### **Exploratory Data Analysis**

**Data Profiling** 

Examine the structure, quality, and characteristics of the dataset. Identify data types, missing values, outliers, and distributional properties.

**Visualizations** 

2

3

Create a variety of visualizations like histograms, scatter plots, and heatmaps to uncover patterns, relationships, and anomalies in the data.

**Hypothesis Generation** 

Formulate informed hypotheses about the data based on your initial explorations. These will guide further analysis and provide direction for the project.

**Feature Engineering** 

Identify and create new features that may provide more predictive power. This could involve transforming, combining, or engineering existing variables.



## Data Analysis and Interpretation

Welcome to the world of data analysis! In this course, we'll explore the exciting realm of understanding, interpreting, and gleaning insights from data. You'll learn how to effectively analyze and manipulate data using various techniques, allowing you to make meaningful conclusions and informed decisions.



#### **Understanding Data Set Variables**

Before diving into the analysis, it's crucial to understand the data set you're working with. Identify the variables, their types, and the relationships between them. This foundation helps you choose the appropriate analytical tools and methods for effective exploration.

#### **Categorical Variables**

Represent characteristics or qualities, often expressed as words or labels. Examples include gender, color, or region.

#### **Numerical Variables**

Represent measurable quantities, often expressed as numbers. Examples include height, weight, or age.

#### Relationships

Examine how variables interact with each other. This could involve correlation, causation, or dependence.

#### **Data Cleaning and Manipulation**

Real-world data often contains inconsistencies, missing values, or outliers. To ensure accurate analysis, you must clean and manipulate the data. This involves removing errors, handling missing information, and transforming variables into a suitable format.



#### **Data Validation**

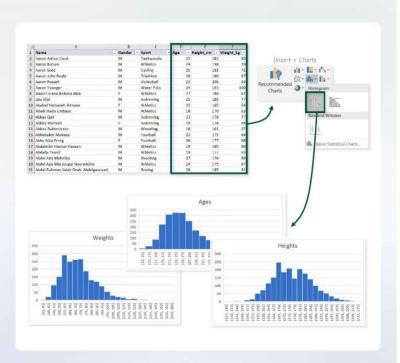
Examine the data for errors or inconsistencies, using techniques like range checks or logical consistency.

#### **Missing Value Imputation**

Handle missing values using methods like mean substitution, last observation carried forward, or more complex imputation techniques.

#### **Data Transformation**

Convert variables into a more suitable format for analysis. This may involve standardizing, normalizing, or creating new variables.



#### **Univariate Data Analysis**

Univariate analysis focuses on understanding individual variables. You'll explore the distribution, central tendency, and dispersion of each variable to gain insights into its characteristics. Techniques include:

1 Frequency Distributions

Visualize the frequency of each value or category within a variable.

- Measures of Central Tendency

  Calculate the mean, median, or mode to understand the central value of a variable.
- 3 Measures of Dispersion

Analyze the spread or variability of a variable using range, variance, or standard deviation.

4 Box Plots

Provide a visual representation of the distribution, quartiles, and outliers of a variable.

#### **Multivariate Data Analysis**

Multivariate analysis explores the relationships between multiple variables. This involves understanding how variables interact, influence, or predict each other. Common techniques include:

#### **Correlation Analysis**

Measure the strength and direction of the linear relationship between two variables. Use a scatter plot to visualize the correlation.

#### **Regression Analysis**

Estimate the relationship between a dependent variable and one or more independent variables. Use regression models to predict the dependent variable based on the independent variables.

#### Principal Component Analysis (PCA)

Reduce the dimensionality of data by finding principal components that capture the most variance. This simplifies the analysis and allows for easier visualization. Use a biplot to visualize the relationship between variables and observations in the reduced dimensional space.



#### **Feature Selection**

From the vast collection of variables, select the most relevant features that contribute significantly to the analysis or prediction task. This involves identifying variables that have strong relationships with the target variable or provide unique insights into the data.

Method	Description
Univariate Feature Selection	Rank variables based on their individual scores (e.g., chi-square test, ANOVA, mutual information).
Recursive Feature Elimination (RFE)	Iteratively remove features that have the least impact on the model's
Feature Importance	performance. Use machine learning algorithms to determine the importance of each feature in a model.

#### **Exploratory Data Analysis (EDA)**

EDA is a crucial step in data analysis, focusing on uncovering patterns, relationships, and anomalies in the data. This involves visualizing the data, testing assumptions, and identifying potential insights.

1

#### **Data Visualization**

Use various charts, plots, and graphs to visualize the data and identify patterns, trends, and outliers.

2

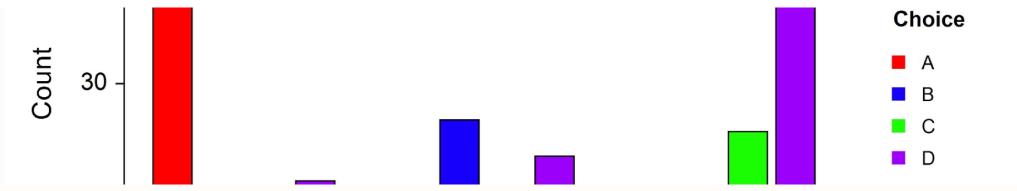
#### **Descriptive Statistics**

Calculate measures of central tendency, dispersion, and distribution to summarize and understand the data.

#### **Hypothesis Testing**

Test assumptions about the data using statistical tests to confirm or reject hypotheses.





#### **Drawing Conclusions and Insights**

Based on your analysis, draw meaningful conclusions and insights from the data. This involves interpreting the results of your analysis, highlighting significant findings, and identifying areas for further exploration.



#### **Identifying Patterns**

Recognize recurring patterns or trends in the data that indicate underlying relationships.



#### **Testing Assumptions**

Validate or reject initial assumptions about the data based on your analysis.



#### **Generating Insights**

Derive meaningful interpretations from the data, answering relevant questions and uncovering hidden insights.



#### Forming Recommendations

Use your insights to formulate actionable recommendations for decision-making or further research.

#### **Assessment and Evaluation**

The assessment of your data analysis journey is crucial to understand its effectiveness and impact. This involves evaluating the quality of your analysis, considering the limitations, and identifying areas for improvement.

1 Accuracy and Reliability

Assess the accuracy and reliability of your findings, ensuring they are grounded in sound methodology and analysis.

Relevance and Applicability

Evaluate the relevance and applicability of your findings to the original research question or problem.

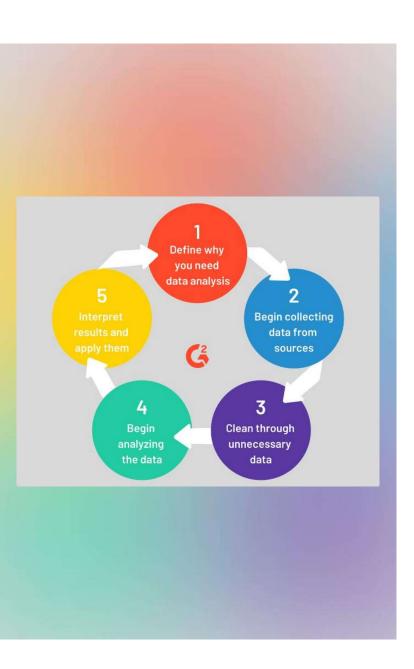
**3** Limitations and Future Future Directions

Acknowledge the limitations of your analysis and identify areas for future research or improvement.



## Introduction to Exploratory Data Analysis

Exploratory Data Analysis (EDA) is a crucial step in the data science workflow. It involves examining and summarizing data to gain insights, uncover patterns, and identify potential issues before building predictive models. This process helps to understand the data's characteristics, identify relevant variables, and uncover hidden relationships. EDA allows data scientists to make informed decisions about data cleaning, feature engineering, and model selection.



#### Importance of EDA in the Data Science Workflow

1 Data Understanding

EDA helps understand the data's underlying structure, distribution, and relationships, providing a foundational understanding for further analysis.

3 Feature Engineering

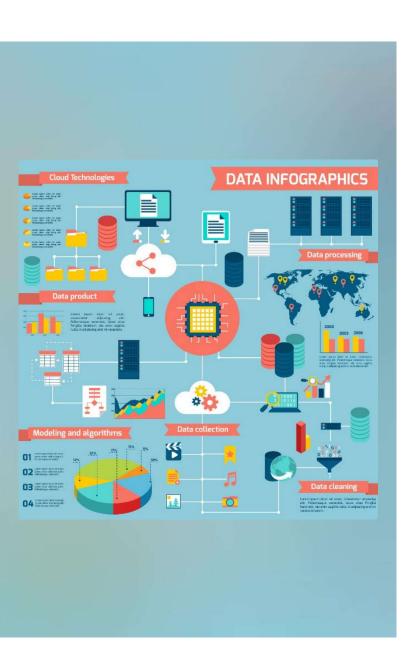
EDA helps discover potential features, relationships, and interactions that can improve the performance of machine learning models.

2 Data Cleaning and Preparation

EDA identifies outliers, missing values, and inconsistencies that need to be addressed before building models, ensuring data quality and reliability.

4 Model Selection

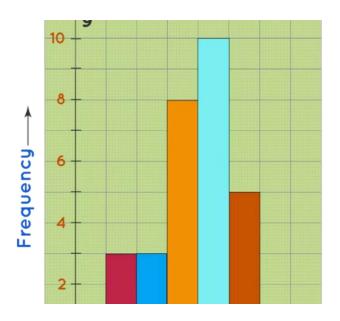
EDA insights guide the choice of appropriate models based on data characteristics, ensuring the model's suitability for the problem.



## Identifying Data Types and Structures

Data Type	Description	Examples
Numeric	Quantitative data, representing measurable quantities.	Age, height, temperature, income.
Categorical	Qualitative data, representing distinct groups or	Gender, color, city, product type.
Ordinal	categories. Categorical data with an inherent order or ranking.	Education level, satisfaction rating, customer feedback.

#### Visualizing Univariate Distributions



# Upper inner fence (24.5) Upper inner fence (24.5) Lower hinge (17.0) Lower adjacent value (13.0)

#### Histograms

Show the distribution of a single variable using bars to represent frequency counts.

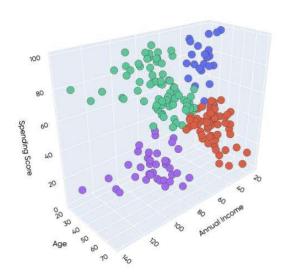
#### **Box Plots**

Display the quartiles, median, and potential outliers of a dataset, providing a compact representation of distribution.

#### Analyzing Bivariate Relationships

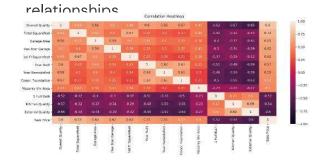
#### **Scatter Plots**

Display the relationship between two continuous variables, showing trends, clusters, and outliers.



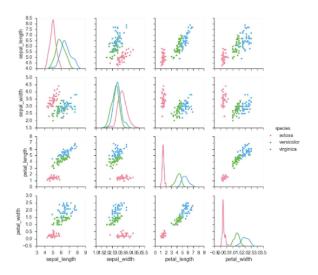
#### Heatmaps

Represent correlation between multiple variables using color intensity to indicate the strength of



#### Pair Plots

Show pairwise relationships between all variables in a dataset, providing a comprehensive view of correlations and distributions.





#### Handling Missing Data

1 2 3

#### Identification

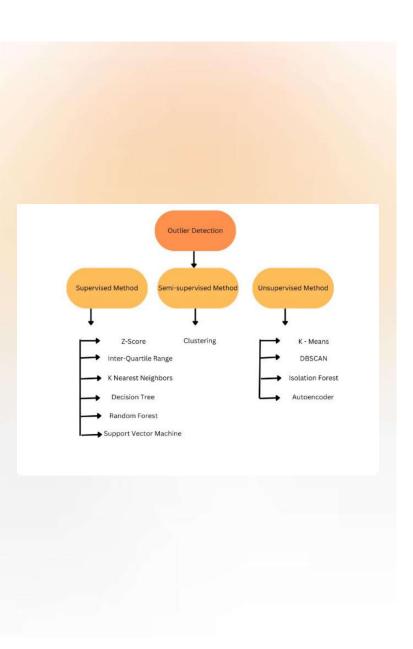
First, identify missing values in the dataset using techniques like summary statistics or visualization.

#### Imputation

Replace missing values with reasonable estimates using methods like mean, median, mode, or more sophisticated algorithms.

#### Removal

Remove rows or columns containing missing values if the amount is negligible or if the missing data cannot be reliably imputed.



#### **Detecting Outliers and Anomalies**



#### **Box Plots**

Identify outliers visually by observing points beyond the whiskers of the box plot.



#### Z-score

Calculate the standardized score for each data point and identify outliers as those with scores exceeding a certain threshold (e.g., 3).



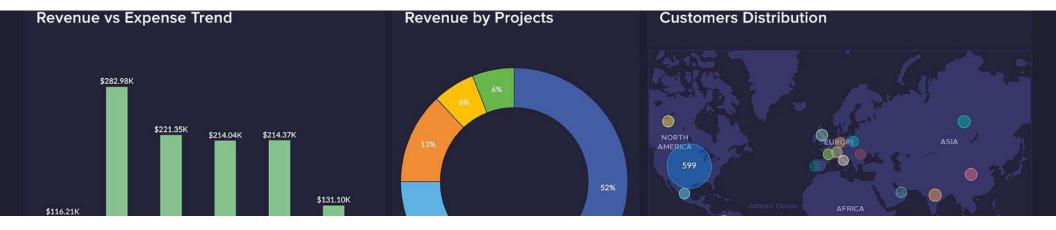
#### **Scatter Plots**

Detect outliers in bivariate relationships by identifying points that deviate significantly from the general pattern.



#### **Clustering Algorithms**

Use clustering algorithms to group data points and identify outliers as those not belonging to any cluster or forming small, distinct clusters.



#### Transforming and Scaling Data

#### Standardization

Scales data to have zero mean and unit variance, making features comparable regardless of their original units.

#### Normalization

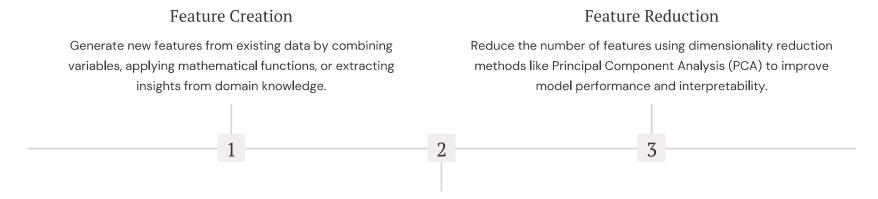
Transforms data to a specific range (e.g., O to 1), useful for algorithms sensitive to scale or distance metrics.

#### Log Transformation

Compresses the range of values and reduces the impact of extreme values, often used for skewed data.



#### Feature Engineering and Selection

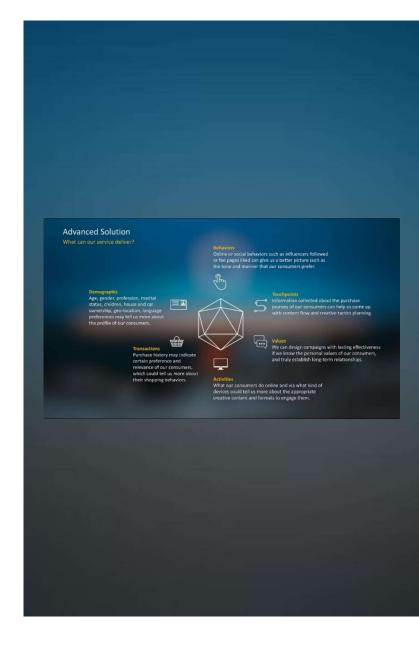


#### **Feature Selection**

Identify the most relevant features for the model using techniques like correlation analysis, feature importance, or statistical tests.

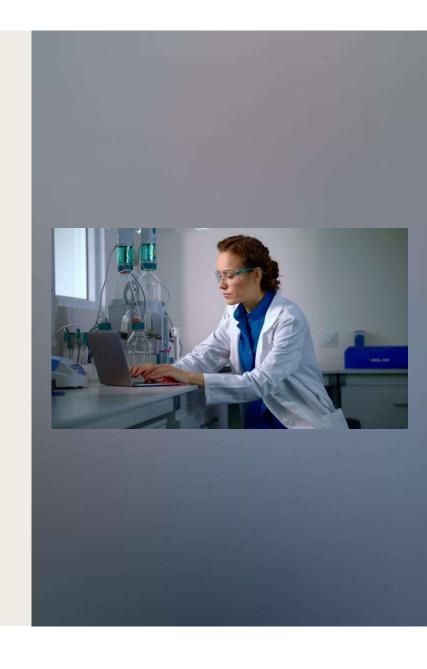
## Communicating Insights from EDA

Clearly communicate the findings from EDA using visualizations, tables, and narrative descriptions. This includes summarizing key insights, highlighting patterns and anomalies, and presenting actionable recommendations for further analysis or model building. Effective communication ensures that stakeholders can understand the data and its implications.



## Introduction to Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is a crucial initial step in any data science project. It's a process of investigating and understanding data before diving into complex modeling or prediction tasks. Think of EDA as a detective's initial examination of a crime scene: you're looking for clues, patterns, and insights that will guide your investigation further. EDA allows you to get familiar with the data, identify its characteristics, and uncover hidden relationships that can be valuable for subsequent analyses.





#### Importance of EDA in Data Science

#### 1 Data Quality

EDA helps identify data quality issues such as missing values, outliers, and inconsistencies.

Addressing these issues early on ensures the reliability and accuracy of subsequent analyses.

#### **3** Feature Engineering

EDA often leads to the creation of new features or transformations of existing features, enhancing the predictive power of machine learning models.

#### Insights and Understanding

EDA uncovers hidden trends, patterns, and relationships in data that might not be immediately obvious. These insights provide a deeper understanding of the data and inform decision-making.

#### **Model Selection**

EDA helps determine the most suitable models for the data by providing insights into the data's structure and distribution.

#### **Objectives of EDA**

#### **Understand Data Structure**

EDA aims to understand the basic structure and organization of the data, including its variables, data types, and relationships between variables.

#### **Assess Data Quality**

EDA helps assess the quality of the data, looking for missing values, outliers, inconsistencies, and other issues that could impact the reliability of analysis.

#### **Prepare Data for Modeling**

EDA transforms and prepares data for subsequent modeling and prediction tasks, ensuring it's in a suitable format for analysis.

#### **Identify Data Patterns**

EDA seeks to identify patterns, trends, and anomalies within the data, uncovering insights that might not be apparent at first glance.

#### **Generate Hypotheses**

EDA generates hypotheses about the data, suggesting potential relationships and areas for further investigation.

#### **Identifying Data Characteristics**

#### **Variables**

EDA begins by identifying the variables in the dataset and their respective data types (e.g., numerical, categorical, textual). This helps understand the nature of the data and its potential uses.

#### **Data Distribution**

EDA analyzes the distribution of each variable, using measures like mean, median, standard deviation, and histograms. This reveals the central tendency and spread of the data.

#### Missing Values

EDA identifies missing values, determining their frequency and potential causes. This is crucial for addressing data quality issues and ensuring accurate analysis.

#### **Detecting Outliers and Anomalies**

#### **Box Plots**

1 Box plots

Box plots visually represent the distribution of data, highlighting potential outliers that fall outside the expected range.

#### **Scatter Plots**

2 Scatter plots

4

Scatter plots reveal relationships between variables and can identify points that deviate significantly from the overall trend, indicating outliers.

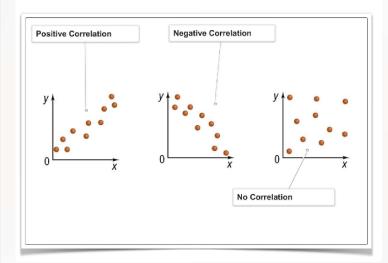
#### **Z-Scores**

Z-scores measure how far a data point is from the mean in terms of standard deviations. Outliers often have significantly high or low Z-scores.

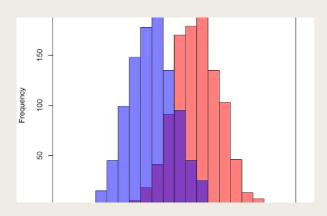
#### Domain Knowledge

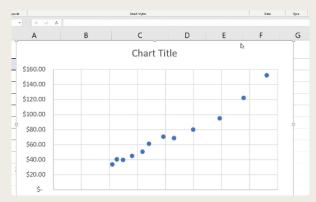
Domain expertise is crucial for interpreting outliers. Sometimes, outliers are not errors but rather legitimate data points that require further investigation.

You can use scatter plots to find trends in data. The scatter plots below show the three types of relationships that two sets of data may have.



#### **Visualizing Data Patterns**







#### Histograms

Histograms depict the frequency distribution of a single variable, showing the concentration of data points across different ranges.

#### **Scatter Plots**

Scatter plots visualize the relationship between two variables, revealing patterns of correlation, clusters, and outliers.

#### **Line Charts**

Line charts represent data trends over time or a continuous variable, highlighting changes and patterns across the data.



#### **Assessing Data Quality**

Missing Values	Identifying and addressing missing values ensures data completeness and reliability for
Outliers	analysis. Detecting and handling outliers prevents skewed results and ensures data accuracy.
Data Types	Confirming data types (e.g., numerical, categorical) ensures correct analysis and
Consistency	interpretation. Checking for inconsistencies within the data helps maintain data integrity and accuracy.

# | The second | The

## Discovering Relationships between Variables

#### Correlation

EDA explores the relationships between variables, measuring correlation coefficients to quantify the strength and direction of linear

#### **Scatter Plots**

Scatter plots visually represent the relationships between variables, revealing patterns of association, trends, and outliers.

#### Heatmaps

Heatmaps provide a visual representation of the correlation matrix, highlighting strong correlations and revealing potential relationships between variables.

## **Generating Hypotheses for Further Analysis**



#### **Insights from EDA**

EDA generates insights that lead to hypotheses about the data, suggesting potential relationships and areas for further exploration.



#### **Questions and Exploration**

EDA helps formulate specific research questions and guide further investigation to confirm or refute the initial hypotheses.



#### **Guided Research**

Hypotheses generated from EDA provide direction for more focused analysis, leading to deeper understanding and potentially new discoveries.





#### **Preparing Data for Modeling and Prediction**

#### 1 Data Cleaning

EDA identifies and addresses data quality issues such as missing values, outliers, and inconsistencies, preparing the data for modeling.

#### 3 Data Transformation

EDA may involve data transformation, such as normalization or scaling, to optimize the performance of machine learning algorithms.

#### 2 Feature Engineering

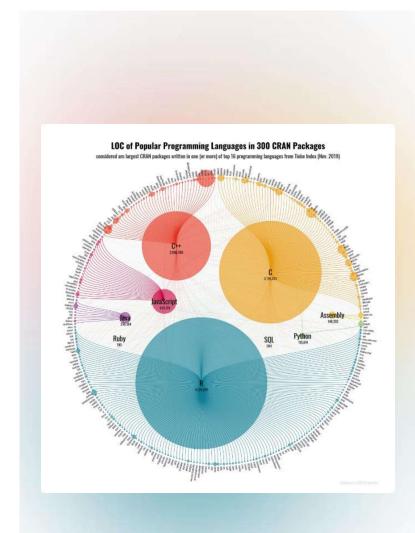
EDA often leads to the creation of new features or transformations of existing features, enhancing the predictive power of machine learning models.

#### 4 Model Selection

EDA helps determine the most suitable models for the data by providing insights into the data's structure and distribution.

## Unveiling the Power of Exploratory Data Analysis

Exploratory Data Analysis (EDA) is a crucial step in any data science project. It involves a systematic and iterative process of examining and summarizing data to gain insights, discover patterns, and formulate hypotheses. EDA helps you understand the underlying structure of your data, identify potential problems, and guide your subsequent data analysis and modeling efforts.



#### Importance of EDA in Data Science

EDA is more than just a preliminary step; it's a foundation for building robust and meaningful insights. It allows you to:

Identify Data Quality Issues

EDA helps you detect missing values, outliers, inconsistent data formats, and other anomalies that can affect your analysis.

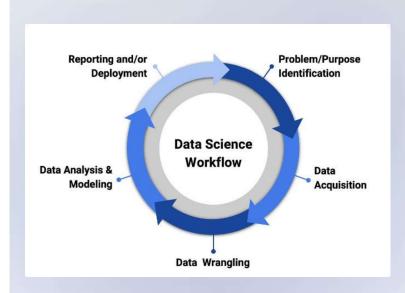
**3** Guide Feature Engineering

By understanding the relationships between variables, you can derive new features that improve the performance of your machine learning models. 2 Discover Hidden Relationships

EDA can reveal unexpected correlations and patterns that may not be immediately apparent from simply looking at the raw data.

**4** Support Hypothesis Generation

EDA provides evidence-based insights that can help you formulate hypotheses that you can later test with more rigorous statistical methods.



#### **Identifying Interesting Datasets for EDA**

The right dataset is the key to a successful EDA project. Look for datasets that are:

#### **Relevant to Your Interests**

Choose a dataset that aligns with your professional goals or personal interests. This will keep you engaged and motivated throughout the project.

#### **Available and Accessible**

Ensure that the dataset is readily available and can be downloaded or accessed through APIs. Consider the file format and size, ensuring it's manageable.

#### **Well-Documented**

Datasets with clear documentation are easier to understand and analyze. Look for information about variable definitions, data sources, and any known biases or limitations.

### Assessing the Complexity and Scope of the Dataset

Once you've chosen a dataset, assess its complexity and scope. This will help you plan your EDA process and determine the appropriate techniques to use.

Size

Is the dataset small enough to be processed on your computer, or will you need cloud computing resources?

Number of Variables

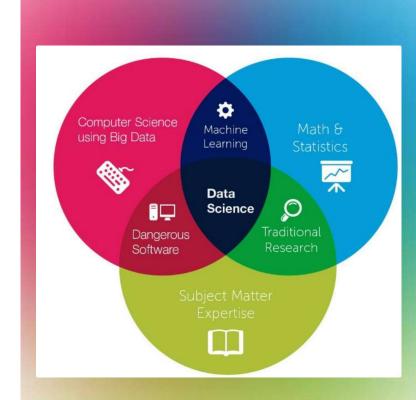
How many variables are in the dataset? A large number of variables may require dimensionality reduction techniques.

Z Data Types

What are the data types of the variables? Are they numerical, categorical, or a mix of both?

Missing Values

How many missing values are there in the dataset? This will affect your analysis and may require imputation techniques.



### Defining Clear Objectives for Your EDA Project

Before diving into the data, define specific objectives for your EDA project. This will provide direction and focus, helping you avoid getting lost in the vast amount of data.

#### Understand the Distribution of Variables

Analyze the central tendency, dispersion, and shape of distributions. This can reveal outliers, skewness, and other characteristics.

#### **Identify Key Insights and Patterns**

Discover interesting trends, anomalies, and relationships that can inform your understanding of the data and guide further analysis.

#### **Explore Relationships Between** Variables

Investigate correlations, dependencies, and interactions between different variables in the dataset.

#### **Generate Hypotheses for Future Investigation**

Formulate hypotheses based on your findings that can be tested using more rigorous statistical methods.

	Question-flow	Real life example: Performance Ratio (PR)
QUALITIVE AREA	What do I want to accomplish?	I want to improve my plant performance
	For what reason?	So it can produce more and make more money!
	Can I state it as a numerical question?	Can I calculate the current performance of my plant (first step to understand how to improve it)?
	How can I calculate it? Do we already have a function, a KPI, a variable addressing the numerical question?	How do I calculate the PR of my plant?
	What is the accuracy and/or the tolerance I need for that numerical result?	For my PR (from 0% to 100%), I'm satisfied with an accuracy of 1% expressed as absolute error.
AREA	Once I have the final result, how can I validate it?	Once I have my result, I can compare it with a similar plant where this analysis was already done by external assessments. Or with my competitor's (remember: The grass is always greener on the other side!).
	Do I have the technical tools to calculate it?	For PR calculation, I need power production data and weather (solar irradiance, temperature) data. Also, I need a calculator, because my math is rusty.
	Do I have the time to do it to meet expectations?	For PR calculation, I would need 2-3 days, and I have time until the next report sharing. (Oh, that's 2 week from now, I'm cool!)
RELATIONSHIP AREA	Can I ask for a double check or for a suggestion to someone I trust?	For double check my activity, I can ask to my Southern branch colleague, they have been doing this for a long time now.



#### **Exploring Data Visualization Techniques**

Data visualization is an essential part of EDA. It allows you to explore and understand data in a more intuitive and engaging way. Use a variety of visualization techniques to explore different aspects of the data.

1 2 3 4

#### **Histograms**

Show the distribution of numerical variables.

#### **Scatter Plots**

Illustrate relationships between two numerical variables.

#### **Box Plots**

Compare the distribution of a variable across different groups or categories.

#### **Heatmaps**

Visualize correlations between multiple variables.



## **Uncovering Insights and Patterns** in the Data

By applying visualization techniques and summary statistics, you can uncover hidden insights and patterns within the data. Look for:



## **Trends and Patterns**

Identify upward or downward trends, cyclical patterns, or other recurring patterns in the data.



## **Outliers and Anomalies**

Investigate data points that deviate significantly from the expected pattern.

These could be errors or interesting insights.



## **Correlations**

Determine how variables are related to each other. Are they positively correlated, negatively correlated, or independent?



## **Clusters and Groups**

Discover groups of data points that share similar characteristics. This could indicate underlying subgroups or categories.

## **Communicating Findings Effectively**

Once you've uncovered insights, you need to effectively communicate them to your stakeholders. Use a combination of text, visualizations, and narratives to tell a compelling story about the data.

Use Clear and Concise Language	Avoid technical jargon and explain concepts in a way that is easy to
Focus on Key Findings	understand. Highlight the most important insights and patterns that you've discovered.
Use High-Quality Visualizations	Create visually appealing and informative graphs, charts, and tables that effectively communicate your
Provide Context and Interpretation	findings. Explain the meaning of your findings and their implications for decision-making.



## LEONARD HOFSTADTER

· +1 (456) 1234567

■ leonard@hiration.com

SF, CA

## Data Scientist

## SUMMARY

7+ years experienced data scientist with a passion to solve real-world business challenges using data analytics. Track record of setting up the Data Science Div. for a leading hospitality firm & rendering consultancy services for a Fortune 500 company. Proficient in deploying complex machine learning and statistical modeling algorithms/techniques for identifying patterns and extracting valuable insights for key stakeholders and organizational leadership.

SE US

## PROFESSIONAL EXPERIENCE

## Positronix Financial Services

**Data Scientist** 

Cct '15 - Present

An Al-based start-up which specializes in delivering solutions to client in the finance domain

Technology Stack: Python, Hadoop, AWS, Pandas, NumPy,

## **Data Visualization & Predictive Analytics** Steering rapid model creation in Python using Pandas.

- NumPy, SciKit-Learn & plot.ly for data visualization Creating NLP models for Sentiment Analysis & MapReduce modules for predictive analytics in Hadoop

## **Key Achievements**

 Established the Data Science division from scratch by recruiting, on-boarding & training a team of 8 Data

## **Epiplace Solutions**

solutions to 1,000+ clients

## Consultant - Data Analytics

feb '13 - Sep '15 With a presence in 20+ cities, Epiplace has delivered cost-effective II

Technology Stack: Python, Pandas, NumPy, SciKit-Learn, Matplotlib, Jupyter Notebook

- Segmentation & Clustering · Applied large scale & low latency machine learning for non-parametric models & high-dimensional data
- · Created multivariate regression-based attribution models & segmentation models using K-means Clustering

 Developed an additive scoring model for OSM and a logistic regression model to yield a K-S statistic of 51.5

## PeopleSoft

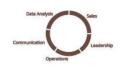
## Data Analyst

Redmond, US

m Jun '11 - Jan '13

## PeopleSoft is one of the largest software developed Data Analytics, Model Development & ML Algorithms

- · Directed model development, validation, testing and implementation of analytical products and applications
- · Deployed advanced text mining algorithms to identify search intent latent in individual keywords



## TECHNICAL SKILLS

- Packages: SciKit-Learn, NumPy, SciPy, Plot.ly, Pandas, NLTK, Beautiful Soup, Matplotlib
- · Big Data Stack: Hadoop, Apache, Pig. Python, PostgreSQL, AWS, Hive, MongoDB, MapReduce, Spark
- Statistics/ML: Linear/Logistic Regression, SVM, Ensemble Trees Random Forests

## **KEY SKILLS**

- Data Analysis
- Stakeholder Management
- · Leadership & Training · Strategy · Project Management & Delivery
- · Process Improvement
- Team Incubation
- Data Visualization • Predictive Modelling & Analytics

## **EDUCATION**

## **UC Berkeley**

## BS in Data Science

m Jul '07 - May '11

Berkeley, US

## **Showcasing Your EDA Skills on Your** Resume

EDA is a valuable skill to highlight on your resume, demonstrating your ability to analyze data and extract meaningful insights. Showcase your EDA skills by:

## **Including Relevant Project** Experience

Describe projects where you applied EDA techniques to solve real-world problems. Highlight the specific techniques you used, the insights you uncovered, and the impact of your analysis.

## **Quantifying Your Results**

Use metrics to demonstrate the value of your EDA work. For example, highlight the improvement in model performance or the number of new insights discovered.

## **Sharing Your Work** 3

Consider showcasing your EDA projects on a portfolio website or on platforms like GitHub. This provides potential employers with a way to see your work firsthand.

## **Highlight Your Skills**

In your resume, list specific EDA skills such as data cleaning, visualization, statistical analysis, and hypothesis testing.

## **Conclusion and Next Steps**

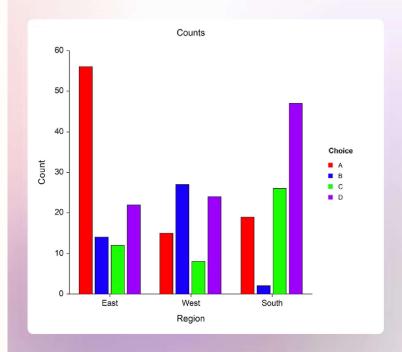
EDA is a powerful tool for exploring data and uncovering hidden insights. By following the steps outlined above, you can conduct effective EDA projects that will enhance your data science skills and make you a more valuable asset to any organization.

As you continue to explore the world of data science, consider these next steps:

- Explore advanced EDA techniques like dimensionality reduction and feature engineering.
- Learn more about different data visualization libraries and tools.
- Practice applying EDA to real-world datasets and participate in data science competitions.
- Network with other data scientists and share your insights and experiences.

## Introduction to Exploratory Data Analysis (EDA) Projects

Exploratory Data Analysis (EDA) is a crucial step in any data science project. It involves analyzing and visualizing raw data to uncover patterns, trends, and insights that can inform further analysis and decision-making.



## EDA for Customer Churn Analysis

## Identifying Key Drivers of Churn

EDA helps identify key factors contributing to customer churn, such as service issues, pricing concerns, or lack of engagement.

- Customer demographics
- Service usage patterns
- Customer feedback

## Understanding Churn Behavior

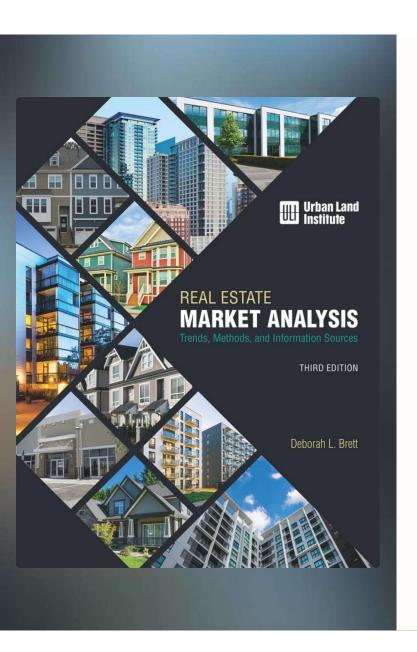
EDA reveals patterns in customer churn, such as churn rates over time, churn profiles, and relationships between churn and other variables.

# TYPES OF DATA VISUALIZATION CHARTS Area Chart Display trends over time A line chart with areas below the lines filled with colors Bubble Chart Show correlation in a dataset Show the labelled circles Map Show data with location as a variable Show magnitude of a phenomenon

## Developing Targeted Retention Strategies

Insights from EDA inform the development of targeted interventions to reduce churn, such as personalized offers, loyalty programs, or proactive customer





## EDA for Predicting Housing Prices

1 Exploring Housing Market Trends

EDA helps analyze historical housing price data to identify trends, seasonality, and relationships between prices and other factors.

3 Building Predictive Models 4

Insights from EDA guide the development of predictive models to estimate housing prices based on relevant features.

2 Identifying Key Features Influencing Prices

EDA determines the most significant features impacting housing prices, such as location, size, amenities, and neighborhood characteristics.

Evaluating Model Performance

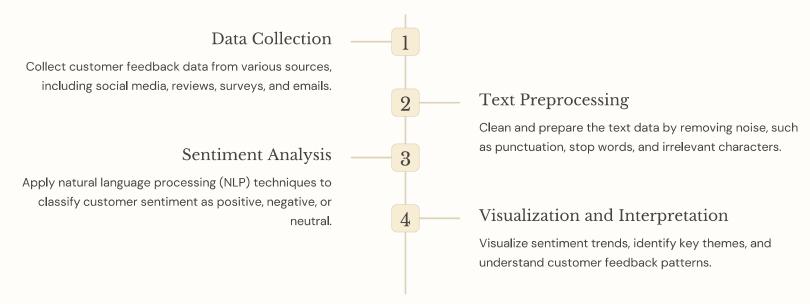
EDA helps evaluate the accuracy and reliability of predictive models by comparing predictions with actual housing prices.







## EDA for Analyzing Customer Sentiment



orientation=horizontal twosided=true



## EDA for Detecting Fraud in Financial Transactions

## Transaction Anomaly Detection

EDA helps identify unusual transaction patterns, such as large transactions, multiple transactions in short intervals, or transactions from unexpected locations.

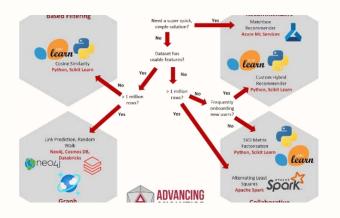
## User Behavior Analysis

EDA analyzes user behavior to detect suspicious activities, such as sudden changes in spending patterns, multiple account accesses from different locations, or unusual login attempts.

## Network Analysis

EDA helps identify fraudulent networks by analyzing connections between accounts, transactions, and individuals.

## EDA for Optimizing E-commerce Recommendations





## Understanding Customer Preferences

EDA analyzes customer browsing history, purchase history, and product ratings to understand their preferences and interests.

## Identifying Product Relationships

EDA helps identify relationships between products, such as frequently purchased together, viewed together, or rated similarly.

## Developing Recommendation Algorithms

Insights from EDA guide the development of personalized recommendation algorithms to suggest products based on customer preferences and product relationships.

## Time Series Modeling Visualize, wrangle, and preprocess time series data Forecast Plot Legend Lege

## EDA for Forecasting Time Series Data

## **Data Preparation**

Prepare the time series data by cleaning, transforming, and aggregating data points.

## Trend and Seasonality Analysis

Identify trends, seasonality, and cyclical patterns in the time series data.

3

## Forecasting Models

Select appropriate forecasting models based on the characteristics of the time series data, such as ARIMA, exponential smoothing, or neural networks.

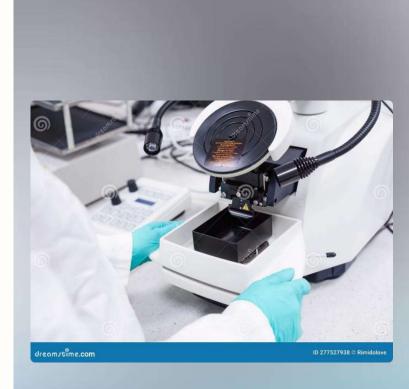
Model Evaluation and Optimization

Evaluate the performance of the forecasting models and optimize them to achieve the best accuracy.

orientation=horizontal

## EDA for Identifying Patterns in Biological Data

Gene Expression Analysis	Identify genes that are differentially expressed between different conditions, such as disease states or treatment groups.
Protein Interaction Networks	Analyze protein-protein interactions to understand biological pathways and identify potential drug targets.
Microbial Community Analysis	Study the composition and diversity of microbial communities to understand their role in health, disease, and environmental processes.





## EDA for Improving Supply Chain Efficiency



## Inventory Management

Analyze inventory levels, demand patterns, and lead times to optimize inventory management and reduce stockouts.



## Distribution Network Optimization

Identify bottlenecks in the distribution network, optimize routing and delivery schedules, and reduce transportation costs.



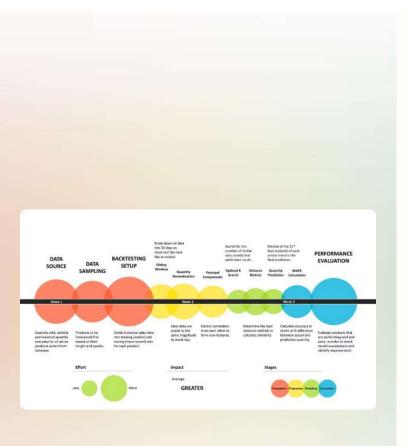
## **Demand Forecasting**

Forecast future demand based on historical sales data, seasonality, and other factors to improve production planning and inventory management.



## **Production Optimization**

Analyze production data to identify inefficiencies, improve production scheduling, and reduce production costs.



## Conclusion and Key Takeaways

EDA is an essential tool for data-driven decision-making. By uncovering hidden patterns and insights, EDA helps businesses understand their data, make informed decisions, and achieve their goals.