# EXPLORATORY DATA ANALYSIS

Imesha Dilshani
*Got help from AI*

# Contents

# Introduction to EDA

Exploratory Data Analysis (EDA) is a must-do step in the data analysis process where one examines and understands what the most important features, characteristics, patterns, or relationships within any given dataset are before going on to do further modeling or other forms of more systematic analyses.

EDA can often be an iterative process of revisiting and refining analyses as new insights emerge. It's a pillar of data analysis, laying the groundwork for intelligent decision-making and problem solving.

EDA is a set of techniques and approaches used to:

- What are the main characteristics of a dataset?
- Discover patterns, trends, and anomalies.
- Identify relationships between variables.
- Visualize data distributions and relationships.
- Generate hypotheses for further testing.
- Check assumptions for statistical analysis.
- Preparing the data for modeling.
- guide feature engineering and model selection.

## Key benefits of EDA:

- Improves understanding of the data
- Uncovers hidden patterns and insights.
- Identifies possible problems with the data (e.g., errors, outliers)
- Informs the selection of features and model building.
- Points the direction for future analysis.
- Guarantees the validity of statistical assumptions.
- Communicates insights effectively through visualizations.
- What are the benefits of doing EDA?
- Describe the usual stages in EDA.

# Data Loading and Understanding

The first step in any Exploratory Data Analysis (EDA) process is data loading and understanding. It means bringing the data into your chosen environment and getting a rough feel for its characteristics. Here's a breakdown of what it entails:

## Importing Libraries

First, import the libraries needed for data manipulation and analysis (pandas), visualization (Matplotlib & Seaborn) as well as NumPy for calculation.

## Loading the Dataset

Use functions provided in the libraries to read data from its source (e.g., CSV file, database, or API). This could mean specifying file paths, database credentials or API keys.

## Initial Inspection

You can use the head and tail functions to get a glimpse of just where there are. Then shape allows you to figure out how much there was-how many rows and columns exist in your data table (number of rows and columns).

## Data Information

Use info to show data types, get non-null counts and memory usage. It shows the basic structure and possible problems such as missing values.

## Descriptive Statistics

For numerical variables, use functions like describe to obtain summary statistics (mean value, median and standard deviation etc.) This explains central tendencies, spread and possible outliers.

## Data Types and Format

Examine the data types of each column (e.g., integer, float or string) and fix any problems that might occur during later processing such as inconsistencies among these columns or formatting issues.

## Missing Values

Identify the extent of missing values in each variable and analyze them. This requires a suitable approach to dealing with them. For instance, substitute imputation; delete or flag outright.

## Exploratory Plotting

Present the distributions of numerical variables with histograms, box plots and density charts. Create bar charts or pie charts for categorical variables to understand their frequencies and proportions.

## Documentation

Document the process of loading and understanding, recorded in such detail as data sources used, selected libraries employed; unexpected difficulties encountered during the course of processing requests; first observations. This will be a point of reference for future analysis and cooperation.

Through careful and thorough loading and understanding of your data, you build a good base from which to launch an enjoyable and revelatory EDA exploration.

# Univariate Analysis

It's important to understand the separate data points while working on Exploratory Data Analysis. Univariate analysis is necessary to make it possible. It is a kind of process like looking for jobs, step by step, examine and analyze each career, their reputation, smartness, etc. Here's what it entails:

## Objectives

- Describe and summarize: The most important descriptive statistics to measure central tendency and spread of data include the mean, mode, median, standard deviation, and range.
- Visualize distribution: Like histograms, box plots, density plots, and violin plots, we need to present the distribution of variable values in a way that represents the shape and frequency of the data.
- Identify outliers: We should identify extreme deviants from the mean and do something to deal with it because it can impact our later analysis.
- Explore relationships: Understand the relationship between variables using scatter plots and correlation coefficients.
- Inform assumptions: Obtain information about data patterns and distributions to test the assumptions required for later statistical tests or modeling.

## Analyzing Numerical Variables

Well, we learn that descriptive statistics like mean, median, standard deviation, quartiles, and range can be used to comprehend the central tendency and variation of values.

We can make the distribution of values visible through visualizations like histograms, box plots and density plots.

We can find the relationship between two variables by utilizing the correlation coefficient which calculates the relationship between columns.

Look for any patterns in the given information by analyzing factors such as time or order of data presentation.

## Analyzing Categorical Variables

To get an idea of how prevalent each category is within the data, calculate frequencies and proportions for them.

Depict category distribution visually in, say, bar charts and pie charts.

Examine contingency tables to find relations between categorical variables.

Examine associations and patterns among various categories.

## Choosing Appropriate Techniques

Choose the best approach according to data set and target skill set. Like, for continuous data, you can use histograms but for discrete data, it's better to use bar charts.

## Iteration and Refinement

repetitive task. Use the knowledge you have gathered so far to get a better understanding of the problem, develop a more specific plan for investigating the issue, and pay more attention to certain aspects of the data you are working with.

## Remember

Univariate analysis is the key to understanding your data properly and making conclusions before going ahead with multivariate analysis.

If you master univariate analysis, you will be able to understand individual variables in your dataset which is the base to understand the overall analysis of the dataset.

*got help from AI*

# Bivariate Analysis

Bivariate analysis is used to understand the relationship between two variable values in the dataset. It's quite important to identify the relationships between variables during EDA to discover any patterns, connections, or impacts. Here's what it entails.

## Objectives

- Discover relationships: Determine connections or correlations among pairs of variables.
- Quantify relationships: Determine the intensity and orientation of connections through the application of correlation coefficients or alternative statistical indexes.
- Visualize relationships: Make graphs and charts that help to understand the data quickly.
- Inform feature selection: Help me choose which features to use in my analysis.
- Generate hypotheses: Determine the possible reasons that cause relationships which can be examined later.

## Techniques and Visualizations

- Scatter plots: Picture how numerical values are connected to each other, so you can see things like a straight line or a group of items that are close together.
- Correlation coefficients: Determine the strength and direction of a linear relationship between two variables. We usually calculate the correlation coefficient Pearson's correlation coefficient(r) for numerical data and Spearman's rank correlation for ordinal data.
- Contingency tables: It is necessary to calculate measures like chi-square independence tests in order to summarize the relationship between two categorical variables. You can do this by using cross-tabulations to calculate the frequencies of each category.
- Grouped bar charts: When comparing numerical variables across categories in a categorical variable, it is important to keep in mind that these groups may have varying characteristics. This makes it difficult to draw meaningful conclusions from the data alone.
- Stacked bar charts: Represent graphically a full picture of the data when comparing two groups or categories.
- Heatmaps: This means we can show the connection between few variable with different intensity of colors.

## Considerations for Different Data Types

Numerical-Numerical: Scatter plots and correlation coefficients can be used to represent the relationship between two variables.

Categorical-Categorical: Do contingency tables, bar charts, and mosaic plots. They are helpful to analyze data.

Numerical-Categorical: We may Also use these visual tools.

## Interpretation and Insights:

In a scatter plot, we can examine the relationship between two variables by plotting them on an x and y axis. By looking at the data points on the graph, we can determine if there is a pattern or trend between the variables. If there are any data points that seem to stand out from the others, these are potential outliers.

The correlation coefficient helps to identify the strength and direction of the relationship between two variables, indicating whether it's a positive, negative, or no correlation.

To determine relationships and absence of relationship between different categorical variables, we have to analyze contingency tables.

Be mindful of possible confounding factors that could affect relationships.

## Key Points:

Univariate analysis investigates individual variables while bivariate analysis goes beyond that by examining how variables correlate with each other.

It is important to use the proper tools and methods to accurately understand and represent the given information.

It is important to be careful when analyzing data and not assume that a relationship exists where none exists.

Put your findings on paper with graphs and data to help others better understand and follow your work.

When you do bivariate analysis correctly, you will see the relationship between the two factors and understand the data better. This will make it easier to create advanced models and do more advanced analysis.

# Multivariate Analysis

Multivariate analysis helps to evaluate the correlation and impact of several variables in a data set. It helps in identifying a lot of hidden relationships and connections among the different variables. Here's a breakdown:

## Objectives

- Discover complex relationships: Identify relationships and connections between numerous factors that help to discover complex structures and interactions within them.
- Identify underlying dimensions: To decrease the complexity of complicated datasets, we need to identify the significant patterns that explain a major portion of the variation.
- Group similar observations: Segmenting groups of data based on common traits.
- Visualize high-dimensional data: Make graphs and diagrams that can help us understand the connections between things.
- Inform feature selection and model building: Help me choose the right variable and build a good model for my task.

## Common Techniques

Principal Component Analysis (PCA): It reduces the number of variables to make it easy analyzing.

Factor Analysis: It works similarly to PCA but its main goal is to find the hidden relationship among the given variables.

Cluster Analysis: Sorting and categorizing things based on the similarity between them is known as clustering. There are several techniques for achieving this like K-means clustering and hierarchical clustering and density-based clustering.

Multidimensional Scaling (MDS): It helps to show large amounts of data on a smaller graph or image, so that the relationships between the different pieces of information are still obvious.

Multivariate Regression: This is a multivariable regression, which looks at how many different factors can affect one thing. It's like linear regression but with more than one predictor variable.

Correlation Matrix: Shows the connections between every pair of variables, which helps us understand the relationships among them.

Heatmaps: The graph represents the relationship among different variables with numerous shades of color and hierarchical sorting technique used to show the clusters.

Pair Plots: A matrix of scatter plots is created to exhibit how the various stock prices are interrelated; this aids us in quickly understanding market trends.

## Key Considerations

Tailor techniques to analysis goals: Select techniques that match with your personal objectives and required outcomes.

Visualize effectively: It is important to present the data in a visual manner so that it becomes easier to understand and interpret the insights.

Interpret cautiously: It's important to realize that there can be biases, restrictions and uncertainties when we apply multivariate methods.

Validate findings: To ensure that the results are reliable and accurate, one can apply more methods or make use of previously acquired domain knowledge.

## Integration with Other EDA Stages

uses univariate and bivariate analysis to find out more complicated pattern.

Helps to inform decisions on feature engineering and model selection for further analysis.

Gives a comprehensive understanding of the data and its hidden properties.

Multivariate analysis is a powerful tool for uncovering hidden structures and patterns in complex datasets. By mastering it, you'll gain the ability to make more informed decisions based on data-driven insights.

# Data Cleaning and Preprocessing

Transforming raw data into a format suitable for further analysis and modelling is important, and is done in two crucial steps, namely data cleaning and preprocessing. This way the data that you obtain is accurate and reliable. Here's a breakdown of their components:

## Data Cleaning

Identifying and handling missing values: Determine on whether to replace missing values with mean or median worth and also decide whether we should remove the missing rows or put it for later discovery.

Correcting inconsistencies and errors: Ensure that the data is clean, correct, and in the right format for better data analysis.

Outlier detection and treatment: So outliers are basically those values that are way different from the rest of the data. If we find such values in our data then we need to handle them carefully. We can either remove them, winsorize them, or transform them. Removing outliers means discarding those values that are considered extreme. Winsorizing involves replacing those extreme values with more moderate ones. Finally, we can transform outliers to make them more like the rest of the data.

Dealing with duplicates: How to remove duplicate rows in excel and handle them?

Encoding categorical variables: Make categorical variable into numerical representation so that model can understand it.

## Data Preprocessing

Feature scaling: Standardize or normalize numerical features to ensure they have similar scales and prevent biases in models.

Feature engineering: Create new features based on existing ones to potentially improve model performance or capture relevant information.

Dimensionality reduction: Reduce the number of features if necessary, using techniques like PCA or feature selection, to improve model efficiency and avoid overfitting.

Data partitioning: Split your data into training, validation, and test sets for model training, evaluation, and unbiased performance assessment.

## Benefits of Data Cleaning and Preprocessing

Improves data quality: Ensures consistent, accurate, and complete data for reliable analysis.

Enhances model performance: Leads to more accurate and robust models by addressing biases and optimizing data for modeling algorithms.

Increases transparency and reproducibility: Documents data cleaning steps for clear communication and repeatability of analysis.

Saves time and resources: Prevents potential issues later in the analysis pipeline and avoids wasting time on unreliable data.

## Key Considerations

Choose data cleaning and preprocessing techniques based on your specific data and analysis goals.

Document your decisions and steps to ensure transparency and reproducibility.

Be cautious about introducing biases through data cleaning or preprocessing.

Validate your cleaning and preprocessing methods to ensure they haven't negatively impacted data quality.

By effectively cleaning and preprocessing your data, you pave the way for successful and insightful analysis, laying the foundation for reliable results and actionable insights.

# Data Visualization

Data visualization is a powerful tool within Exploratory Data Analysis (EDA) that brings insights to life through visual representations. It empowers you to communicate patterns, trends, and relationships within your data in a clear, engaging, and often intuitive way. Here's a breakdown:

## Objectives

- Reveal hidden patterns: Uncover trends, anomalies, and relationships that might not be apparent in raw data tables.
- Communicate insights effectively: Share findings with diverse audiences, fostering understanding and driving action.
- Explore data interactively: Enable dynamic exploration of data through interactive visualizations, fostering deeper insights and understanding.
- Support decision-making: Provide visual evidence to aid in informed decisions and problem-solving.
- Generate hypotheses: Guide the formation of testable hypotheses for further investigation.

## Types of Visualizations

- Univariate visualizations: Histograms, box plots, density plots, bar charts, pie charts, scatter plots.
- Bivariate visualizations: Scatter plots, correlation matrices, grouped bar charts, stacked bar charts, heatmaps.
- Multivariate visualizations: Pair plots, parallel coordinates, scatter plot matrices, heatmaps with clustering.

## Key Considerations

Choose visualizations wisely: Select visualizations that align with your data types, analysis goals, and audience.

Design for clarity: Ensure visualizations are easy to read, interpret, and visually appealing.

Use visual cues effectively: Employ color, size, shape, and position to highlight patterns and guide attention.

Label clearly: Provide descriptive titles, axis labels, and legends for context and understanding.

Prioritize interactivity: Consider interactive visualizations for dynamic exploration and engagement.

Integrate with storytelling: Use visualizations to weave compelling narratives that drive action.

## Tools and Libraries

Python: Matplotlib, Seaborn, Plotly, Bokeh, Altair.

R: ggplot2, plotly, shiny.

*got help from AI*

JavaScript: D3.js, Chart.js, Vega-Lite.

Dashboarding tools: Tableau, Power BI, Qlik Sense.

## Best Practices

Experiment with different visualizations: Explore various options to find those that best reveal insights.

Iterate and refine: Adjust visualizations based on feedback and evolving understanding.

Contextualize visualizations: Explain key takeaways and implications for decision-making.

Consider ethical implications: Ensure visualizations represent data fairly and accurately, avoiding biases or misleading representations.

By harnessing the power of data visualization, you can communicate complex data stories in a captivating and impactful way, leading to better understanding, informed decisions, and actionable insights.

# Ethical Implications of EDA

While Exploratory Data Analysis (EDA) is a powerful tool for uncovering insights, it's crucial to consider its ethical implications. Here are some key points to keep in mind:

## Bias and Fairness

Data may inherently contain biases. EDA can unknowingly amplify these biases if not approached cautiously.

Be mindful of how variables are chosen, grouped, and visualized. Ensure your analysis doesn't reinforce or perpetuate existing discriminatory or unfair patterns.

Consider conducting fairness analyses and counterfactual reasoning to identify and mitigate potential bias in your findings.

## Privacy and Security

Anonymize or aggregate data, when possible, especially when dealing with sensitive information.

Implement appropriate security measures to protect data from unauthorized access or misuse.

Be transparent about data sources and anonymization techniques used in your analysis.

## Transparency and Reproducibility

Document your EDA process clearly, including data sources, cleaning methods, and chosen visualizations.

This allows others to understand your findings and validate your analysis.

Avoid cherry-picking results or manipulating data to fit a predetermined narrative.

## Explainability and Interpretation

Ensure your visualizations and interpretations are clear and understandable, avoiding jargon or misleading representations.

Contextualize your findings within the broader domain and avoid drawing unwarranted conclusions from limited data.

Be upfront about assumptions, limitations, and uncertainties associated with your analysis.

## Algorithmic Bias

Algorithms used in EDA may themselves be biased, potentially leading to biased insights.

Choose well-established and validated algorithms with documented fairness considerations.

Monitor for potential algorithmic bias in your analysis and take steps to mitigate its impact.

## Societal Impact

Consider the potential societal implications of your findings and how they might be used.

Avoid analyses that could perpetuate harmful stereotypes or contribute to societal inequalities.

Use your research ethically and responsibly, promoting positive societal change.

By staying mindful of these ethical considerations, you can conduct EDA in a responsible and transparent manner, ensuring your findings contribute to a fairer, more equitable, and inclusive data-driven future.