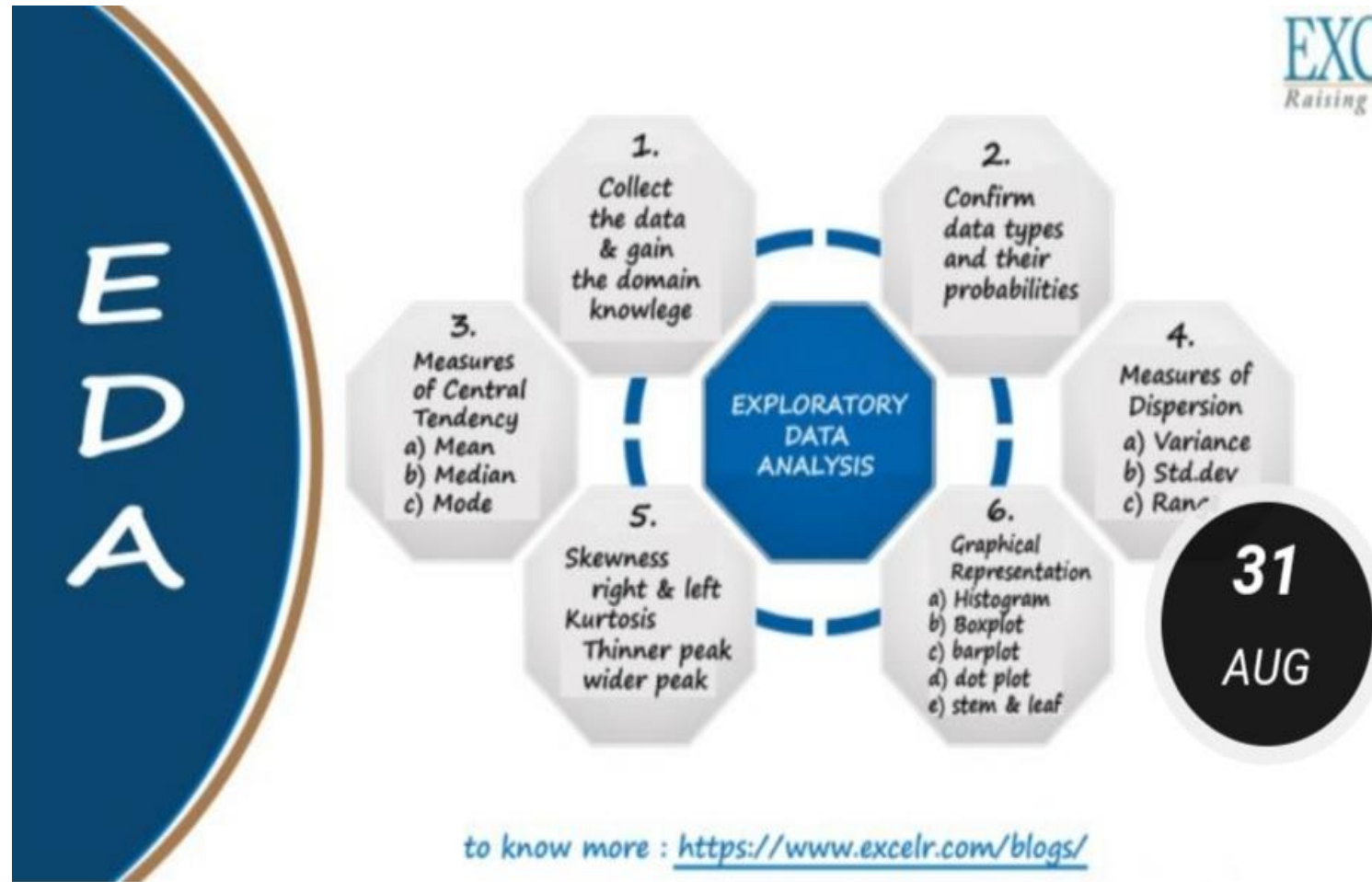
# 10. 3D Scatter Plot


3D Scatter Chart (1)
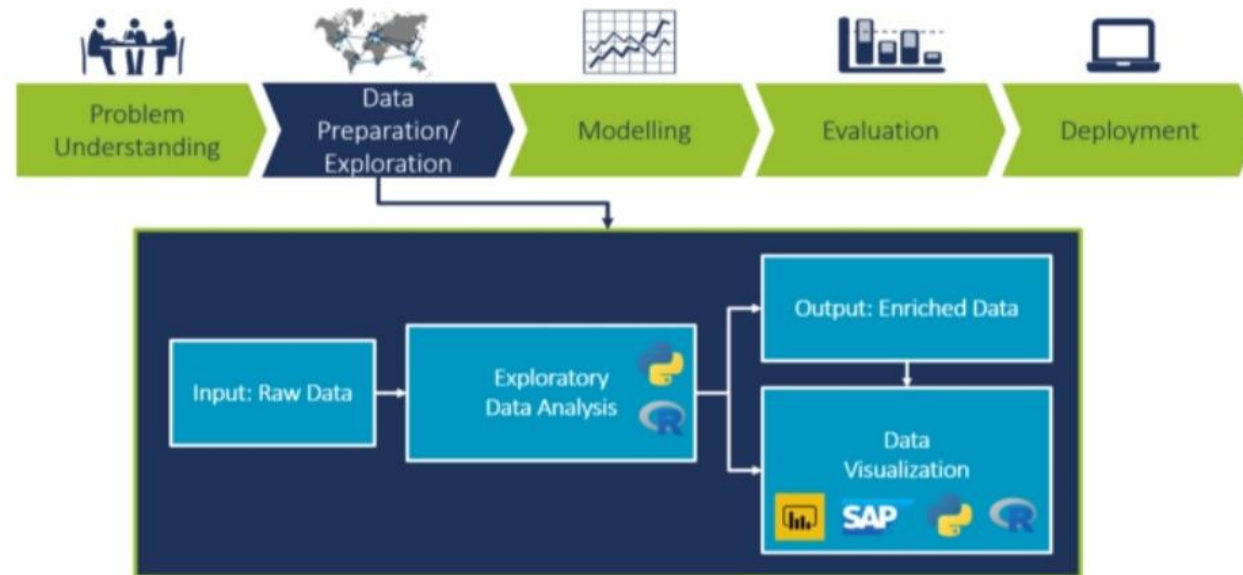
# EDA steps and Visualization

# EDA steps and Assumptions



Fig. 1: Data Science Project Flow

# EDA Analysis Process



Exploratory Data Analysis Process

Goal and Strategy

Modeling or Insights

Data Collection

Feature Engineering
- Extraction
- Augmentation
- Selection

Attribute Classification
- Dimensions
- Metrics
- Data Types

Copyright © 2019 Christopher S. Penn | @cspenn | cspenn.com | @TrustInsights | TrustInsights.ai

Preparation
- Centering
- Scaling
- Cleaning

Initial Analysis
- Univariate
- Multivariate
- Quality
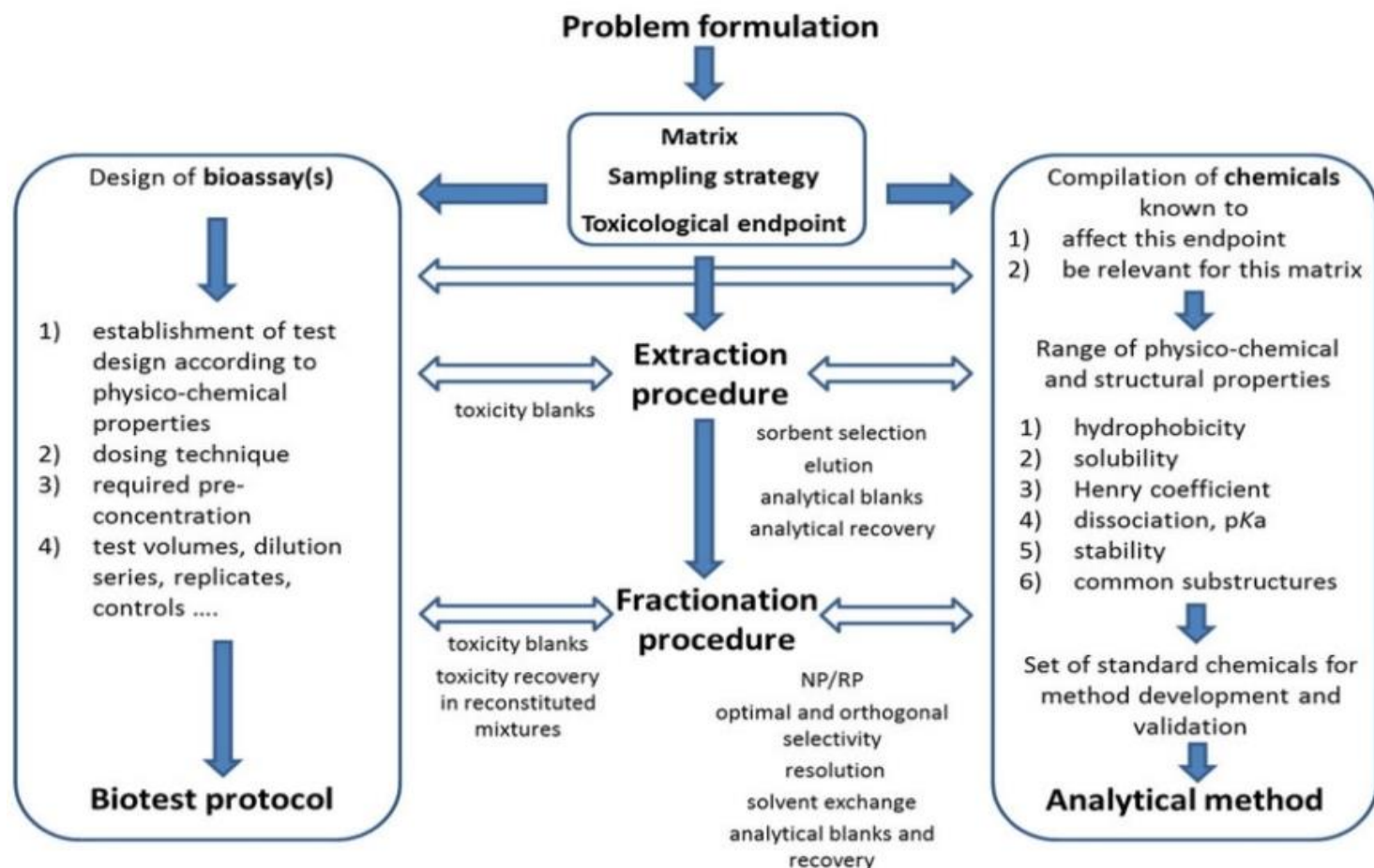- Anomaly Detection

Requirements Verification

# EDA Concept

# EDA Method

# EDA Conclusion

## Conclusion

1. This paper presents an approach for automating the exploratory data analysis step in the knowledge discovery in data.

2. This EDA process identifies inappropriate and suspicious attributes, selects the most appropriate attribute representation, create univariate and multivariate derived attributes, and chooses an optimal subset of attributes to retain for the model.

3. Using the resultant simplified attribute subset reduces elapsed CPU time for building and using a model, increases model accuracy, and improves the explanatory power of the model.

Exploratory Data Analysis

www.educba.com