

Contents

Introduction to EDA	2
Key benefits of EDA:	2
Data Loading and Understanding.....	2
Importing Libraries.....	2
Loading the Dataset	3
Initial Inspection.....	4
Data Information.....	4
Descriptive Statistics	5
Data Types and Format	5
Missing Values.....	5
Exploratory Plotting.....	6
Documentation	7
Univariate Analysis.....	7
Objectives.....	7
Analyzing Numerical Variables.....	8
Analyzing Categorical Variables	8
Choosing Appropriate Techniques	8
Iteration and Refinement.....	8
Remember.....	8

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.

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.ticker as mtick
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [ ]:
```

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.

```
In [2]: telco_base_data = pd.read_csv('WA_Fn-UseC_-Telco-Customer-Churn.csv')
```

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).

```
In [3]: telco_base_data.head()
```

```
Out[3]:
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	DeviceProtection	TechSupp
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No	
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...	Yes	
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...	No	
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...	Yes	
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No	

5 rows x 21 columns

```
In [4]: telco_base_data.shape
```

```
Out[4]: (7043, 21)
```

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.

```
In [11]: telco_base_data.info(verbose = True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   customerID            7043 non-null   object
1   gender                 7043 non-null   object
2   SeniorCitizen          7043 non-null   int64
3   Partner                7043 non-null   object
4   Dependents             7043 non-null   object
5   tenure                 7043 non-null   int64
6   PhoneService           7043 non-null   object
7   MultipleLines           7043 non-null   object
8   InternetService         7043 non-null   object
9   OnlineSecurity          7043 non-null   object
10  OnlineBackup            7043 non-null   object
11  DeviceProtection        7043 non-null   object
12  TechSupport             7043 non-null   object
13  StreamingTV             7043 non-null   object
14  StreamingMovies         7043 non-null   object
15  Contract                7043 non-null   object
16  PaperlessBilling        7043 non-null   object
17  PaymentMethod           7043 non-null   object
18  MonthlyCharges          7043 non-null   float64
19  TotalCharges            7043 non-null   object
20  Churn                   7043 non-null   object
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB
```

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.

```
In [7]: telco_base_data.describe()
Out[7]:
```

	SeniorCitizen	tenure	MonthlyCharges
count	7043.000000	7043.000000	7043.000000
mean	0.162147	32.371149	64.761692
std	0.368612	24.559481	30.090047
min	0.000000	0.000000	18.250000
25%	0.000000	9.000000	35.500000
50%	0.000000	29.000000	70.350000
75%	0.000000	55.000000	89.850000
max	1.000000	72.000000	118.750000

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.

```
In [6]: telco_base_data.dtypes
Out[6]:
```

customerID	object
gender	object
SeniorCitizen	int64
Partner	object
Dependents	object
tenure	int64
PhoneService	object
MultipleLines	object
InternetService	object
OnlineSecurity	object
OnlineBackup	object
DeviceProtection	object
TechSupport	object
StreamingTV	object
StreamingMovies	object
Contract	object
PaperlessBilling	object
PaymentMethod	object
MonthlyCharges	float64
TotalCharges	object
Churn	object
dtype:	object

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.

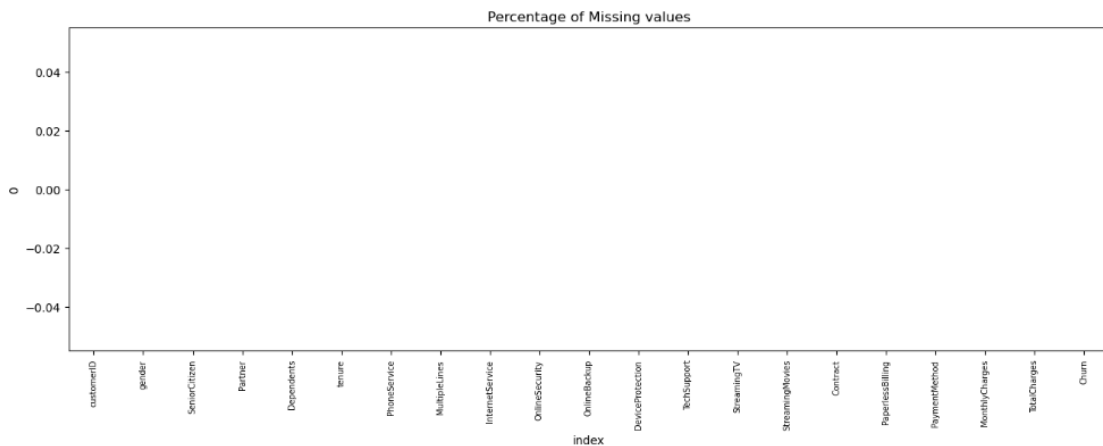
```
In [17]: # Create a DataFrame with the percentage of missing values
missing = pd.DataFrame((telco_base_data.isnull().sum()) * 100 / telco_base_data.shape[0]).reset_index()

# Set up the plot
plt.figure(figsize=(16, 5))

# Use a bar plot to visualize the percentage of missing values
ax = sns.barplot(x='index', y=0, data=missing)

# Customize the plot
plt.xticks(rotation=90, fontsize=7)
plt.title("Percentage of Missing values")

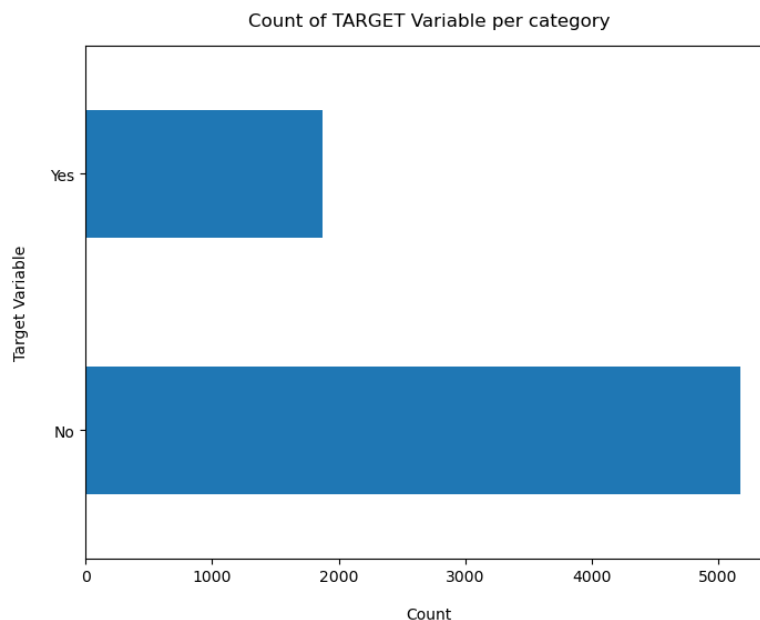
# Show the plot
plt.show()
```



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.

```
In [8]: telco_base_data['Churn'].value_counts().plot(kind='barh', figsize=(8, 6))
plt.xlabel("Count", labelpad=14)
plt.ylabel("Target Variable", labelpad=14)
plt.title("Count of TARGET Variable per category", y=1.02);
```



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

One key aspect of Exploratory Data Analysis (EDA) is gaining an understanding of individual variables within your dataset. The first step to making this possible is univariate analysis. Think of it as the act of diving deep into each variable one at a time, looking for its patterns, tendencies and characteristics. Here's what it entails:

Objectives

- Describe and summarize: Understanding the central tendency and spread of data Point or measure the most important descriptive statistics are mean, mode median, standard deviation range etc.
- Visualize distribution: Conform your visualizations to those like histograms, box plots, density plots and violin plots to display the form and frequency of values for a variable.

- Identify outliers: Extract and process possible outliers that are extreme deviants from the main distribution, which might interfere with later analysis.
- Explore relationships: Investigate how the variable is related to other variables in the dataset through 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

You can understand the central tendency and spread of values with descriptive statistics such as mean, median, standard deviation quartiles and range.

The distribution of values can be made visible through visualizations such as histograms, box plots and density plots.

Explore potential relationships by calculating correlation coefficients with other variables.

Interpret trends and patterns in the data, examining aspects such as time or order.

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

Pick the right techniques based on data type and what you want to learn. For instance, histograms are more suitable for continuous variables and bar charts work well on categorical data.

Iteration and Refinement

Univariate analyses can be an iterative process. Rely on your preliminary findings to put some structure into further probing, refine the questions you are asking yourself, and take a closer look at specific parts of the data.

Remember

This univariate analysis offers a solid basis from which to grapple with your data, before jumping headlong into more complex relationships and models.

Should you master univariate analysis, then not only will you come to understand each variable in your dataset, but this lays the foundation upon which one can explore and thoroughly comprehend the big picture.

4. Bivariate Analysis:

Bivariate analysis delves into the relationships between pairs of variables within your dataset. It's a crucial step in Exploratory Data Analysis (EDA) that unveils how variables interact, revealing potential patterns, associations, and influences. Here's what it entails:

1. Objectives:

Discover relationships: Identify potential associations, patterns, or dependencies between pairs of variables.

Quantify relationships: Measure the strength and direction of relationships using correlation coefficients or other statistical measures.

Visualize relationships: Create informative visualizations to depict relationships clearly and intuitively.

Inform feature selection: Guide the choice of relevant features for subsequent modeling and analysis.

Generate hypotheses: Develop testable hypotheses about the nature and causes of relationships for further investigation.

2. Techniques and Visualizations:

Scatter plots: Visualize relationships between numerical variables, revealing patterns like linear trends, clusters, or nonlinear associations.

Correlation coefficients: Measure linear relationships, commonly using Pearson's correlation coefficient (r) for numerical variables and Spearman's rank correlation for ordinal variables.

Contingency tables: Summarize relationships between two categorical variables, cross-tabulating frequencies and calculating measures like chi-square independence tests.

Grouped bar charts: Compare numerical variables across categories of a categorical variable.

Stacked bar charts: Visualize the composition of a whole across categories of two categorical variables.

Heatmaps: Represent relationships between multiple variables simultaneously, often using color intensity to indicate strength of association.

3. Considerations for Different Data Types:

Numerical-Numerical: Use scatter plots and correlation coefficients.

Categorical-Categorical: Use contingency tables, bar charts, and potentially mosaic plots.

Numerical-Categorical: Use grouped bar charts, box plots, and potentially violin plots.

4. Interpretation and Insights:

Analyze scatter plots to identify patterns, trends, and potential outliers.

Interpret correlation coefficients to understand strength and direction of linear relationships (positive, negative, or no correlation).

Analyze contingency tables to identify associations and independence between categorical variables.

Consider potential confounding variables that might influence relationships.

5. Key Points:

Bivariate analysis builds upon univariate analysis, exploring relationships between variables rather than individual characteristics.

Choose appropriate techniques and visualizations based on data types and analysis goals.

Beware of spurious correlations and consider potential causal relationships cautiously.

Document your findings, including visualizations and statistical measures, for clarity and reproducibility.

By effectively conducting bivariate analysis, you'll uncover meaningful connections and interactions between variables, deepening your understanding of the data and setting the stage for more sophisticated modeling and analysis.

5. Multivariate Analysis:

Multivariate analysis takes EDA to a higher level, exploring relationships among multiple variables simultaneously. It unveils complex patterns, associations, and structures that might not be apparent when examining variables in isolation. Here's a breakdown:

1. Objectives:

Discover complex relationships: Uncover patterns and associations involving multiple variables, often revealing hidden structures and interactions.

Identify underlying dimensions: Reduce the dimensionality of complex datasets by finding patterns that explain a large portion of the variation.

Group similar observations: Cluster data points based on shared characteristics, revealing natural groupings or segments within the dataset.

Visualize high-dimensional data: Create informative visualizations to represent complex relationships and patterns in a visual format.

Inform feature selection and model building: Guide variable selection and model development for predictive or descriptive tasks.

2. Common Techniques:

Principal Component Analysis (PCA): Reduces dimensionality by identifying new, uncorrelated variables (principal components) that capture most of the variance in the data.

Factor Analysis: Similar to PCA, but focuses on identifying latent factors that explain correlations between variables.

Cluster Analysis: Groups data points based on similarity, using techniques like K-means clustering, hierarchical clustering, and density-based clustering.

Multidimensional Scaling (MDS): Visualizes high-dimensional data in lower-dimensional space, preserving relative distances between points.

Multivariate Regression: Analyzes relationships between multiple predictor variables and a response variable, extending linear regression to handle multiple predictors.

Correlation Matrix: Displays correlations between all pairs of variables in a table format, revealing patterns of association.

Heatmaps: Visualizes relationships between multiple variables using color intensity, often with hierarchical clustering to reveal groups.

Pair Plots: Creates a matrix of scatter plots to visualize pairwise relationships between multiple variables simultaneously.

3. Key Considerations:

Tailor techniques to analysis goals: Choose methods that align with your specific questions and desired insights.

Visualize effectively: Use appropriate visualizations to communicate complex relationships clearly and intuitively.

Interpret cautiously: Consider potential biases, limitations, and assumptions of multivariate techniques.

Validate findings: Use additional methods or domain knowledge to confirm insights and avoid overfitting.

4. Integration with Other EDA Stages:

Builds upon univariate and bivariate analyses to uncover more comprehensive patterns.

Informs feature engineering and model selection for further analysis.

Contributes to a holistic understanding of the dataset and its underlying structure.

By mastering multivariate analysis, you'll unlock the power to uncover hidden structures, patterns, and relationships within complex datasets, leading to deeper insights and more informed decisions.

6. Data Cleaning and Preprocessing:

Data cleaning and preprocessing are crucial steps in the EDA process, transforming your raw data into a format suitable for further analysis and modeling. They ensure the quality and integrity of your data, leading to more accurate and reliable results. Here's a breakdown of their components:

1. Data Cleaning:

Identifying and handling missing values: Decide whether to impute missing values (with strategies like mean or median imputation), remove rows with missing values, or flag them for further consideration.

Correcting inconsistencies and errors: Address typos, data entry mistakes, and formatting issues that might hinder analysis.

Outlier detection and treatment: Analyze and decide how to handle outliers that deviate significantly from the rest of the data (e.g., remove, winsorize, or transform).

Dealing with duplicates: Identify and remove duplicate rows or handle them based on your analysis needs.

Encoding categorical variables: Convert categorical variables into numerical representations for compatibility with modeling algorithms.

2. 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.

3. 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.

4. 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.

7. 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:

1. 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.

2. 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.

3. 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.

4. Tools and Libraries:

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

R: ggplot2, plotly, shiny.

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

Dashboarding tools: Tableau, Power BI, Qlik Sense.

5. 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.

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:

1. 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.

2. 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.

3. 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.

4. 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.

5. 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.

6. 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.