# Exploratory Data Analysis

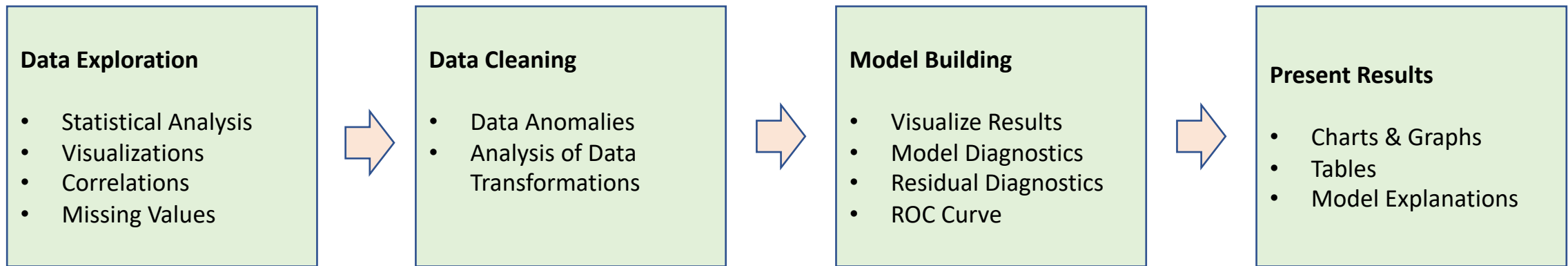## Manas Dasgupta

# Exploratory Data Analysis (EDA)

# Introduction to EDA

Exploratory Data Analysis (EDA) is used by data scientists to analyze and investigate data sets and summarize their main characteristics, often employing data visualization methods. It helps determine how best to manipulate data sources to get the answers you need, making it easier for data scientists to discover patterns, spot anomalies, test a hypothesis, or check assumptions.

EDA is primarily used to see what data can reveal beyond the formal modelling and provides a provides a better understanding of data set variables and the relationships between them. It can also help determine if the statistical techniques you are considering for data analysis are appropriate.

# EDA in various Phases of DS/ML

**Data Exploration**

- Statistical Analysis
- Visualizations
- Correlations
- Missing Values

**Data Cleaning**

- Data Anomalies
- Analysis of Data Transformations

**Model Building**

- Visualize Results
- Model Diagnostics
- Residual Diagnostics
- ROC Curve

**Present Results**

- Charts & Graphs
- Tables
- Model Explanations

# Goals of EDA

**For Direct Business Purposes**

Analysis and Finding Patterns and Insights

Creating Reports and Visualizations

**For Machine Learning**

Pre-processing Data for Machine Learning

Establishing Assumptions about Data

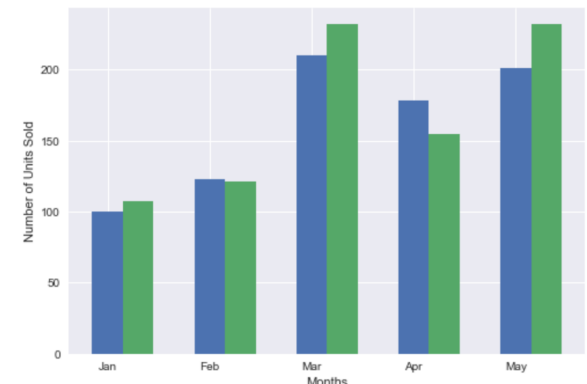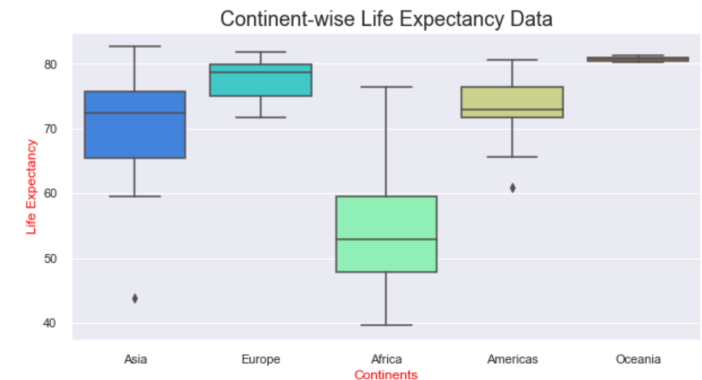# Types of EDA

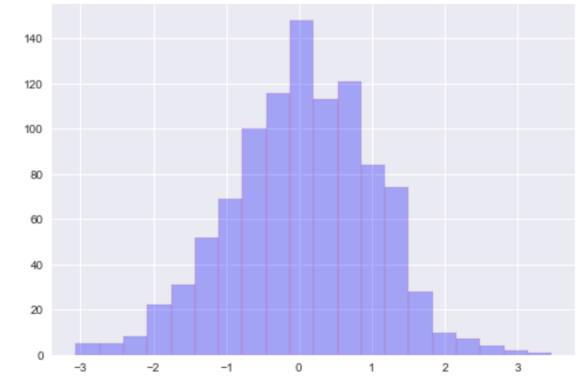**There are four primary types of EDA:**

**Univariate non-graphical:** This is simplest form of data analysis, where the data being analysed consists of just one variable. Since it's a single variable, it doesn't deal with causes or relationships. The main purpose of univariate analysis is to describe the data and find patterns that exist within it.

**Univariate graphical:** Non-graphical methods don't provide a full picture of the data. Graphical methods are therefore required. Common types of univariate graphics include:
- Histograms, a bar plot in which each bar represents the frequency (count) or proportion (count/total count) of cases for a range of values.
- Box plots, which graphically depict the five-number summary of minimum, first quartile, median, third quartile, and maximum.

**Multivariate nongraphical:** Multivariate data arises from more than one variable. Multivariate non-graphical EDA techniques generally show the relationship between two or more variables of the data through cross-tabulation or statistics.
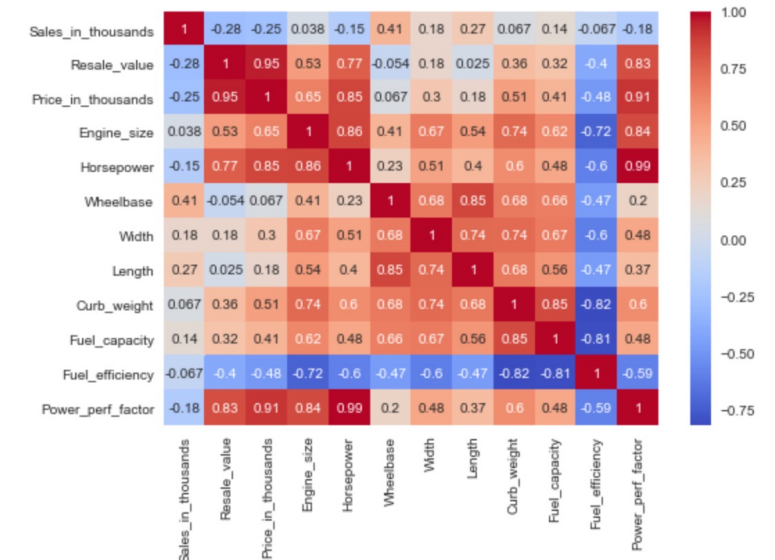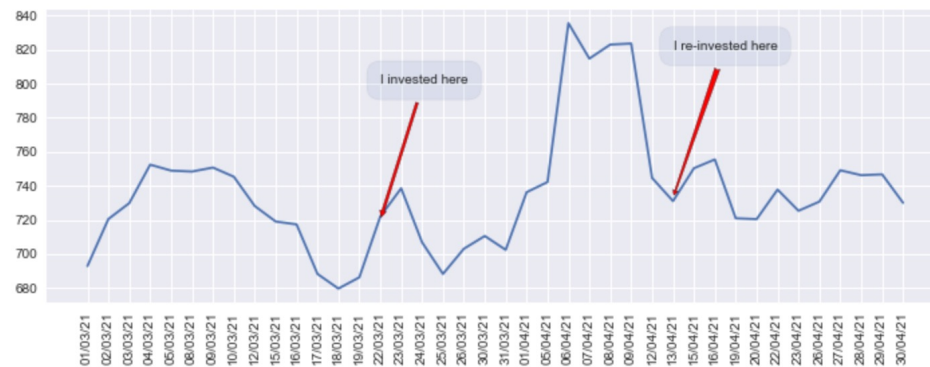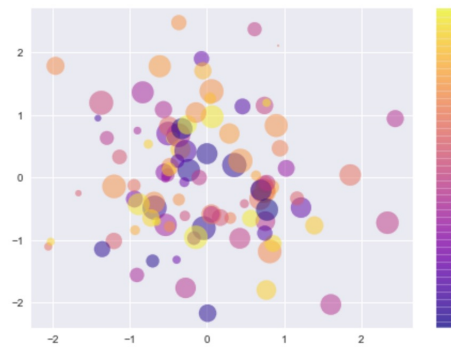
**Multivariate graphical:** Multivariate data uses graphics to display relationships between two or more sets of data. The most used graphic is a grouped bar plot or bar chart with each group representing one level of one of the variables and each bar within a group representing the levels of the other variable.





Continent-wise Life Expectancy Data

# Types of EDA

Other common types of multivariate graphics include:

- **Scatter plot**, which is used to plot data points on a horizontal and a vertical axis to show how much one variable is affected by another.

- **Multivariate chart**, which is a graphical representation of the relationships between factors and a response.

- **Run chart**, which is a line graph of data plotted over time.

- **Bubble chart**, which is a data visualization that displays multiple circles (bubbles) in a two-dimensional plot.

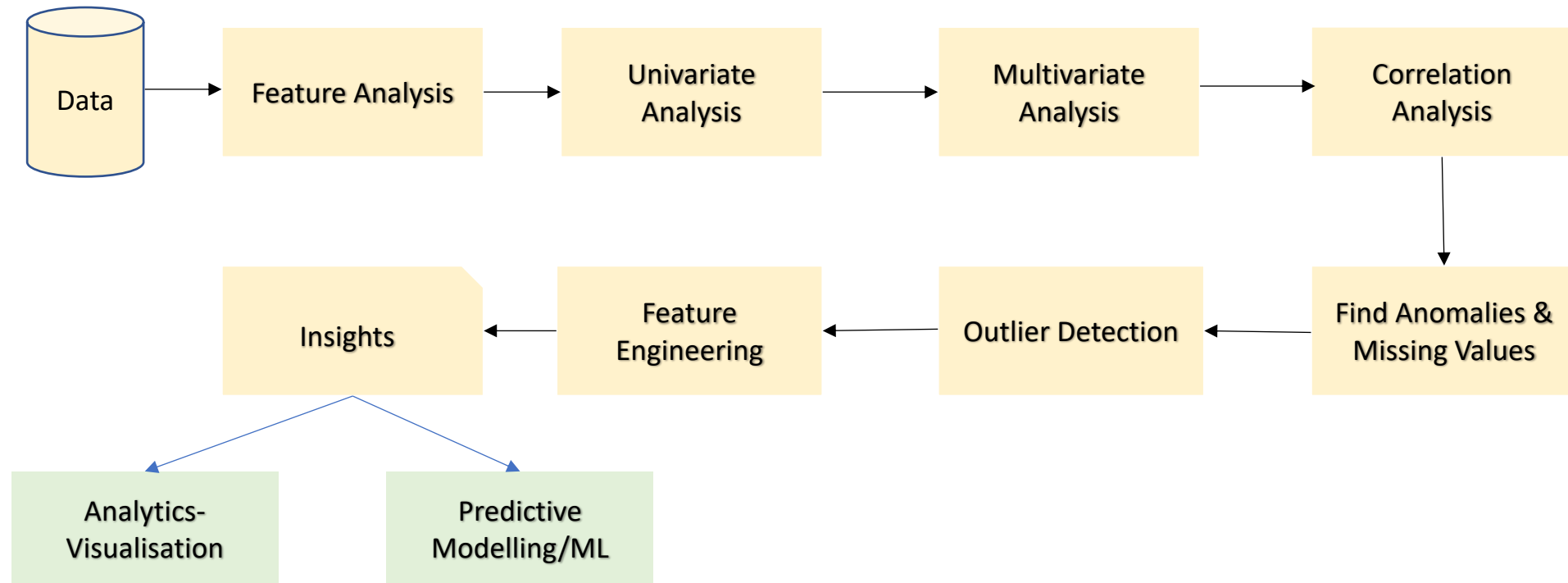- **Heat map**, which is a graphical representation of data where values are depicted by color.

# Tools of EDA

**Various Statistical Techniques you are used to perform EDA:**

1. Univariate Analysis and visualizations.

2. Bivariate and Multi-variate visualizations and visualizations for mapping and understanding interactions between different features in the data.

3. Predictive models, such as Linear Regression, use statistics and data to predict outcomes.

4. Classification Modelling for correctly predicting data labels.

5. Unsupervised techniques such as Clustering (e.g. K-Means) that help in finding similar behavioural patterns in Data.

6. Dimensionality Reduction Techniques (e.g. PCA) that help bringing down complex high dimensional data into lower dimensional data without losing (much) information.

# High Level EDA Process

# Why EDA in Machine Learning

**EDA is a very critical task in Machine Learning for the following purposes**

- Find out Missing Data

- Finding Data Anomalies

- Finding Outliers

- Finding class imbalances

- Finding correlations between data

- Finding the right features for inclusion in the model

- Ascertain data types/values and transformation needs

- Ascertain Data Scales for scale transformations

- Ascertain need for derived features which might be more useful
  and help reduce dimensionality

- Ascertain need for additional data and external data integrations

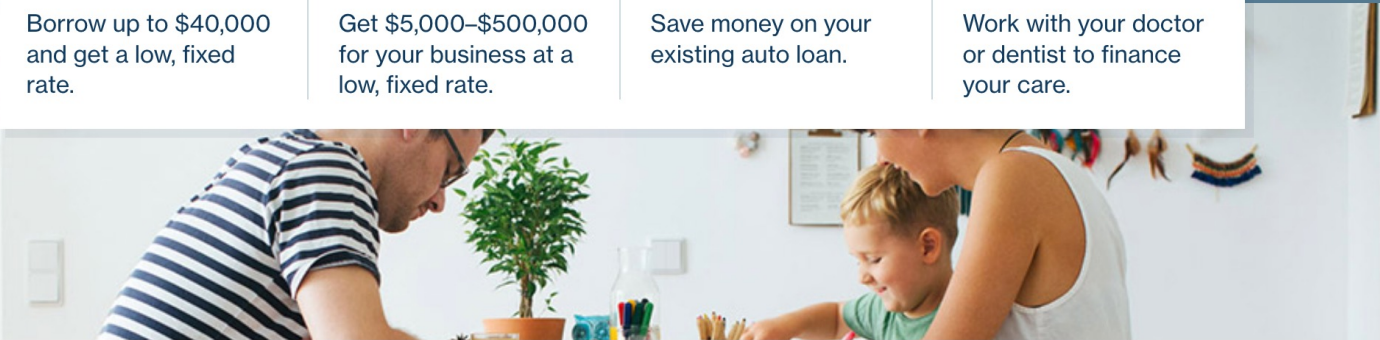- Ascertain Noises in Data

# Lending Club EDA Case Study



Lending Club is an US based Consumer Finance Company. They have shared their Loan Transaction data in public for the purpose of Data Science research. The most popular usage of Lending Club Data has become the EDA Case Study.

# Lending Club Project Problem Statement

When Lending Club receives a loan application from an individual, the company has to make a decision for loan approval based on the Applicant's profile and Credit History. This is called Credit Decision. A correct Credit Decision is important to the company as a rejected Credit is loss of business, at the same time, credit given to an un-credit-worthy person leads to risk of default.

To make the Credit Decision process efficient and correct, Lending Club wants to undertake an exercise to analyse past loan data with the incidences of 'default's and 'non-default's. Through this analysis, Lending Club wants to find what factors in the Applicants Profile must they focus on to ensure Credit is given to the right applicants and reduce chance of Credit Decision errors.

Historically, a manual effort with some programming/tool intervention would have helped in the process with large amount of effort and unreliability in results.

Luckily for us, with Data Science techniques at our disposal, we will look at the automating the entire process and bring out the required insights.

# EDA Process

Statistical Analysis and Summarisation of the Data

Univariate Analysis

Bivariate Analysis

Conclusions

As part of the Analysis, we will use a number EDA techniques including abundant amount of Visualisations.

The final goal will be to draw conclusions from the data that would help Lending Club in better profiling of customer for better Credit Decisions.