



AI-Powered Analytics - Generating interactive insights over unprocessed data

ANA PATRÍCIA CORREIA PADEIRO

(Licenciada)

Trabalho de Projeto para obtenção do Grau de Mestre em
Engenharia Informática e Multimédia

Orientadores: Mestre Ana Luísa Gonçalves Neves Gomes
Doutor Gonçalo Caetano Marques

Júri:

Presidente: Doutor Pedro Mendes Jorge

Vogais: Doutor André Ribeiro Lourenço

Doutor Gonçalo Caetano Marques

Novembro 2024



AI-Powered Analytics - Generating interactive insights over unprocessed data

ANA PATRÍCIA CORREIA PADEIRO

(Licenciada)

Trabalho de Projeto para obtenção do Grau de Mestre em
Engenharia Informática e Multimédia

Orientadores: Mestre Ana Luísa Gonçalves Neves Gomes, Millennium BCP
Doutor Gonçalo Caetano Marques, ISEL

Júri:

Presidente: Doutor Pedro Mendes Jorge, ISEL

Vogais: Doutor André Ribeiro Lourenço, ISEL
Doutor Gonçalo Caetano Marques, ISEL

Novembro 2024

"Success is not final, failure is not fatal: It is the courage to continue that counts."– Winston Churchill

Acknowledgements

I would like to start by thanking my supervisors, Luísa Gomes and Gonçalo Marques, for their tireless support throughout this process. Your availability and willingness to assist were essential to the success of this study. The experience and knowledge you shared with me were truly crucial. Thank you for all your dedication and for being so committed to this project; I couldn't have asked for better supervisors.

To my company, for allowing me to balance my master's studies with work and for giving me the opportunity to develop a project aligned with my research interests. To my team, who consistently showed concern for my well-being and the progress of my work, and to my colleagues who have been incredibly supportive — Vanessa Bernardes, Leonor Santos, Afonso Machado, Yamin Yassin, Sofia Cotter, Hugo Brito, and Bruno Lampreia — my deepest thanks for your patience and encouragement.

To my family, who always believed in me and supported my dreams. My parents, who have been an essential pillar, always offering care and concern throughout my life. Thank you for allowing me to study what I love most. To my sister, who, despite the headaches (as only siblings can give!), is the person who endures my crises and celebrates my achievements. Once again, you were fundamental on this journey.

To the friends I gained at ISEL, especially Ana Ferreira, who started this adventure with me at ISEL and joined me in the master's program. Here we are, about to conclude this great chapter. To Rita Marques and Beatriz Santos, who were always there for me during my moments of crisis — many tears were shed at ISEL, and because of ISEL. To the many others I met, such as João Alves, João Novo, Joana Pires, Pedro Tomar, Bernardo Jaco, Cláudia Fernandes, Joana Carrilho — I could keep naming people because, fortunately, this institute not only helped me grow as a person but also allowed me to meet incredible individuals.

To my long-time friends — João, Érica, Joana, Rafael, Filipe, and Pedro — your friendship has undoubtedly been my escape from my studies. When I needed to pretend nothing existed, I knew you were always there to support me. And you put up with me — which was certainly no small feat this year. Thank you for being my refuge.

Lastly, to ISEL, which has been my home over these past years and where I grew not only academically but also personally. People often say that college years are the hardest, and while I don't deny the challenges, for me, they were also the best years of my life. Thank you to all the people — both faculty and staff — who were part of this journey and supported me along the way.

Abstract

In the current context, where the amount of data grows exponentially, organizations face increasing challenges in managing large volumes of unprocessed data. With this increase and technological advances, traditional data analysis methods have become insufficient to extract insights efficiently and in real time.

The main objective of this project is to develop an analysis pipeline based on Artificial Intelligence (AI), capable of transforming raw data into insights that support informed decision-making. This pipeline covers the entire data lifecycle, from its ingestion and preparation, through the creation of Machine Learning (ML) models, to the generation of interactive visualizations. Furthermore, Natural Language Processing (NLP) techniques and an easy-to-use interface are integrated, with the aim of maximizing the efficiency and accessibility of results.

An essential part of this work was the study of the integration of the BERTopic model, an advanced topic modeling technique that uses language embeddings for the automatic extraction of significant topics in large volumes of unstructured textual data. This approach makes it possible to identify linguistic patterns and trends, providing valuable insights into areas such as customer service, risk management and anomaly detection.

To ensure maximum use of data, features such as automatic identification of data types (numeric, categorical, boolean, etc.), treatment of missing values, vectorization, detection of outliers and automated generation of statistical summaries were implemented. Furthermore, normalization and automatic selection techniques for features were used, allowing the complexity of the data to be reduced and ensuring that the insights generated were robust and accurate.

In order to carry out a detailed analysis of patterns and trends in the data, ML techniques were applied, specifically unsupervised learning, through clustering algorithms. This enabled practical recommendations that leverage the extracted insights, with the aim of optimizing areas such as customer service, risk management and anomaly detection.

Keywords: Machine Learning, Big Data, Data Analytics, Unsupervised Learning, AI-Powered Analytics, Data Preprocessing, BERTopic, Natural Language Processing, Clustering Algorithms, Feature Selection, Decision-Making, Operational Efficiency, Data Insights.

Resumo

No contexto atual, onde a quantidade de dados cresce exponencialmente, as organizações enfrentam desafios cada vez maiores na gestão de grandes volumes de dados não processados. Com esse aumento e os avanços tecnológicos, os métodos tradicionais de análise de dados tornaram-se insuficientes para extrair *insights* de forma eficiente e em tempo real.

O principal objetivo deste projeto é desenvolver uma *pipeline* de análise baseada em Inteligência Artificial (IA), capaz de transformar dados brutos em *insights* que apoiem a tomada de decisões informadas. Esta *pipeline* abrange todo o ciclo de vida dos dados, desde a sua ingestão e preparação, passando pela criação de modelos de *Machine Learning* (ML), até à geração de visualizações interativas. Além disso, integram-se técnicas de Processamento de Linguagem Natural (PLN) e uma interface de fácil utilização, com o objetivo de maximizar a eficiência e acessibilidade dos resultados.

Uma parte essencial deste trabalho foi o estudo sobre a integração do modelo BERTopic, uma técnica avançada de modelagem de tópicos que utiliza *embeddings* de linguagem para a extração automática de tópicos significativos em grandes volumes de dados textuais não estruturados. Esta abordagem possibilita a identificação de padrões e tendências linguísticas, fornecendo *insights* valiosos para áreas como atendimento ao cliente, gestão de riscos e detecção de anomalias.

Para garantir o máximo aproveitamento dos dados, foram implementadas funcionalidades como a identificação automática de tipos de dados (numéricos, categóricos, booleanos, etc.), tratamento de valores ausentes, vetorização, detecção de *outliers* e geração automatizada de resumos estatísticos. Além disso, foram utilizadas técnicas de normalização e seleção automática de *features*, permitindo reduzir a complexidade dos dados e assegurar que os *insights* gerados fossem robustos e precisos.

Com o intuito de realizar uma análise detalhada de padrões e tendências nos dados, foram aplicadas técnicas de ML, especificamente aprendizagem não-supervisionada, através de algoritmos de *clustering*. Isso possibilitou recomendações práticas que aproveitam os *insights* extraídos, com o objetivo de otimizar áreas como atendimento ao cliente, gestão de riscos e detecção de anomalias.

Palavras-chave: Machine Learning, Big Data, Análise de Dados, Aprendizagem Não Supervisionada, Análise Impulsionada por IA, Pré-processamento de Dados, BERTopic, Processamento de Linguagem Natural, Algoritmos de Clustering, Seleção de Variáveis, Tomada de Decisões, Eficiência Operacional, Insights de Dados.

Contents

List of Figures	xvii
Listings	xix
Acronyms	xxi
1 Introduction	1
1.1 Context	2
1.2 Main Goals	5
1.3 Contributions	5
1.4 Document Organization	6
2 Background and Related Work	9
2.1 Data Science and Big Data	10
2.1.1 Introduction to Big Data	11
2.1.2 Challenges of Big Data	12
2.2 Machine Learning	12
2.2.1 Supervised Learning	14
2.2.2 Unsupervised Learning	15
2.2.3 Reinforcement Learning	15
2.2.4 Semi-Supervised Learning	16
2.3 Definition of Insight	16
2.3.1 Role of Insights in Modern Data-Driven Environments	16
2.3.2 How Insights are Derived	17
2.3.3 Example of Insight Generation	18
2.4 Data Pre-processing Techniques	19
2.4.1 Handling Missing Data	19
2.4.2 Outlier Detection	23
2.4.3 Normalization and Standardization	28
2.4.4 Vectorization	29
2.4.5 Feature Selection	31
2.4.6 Dimensionality Reduction	35
2.5 Natural Language Processing and Topic Modeling	37
2.6 Interpretability and Explainability	39
2.7 Clustering Methods	40

2.7.1	Centroid-based Clustering	41
2.7.2	Density-based Clustering	44
2.7.3	Distribution-based Clustering	46
2.7.4	Hierarchical Clustering	46
2.8	Visualization Libraries	49
2.9	Visualization of Model Results	51
3	Data Description	53
3.1	Data Source and Nature	53
3.2	Data Structure and Analysis	53
3.2.1	Structure and Analysis of the Retail Transactional Dataset	53
3.2.2	Structure and Analysis of the Credit Card Dataset	55
3.2.3	Structure and Analysis of the Investment Preferences Dataset	56
4	Implementation	59
4.1	Apache Spark	59
4.2	Databricks	59
4.3	Functions and Methodologies	60
4.3.1	File Handling and Conversion	60
4.3.2	Data Type Identification and Conversion	61
4.3.3	Handling Null Values	62
4.3.4	Vectorization	63
4.3.5	Outlier Detection	64
4.3.6	Correlation	65
4.3.7	Feature Selection	65
4.3.8	Clustering	66
4.3.9	Visualization Methods	68
5	Testing and Validation	71
5.1	Testing Scope and Challenges	71
5.2	Data Preprocessing	71
5.2.1	Handling Missing Data	72
5.2.2	Outlier Detection and Removal	72
5.3	Feature Selection	74
5.4	Cluster Analysis Using Various Methods	76
5.4.1	Optimal Parameter Determination	76
5.4.2	Clustering Methods and Results	77
5.5	BERTopic Analysis and Results	83
5.5.1	Results and Insights	83
5.5.2	Interpretation of Results	84
5.6	Summary	85
6	Conclusions and Future Work	87
6.1	Key Insights and Challenges	87
6.2	Future Work	88

6.3 Final Reflection	89
Bibliography	91

List of Figures

1.1	Applications of ML in Various Sectors	1
1.2	Overview of the ML Process	2
1.3	Comparison of Structured and Unstructured Data	4
1.4	Example of Interactive Data Visualization Dashboard	4
2.1	The 5Vs of Big Data	11
2.2	Overview of Machine Learning Categories: Supervised, Unsupervised, Reinforcement, and Semi-Supervised Learning	14
2.3	Schematic Representation of Supervised Learning Process: Labeled Data Leading to Predictions	14
2.4	Clustering Process in Unsupervised Learning: Grouping Data Points Based on Similarities	15
2.5	Process of transforming data into insights	18
2.6	Correlation between educational level and loan default rates	18
2.7	Forward and Backward Fill Imputation	20
2.8	KNN Imputation	21
2.9	Regression Imputation	22
2.10	Autoencoders Imputation	23
2.11	Isolation Forest method	26
2.12	Outlier detection using DBSCAN: points outside of dense clusters are considered outliers	27
2.13	The creation of vectors	30
2.14	Example of PCA	36
2.15	Example of t-SNE technique	36
2.16	Clustering Process	40
2.17	Examples of Different Clustering Algorithms	41
2.18	Initial Data and Centroids	42
2.19	Assigning Points to Clusters	42
2.20	Updating Centroids	42
2.21	Comparison of standard K-Means and Mini-Batch K-Means	43
2.22	Types of points in the DBSCAN algorithm	45
2.23	Example of the DBSCAN algorithm in action	45
2.24	Example of a Dendrogram	46
2.25	Comparison of Agglomerative and Divisive Clustering	47
2.26	BIRCH Clustering Process	49
2.27	Graph visualization using the Matplotlib library	49
2.28	Graph visualization using the Plotly library	50

2.29	Graph visualization using the Bokeh library	50
2.30	Graph visualization using the Altair library	51
3.1	Distribution of Data Types in the Retail Transaction Dataset	54
3.2	Distribution of Data Types in the Credit Card Dataset	55
3.3	Distribution of Data Types in the Investment Preferences Dataset	56
5.1	Outlier Detection and Removal Process using Box Plot Example	73
5.2	Histograms of BALANCE and PURCHASES Variables Before and After IQR-Based Outlier Removal	74
5.3	Optimal K Determination using the Elbow Method	76
5.4	Clusters generated by K-means	78
5.5	Quantile vs. Number of Clusters	78
5.6	Clusters from the Mean Shift algorithm	79
5.7	Clusters from the Birch algorithm	80
5.8	Dendrogram Showing Hierarchical Clustering and Optimal Number of Clusters	81
5.9	Clusters from the Agglomerative Clustering algorithm	82
5.10	Affinity Propagation Clustering	83

Listings

4.1	Pseudocode for function <code>get_file_type()</code>	60
4.2	Pseudocode for function <code>to_spark()</code>	61
4.3	Pseudocode for function <code>get_type_of_column()</code>	61
4.4	Pseudocode for function <code>convert_column_types()</code>	62
4.5	Pseudocode for function <code>handle_missing_data()</code>	63
4.6	Pseudocode for function <code>vectorize_column()</code>	64
4.7	Pseudocode for function <code>detect_outliers()</code>	64
4.8	Pseudocode for function <code>select_features()</code>	66
4.9	Pseudocode for function <code>find_optimal_k()</code>	67
4.10	Pseudocode for function <code>apply_clustering_methods()</code>	68

Acronyms

AI	Artificial Intelligence 1, 2, 4, 12, 31, 37, 39
BERTopic	Bidirectional Encoder Representations from Transformers Topic Modeling ix, xi, 9, 38, 39, 69
BFILL	Backward Fill 20, 21
BoW	Bag of Words 29, 30
CBOW	Continuous Bag of Words 31
DBSCAN	Density-Based Spatial Clustering of Applications with Noise xvii, 9, 27
FFILL	Forward Fill 20, 21
GloVe	Global Vectors for Word Representation 31
IoT	Internet of Things 2, 10, 11
IQR	Interquartile Range 24, 28
KNN	K-Nearest Neighbors xvii, 21, 28, 62
L2	L2 Normalization 29, 35
ML	Machine Learning ix, xi, xvii, 1, 2, 4, 5, 6, 9, 10, 12, 13, 14, 17, 18, 19, 23, 27, 28, 29, 31, 32, 33, 35, 39, 40, 51, 59, 60, 63, 64, 87
MSE	Mean Squared Error 15
NLP	Natural Language Processing 2, 4, 5, 6, 9, 18, 29, 31, 37, 38, 87
PCA	Principal Component Analysis 35, 69
TF-IDF	Term Frequency-Inverse Document Frequency 30

Chapter 1

Introduction

In today’s digital landscape, the exponential growth of data has transformed how organizations operate across various sectors. Modern enterprises increasingly rely on data-driven strategies to maintain a competitive edge and foster innovation [59]. For example, in the financial sector, transactional systems generate massive volumes of data daily, making fraud detection a critical challenge. Similarly, in healthcare, patient data must be analyzed in real-time to improve outcomes. With the rise of technologies such as the Internet of Things (IoT), mobile platforms, and cloud computing, the generation and collection of vast amounts of data have become ubiquitous. However, managing and extracting value from this data—particularly unstructured formats such as text, images, or videos—remains a significant challenge [73].

Data has become an indispensable asset in sectors like finance, healthcare, and retail, where critical decisions are based on large volumes of information collected from transactional systems, sensors, and user-generated content. Although gathering data is easier than ever, turning raw, unprocessed data into actionable insights is complex and resource-intensive. Traditional data analysis methods often struggle to handle the growing volume, velocity, and variety of data, especially when dealing with large-scale real-time data streams or complex unstructured data such as text from customer reviews, social media posts, and emails [16].



Figure 1.1: *Applications of ML in Various Sectors*

As shown in Figure 1.1, ML has a broad range of applications across industries, from finance to healthcare, showcasing its ability to handle diverse datasets and uncover valuable insights.

Unstructured data, which constitutes the majority of the information generated today, presents unique challenges due to its lack of predefined structure. To address these challenges, organizations are turning to AI and ML solutions to automate data processing and generate insights in real-time. These technologies allow businesses to efficiently process large datasets, identify patterns, predict outcomes, and generate valuable insights [35].

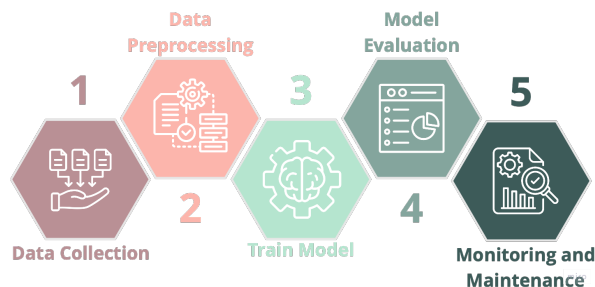


Figure 1.2: *Overview of the ML Process*

In Figure 1.2, an overview of the machine learning process is shown, which illustrates the core stages, including data collection, preparation, model training, and evaluation. These steps are critical for transforming raw data into actionable insights.

AI-powered systems streamline data analysis by automating many labor-intensive tasks [30]. They are particularly effective at handling both structured and unstructured data, helping organizations make faster, more informed decisions. The use of ML models enables companies to process data with greater accuracy, allowing for the identification of trends and the generation of actionable insights that drive strategic initiatives [35].

Despite the potential of AI and ML, developing AI-powered analytics pipelines is a multifaceted challenge [30]. These systems require expertise in a range of areas, including data ingestion, cleaning, model development, and visualization. Moreover, the integration of NLP is crucial for extracting insights from unstructured data sources [11], such as text, which are often difficult to analyze using traditional techniques. For example, applying NLP to customer feedback can uncover recurring complaints, enabling businesses to optimize services. However, this process demands significant computational resources and scalability to handle large datasets efficiently. Another key challenge lies in making the insights generated by these systems accessible to users, regardless of their technical background. To address this, the pipeline will incorporate interactive dashboards, enabling non-technical stakeholders to explore insights with ease.

Given these complexities, there is a growing need for scalable, automated solutions that streamline the data analysis workflow, from ingestion and preparation to insight generation and visualization. This is especially important in industries like e-commerce, where analyzing customer behavior patterns can drive personalized recommendations, or healthcare, where real-time data can improve patient care. This research focuses on developing an AI-powered analytics pipeline capable of transforming raw data into interactive, actionable insights that support informed decision-making across various domains [11].

1.1 Context

The ability to process and analyze large datasets has become a critical driver of success in many industries. Advances in IoT, mobile platforms, and cloud-based infrastructure have dramatically increased the volume of data generated. For example, in manufacturing, IoT sensors generate real-time data that requires swift analysis to optimize production lines. Similarly, in retail, unstructured data from social media can provide insights into customer sentiment and emerging trends. However, while access to data has become widespread, traditional analysis techniques often fall short when confronted with the scale and complexity of modern datasets, particularly unstructured data such as free text,

images, and videos.

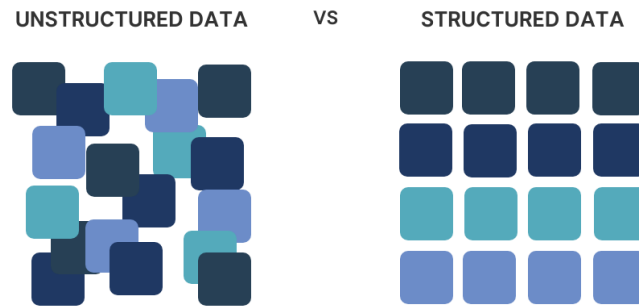


Figure 1.3: *Comparison of Structured and Unstructured Data*

As shown in Figure 1.3, structured data typically follows a predefined schema (e.g., tables or spreadsheets), while unstructured data can take many forms, such as text, images, or audio files, making it more challenging to process and analyze.

To address these challenges, organizations are adopting AI-powered analytics solutions. These systems leverage ML algorithms to automate much of the data processing workflow, allowing businesses to extract real-time insights and respond more quickly to changing conditions. The integration of NLP techniques further enhances the ability to analyze unstructured text data, such as customer feedback or social media posts, by identifying patterns in language, detecting sentiment, and categorizing themes automatically. This allows organizations to not only understand customer sentiment more accurately but also to anticipate emerging trends and adjust their strategies proactively [13].



Figure 1.4: *Example of Interactive Data Visualization Dashboard*

An example of an interactive data visualization dashboard is shown in Figure 1.4. Such dashboards allow non-technical users to explore data and gain insights through intuitive interfaces, enabling better decision-making based on complex data analytics.

Although AI and ML provide significant advantages, substantial challenges remain in developing systems that can manage and analyze unstructured data in real time. These challenges include the need for scalable infrastructure, real-time processing capabilities, and the ability to extract meaningful information from diverse data formats. This research proposes a comprehensive AI-powered analytics

pipeline that integrates ML and NLP techniques to address these challenges and generate actionable insights that can be used to inform business strategy, improve customer service, and enhance operational efficiency [47].

1.2 Main Goals

The primary objective of this study is to explore how advanced MLML techniques can be applied to large datasets to extract actionable insights, transforming raw data into a strategic asset. Specifically, the research aims to develop an **AI-powered analytics pipeline** capable of generating interactive insights that support informed decision-making across various industries. The core goals of this research are outlined as follows:

1. **Data Ingestion and Preparation:** Automate the collection, cleaning, and integration of data from multiple sources to ensure high-quality, consistent datasets [2]. The system is designed to function independently of specific datasets, thus allowing flexibility in processing diverse data inputs across structured and unstructured formats.
2. **Machine Learning Model Development:** Designing robust ML models capable of identifying complex patterns, making accurate predictions, and generating valuable insights from both structured data (e.g., transactional data) and unstructured data (e.g., text-based feedback).
3. **Interactive Visualization and Insights Generation:** Create dynamic visualization tools that allow users to explore and interpret insights in real time, making complex analytical results accessible and actionable[54].
4. **NLP Integration:** Apply advanced NLP techniques to extract relevant topics, trends, and sentiment from large volumes of unstructured text data, enabling deeper insights into qualitative information [23].
5. **User-Friendly Interface:** Design an intuitive interface that enhances accessibility to the analytics pipeline, making complex data insights understandable and actionable for users with varying levels of technical expertise.
6. **Continuous Learning and Improvement:** Build mechanisms for continuous learning, enabling the system to adapt and improve by learning from new data over time, thereby enhancing model accuracy and relevance.
7. **Deployment and Support:** Deploy the system in real-world environments and establish on-going support to address operational challenges, ensuring smooth functionality and reliability in production settings.
8. **Realizing Business Impact:** Demonstrate the practical applications of insights generated by the pipeline, including enhanced decision-making capabilities, improved customer service, and effective risk management [47].

1.3 Contributions

This research contributes substantially to the fields of Data Science and Business Analytics by introducing:

- **Comprehensive AI-Powered Pipeline:** An end-to-end system that automates data ingestion, cleaning, and analysis, enabling faster, more efficient decision-making.
- **Advanced NLP Integration:** State-of-the-art NLP techniques for extracting actionable insights from unstructured text data, which empower organizations to understand qualitative trends and feedback.
- **Real-World Validation:** Validation of the AI-powered pipeline through practical, real-world case studies, demonstrating its effectiveness in anomaly detection, customer service optimization, and more.
- **Scalability and Adaptability:** A scalable architecture that accommodates growing datasets and adapts to evolving business requirements, ensuring long-term usability and relevance.
- **Enhanced Decision-Making Tools:** Advanced visualization tools for exploring trends and patterns, fostering data-driven strategies and more informed decision-making.

1.4 Document Organization

This thesis is divided into six chapters, each addressing a crucial aspect of the research.

In Chapter 2, a review of the fundamental concepts related to Data Science, Big Data, ML, NLP is provided. Additionally, this chapter explores the challenges involved in processing large datasets and presents a survey of existing methodologies for data analysis, preprocessing, and topic modeling. A comprehensive review of related work is also included, examining the latest advancements and identifying gaps that this research aims to fill.

Chapter 3 focuses on the data used throughout this research, providing a detailed description of the data sources, structures, and characteristics. This chapter analyzes three primary datasets used in the study, highlighting their key attributes and explaining how these attributes align with the ML models applied in the analysis. Specific attention is given to both structured and unstructured data types, and their relevance to the objectives of this research. The chapter provides an in-depth examination of each dataset, covering topics such as data preprocessing, handling missing values, and preparing the data for model training and evaluation.

In Chapter 4, the technical implementation of the AI-powered analytics pipeline is described in depth. This chapter outlines how the datasets from Chapter 3 are ingested, processed, and used at various stages of the pipeline. The chapter also covers the tools and frameworks employed, such as Apache Spark and Databricks, and provides detailed explanations of custom functions and methodologies developed for data processing, outlier detection, vectorization, clustering, and visualization [15].

Chapter 5 presents the testing and validation methodologies used to evaluate the performance of the AI-powered analytics pipeline with real-world datasets. This chapter assesses the effectiveness of the ML models in extracting meaningful insights, identifying patterns, and supporting decision-making. Special emphasis is placed on the clustering and topic modeling results, as well as the challenges encountered during testing, particularly the computational limitations and stability issues with the Databricks platform when handling large-scale data.

Finally, Chapter 6 concludes the research by summarizing the key findings and reflecting on the contributions made. This chapter also outlines potential improvements and suggests future research directions,

particularly focusing on enhancing the scalability, adaptability, and accuracy of the AI-powered pipeline. The chapter also highlights the practical applications of the pipeline, including its deployment in a web-based interface for delivering interactive insights to end-users, further demonstrating the system's impact and potential for real-world usage.

Chapter 2

Background and Related Work

In today's data-driven world, organizations rely heavily on advanced analytical techniques to transform vast amounts of raw data into actionable insights that support informed decision-making [52]. However, the process of converting unprocessed data into meaningful information is both challenging and multifaceted. For instance, industries like healthcare and finance generate large-scale data streams that require real-time analysis to uncover trends, detect anomalies, and enable quick decisions. Traditional analytical methods often struggle with scalability, accuracy, and the ability to process data in real-time. These challenges underscore the growing importance of modern AI-powered analytics pipelines, which integrate advanced Machine Learning (ML) and Natural Language Processing (NLP) techniques to address these issues.

This chapter provides a comprehensive exploration of the core components that form these pipelines, beginning with the concept of **insights**—the valuable, data-derived knowledge that empowers businesses to make strategic decisions. We will delve into essential data pre-processing techniques, such as **handling missing values**, **outlier detection**, and **feature scaling**, which are critical for ensuring data quality and preparing datasets for further analysis [9].

Next, we explore ML algorithms—both supervised and unsupervised—which form the backbone of modern analytics [71]. These include clustering algorithms such as **K-Means** and **DBSCAN**, which help identify hidden patterns within unlabeled datasets [55].

With the increasing significance of NLP, this chapter examines how NLP techniques process unstructured text data, focusing on models like **BERTopic**, which generate coherent topics from large volumes of text [23]. Such models enable businesses to derive insights from unstructured data sources, such as customer feedback or social media content.

Data visualization is another key element, making complex analyses accessible and actionable. Tools like **Matplotlib** and **Seaborn** allow businesses to present complex datasets and analytical results in visually intuitive formats [54]. These tools help uncover trends, patterns, and insights that facilitate real-time decision-making.

Moreover, the chapter addresses the challenges posed by **Big Data**, including the "5 Vs"—**Volume**, **Velocity**, **Variety**, **Veracity**, and **Value**—and how these factors influence the need for sophisticated data management and analytics systems [70]. Understanding these dynamics is crucial for developing scalable, robust, and real-time analytics systems.

Finally, we discuss advanced **clustering methods** and their role in discovering natural patterns in

large datasets, even when predefined labels are absent. By combining these components, organizations can build powerful analytics systems that not only handle the complexities of Big Data but also deliver timely and insightful results [52].

2.1 Data Science and Big Data

Data **Science** [75] is an interdisciplinary field that combines scientific methods, algorithms, and systems to extract insights and actionable knowledge from both structured and unstructured data. It integrates principles from statistics, computer science, and domain-specific expertise to address complex data challenges. For example, in retail, Data Science enables the analysis of customer purchasing behaviors, while in healthcare, it assists in predicting patient outcomes based on historical data. This transformative field has become indispensable for organizations across sectors, empowering them to make informed, data-driven decisions. By leveraging advanced technologies, Data Science facilitates not only predictive analysis but also prescriptive insights, allowing businesses to optimize operations and enhance customer satisfaction.

With advancements in computational power and data availability, Data Science has evolved into a critical tool for industries looking to leverage their data for competitive advantage. The field combines traditional statistical analysis with cutting-edge ML and AI techniques to uncover patterns, make predictions, and inform decision-making.

At the heart of Data Science are several key processes that transform raw data into valuable insights:

- **Data Collection:** Gathering data from multiple sources, such as transactional databases, IoT devices, and social media platforms. Effective data collection ensures that data from various formats (structured, semi-structured, and unstructured) can be processed together for deeper analysis.
- **Data Processing and Cleaning :** Preparing data by handling missing values, identifying outliers, and transforming variables to ensure accuracy and reliability [9]. Data cleaning is vital for ensuring the integrity of any analysis, especially when large datasets are involved.
- **Exploratory Data Analysis (EDA):** Before any advanced analysis, data scientists conduct EDA to understand the dataset's underlying structure. Through statistical summaries and visualizations, EDA helps reveal trends, outliers, and relationships within the data, guiding the next steps in analysis and model development.
- **Predictive and Prescriptive Analytics:** Moving beyond exploration, Data Science employs predictive analytics to forecast future outcomes based on historical data, often utilizing ML algorithms. Predictive analytics focuses on identifying what is likely to happen in the future, providing insights into potential trends and events. In contrast, prescriptive analytics goes a step further by not only predicting outcomes but also providing actionable recommendations [60]. It helps organizations optimize decisions and determine the best course of action by suggesting specific actions to take in order to achieve desired results, considering various scenarios and constraints.

- **Data Visualization:** Presenting complex analytical results through clear, intuitive visuals (e.g., charts, graphs, and interactive dashboards) is crucial. Data visualization enhances understanding, allowing stakeholders to quickly grasp key insights and make informed decisions [56].

In industries like e-commerce, healthcare, and manufacturing, Data Science helps organizations optimize operations, predict customer behavior, and drive innovation, making it an indispensable tool in today's competitive landscape.

2.1.1 Introduction to Big Data

Big Data [70] refers to datasets that are so large, complex, or fast-moving that they exceed the capabilities of traditional data processing tools. These datasets are characterized by the **5 Vs**—**Volume**, **Velocity**, **Variety**, **Veracity**, and **Value**—each representing a critical aspect of handling and processing large-scale data.

Big Data is often characterized by the **5 Vs**:

- **Volume:** The sheer amount of data generated globally, often measured in terabytes or petabytes. For example, social media platforms, IoT devices, and transactional systems create vast amounts of data that require advanced storage and processing capabilities.
- **Velocity:** The speed at which data is produced and must be processed. In industries such as finance or healthcare, real-time data processing is essential for timely decision-making and responding to dynamic environments.
- **Variety:** Big Data comes in many forms, from structured data (like tables) to unstructured data (such as text, images, and videos). Managing this diverse range of formats presents a significant challenge for data integration and analysis.
- **Veracity:** The reliability and quality of data. In the context of Big Data, it can be difficult to ensure that data is accurate and free from inconsistencies, as noise and incomplete records can distort analysis.
- **Value:** Ultimately, the goal of Big Data is to extract meaningful insights that drive business decisions. If data does not provide value, the resources required to manage it are not justified.

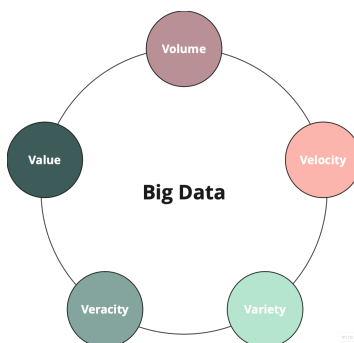


Figure 2.1: *The 5Vs of Big Data*

Big Data’s ability to deliver value comes from its potential to reveal hidden patterns, trends, and correlations in datasets too large or complex for traditional tools. Organizations that can harness Big Data effectively are more likely to gain a competitive edge, optimizing their operations, understanding customer behaviors better, and responding faster to market changes.

2.1.2 Challenges of Big Data

Handling Big Data introduces unique challenges, primarily due to the scale, complexity, and speed at which data must be processed:

- **Data Integration:** Big Data often originates from multiple sources and exists in a variety of formats—structured, semi-structured, and unstructured. Merging these disparate data types into a unified and consistent format is a complex process requiring advanced ETL (Extract, Transform, Load) workflows. The diversity of data formats makes it difficult to create cohesive datasets.
- **Scalability:** Traditional data management systems are not equipped to handle the sheer size of Big Data. Technologies such as **Apache Hadoop** and **Apache Spark** have become essential for distributed processing, enabling organizations to scale their data processing across multiple servers and efficiently manage increasing data volumes.
- **Data Quality:** Ensuring that Big Data is clean, accurate, and complete is vital for generating reliable insights. Poor data quality can lead to flawed analyses, which in turn results in misguided decisions. Robust cleaning, validation, and quality control methods are necessary to maintain data integrity.
- **Privacy and Security:** With the exponential growth of data, particularly in sensitive areas such as personal information, the risk of data breaches increases. Ensuring compliance with regulations like the **General Data Protection Regulation (GDPR)** and implementing strong security measures are critical for protecting data and maintaining trust.

Effectively addressing these challenges is crucial for organizations aiming to leverage Big Data for strategic advantages. By overcoming data integration issues, ensuring scalability, and maintaining high standards of data quality, businesses can turn vast amounts of raw data into actionable insights [60]. For instance, overcoming data integration challenges can lead to more cohesive and comprehensive analyses, while ensuring scalability enables businesses to handle growing data volumes without performance degradation. Maintaining data quality is essential to avoid flawed analyses and misguided decisions.

2.2 Machine Learning

Machine Learning [14] is a fundamental area of AI that enables systems to learn from data and improve their performance without explicit programming. Unlike traditional algorithms that follow predefined rules, ML models identify patterns and relationships within data, allowing them to make predictions, classifications, or decisions. This ability to adapt and improve over time makes ML a powerful tool for solving complex, data-driven challenges in various industries [57].

With the exponential growth of data sources such as transaction logs, sensor data, and unstructured text, ML has become indispensable in deriving insights, making predictions, and enhancing decision-making processes. Its versatility extends across industries like healthcare, finance, marketing, and manufacturing, where its ability to scale with data and continuously improve is invaluable [34].

ML [71] can be divided into three main categories, based on how models learn from the data they are given:

- **Supervised learning:** Supervised learning involves training a model on labeled data, where each input has a corresponding output. The model learns to predict the output for new, unseen data by finding patterns in the labeled examples. In email spam detection, a model is trained on a dataset where emails are labeled as either "spam" or "not spam." The model learns to identify key features, such as certain keywords or sender information, that are common in spam emails. Once trained, the model can predict whether new incoming emails are likely to be spam or legitimate, helping users keep their inboxes clean.
- **Unsupervised learning:** Unsupervised learning deals with data that is not labeled. The model's goal is to find hidden structures or patterns in the data without specific outcomes being given [69]. In retail, unsupervised learning can be used to analyze customer purchasing behaviors. A clustering algorithm might group customers based on their buying habits, identifying different segments such as frequent shoppers, occasional buyers, or discount hunters. These insights allow businesses to personalize marketing efforts without predefined categories.
- **Reinforcement Learning:** Reinforcement learning (RL) involves an agent that learns to make decisions by interacting with its environment, receiving feedback in the form of rewards or penalties based on its actions. In video games, reinforcement learning can be used to train AI opponents that learn how to beat human players. The AI makes moves and receives feedback—such as winning or losing a round—and over time, it learns the strategies that are most effective for winning the game.
- **Semi-Supervised Learning:** Semi-supervised learning combines a small amount of labeled data with a large amount of unlabeled data. This approach is especially useful when obtaining labeled data is expensive or time-consuming. In facial recognition, labeling every face in a large dataset can be impractical. A semi-supervised model might be trained with a small number of labeled face images and a larger set of unlabeled ones. The labeled data helps the model learn basic patterns, while the unlabeled data allows it to generalize and improve recognition accuracy on a broader set of faces.

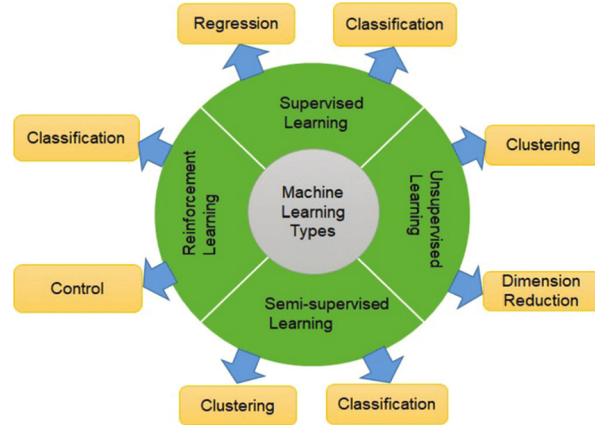


Figure 2.2: Overview of Machine Learning Categories: Supervised, Unsupervised, Reinforcement, and Semi-Supervised Learning

2.2.1 Supervised Learning

Supervised learning is the most commonly used type of ML. In this approach, the model is provided with a labeled dataset where each input (feature) is associated with a specific output (label). The goal of the model is to learn the mapping between inputs and outputs so that when new, unseen data is introduced, it can accurately predict the correct output.

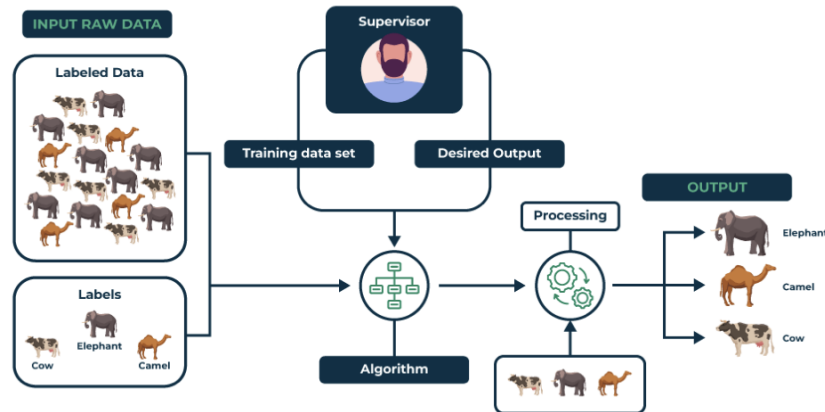


Figure 2.3: Schematic Representation of Supervised Learning Process: Labeled Data Leading to Predictions

In practice, supervised learning involves the model comparing its predictions to the actual outputs (labels) and adjusting its internal parameters to minimize errors. Over time, the model becomes more accurate at generalizing to new data. This type of learning is particularly useful for tasks like classification (e.g., determining if an email is spam or not) and regression (e.g., predicting house prices based on features like size and location).

Key Characteristics of Supervised Learning:

- **Labeled Data:** The model is trained on a dataset where the correct output is provided for each input. This allows the model to learn from the examples, adjusting its predictions based on the labeled outcomes.

- **Training Process:** Supervised models learn iteratively using optimization techniques like gradient descent, which minimizes the error between predicted and actual outcomes. Common loss functions include MSE for regression tasks and cross-entropy for classification.
- **Predictive Modeling:** Once trained, the model can make predictions on new, unseen data. Supervised learning is ideal for tasks that require forecasting or categorizing information based on historical examples.

2.2.2 Unsupervised Learning

Unsupervised learning involves training a model on a dataset without labeled outputs. The model's task is to explore the data and identify hidden patterns, groupings, or structures without any prior guidance on what the “correct” answers are. This is particularly valuable for exploratory analysis where the goal is to find natural patterns in data, such as clustering or anomaly detection.

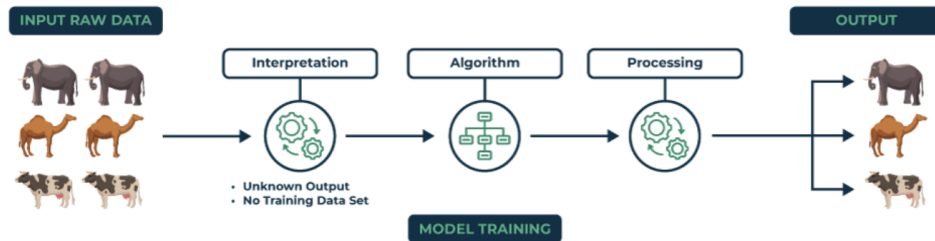


Figure 2.4: *Clustering Process in Unsupervised Learning: Grouping Data Points Based on Similarities*

In this type of learning, the model organizes data by finding similarities between data points. Over time, it becomes better at identifying clusters or detecting outliers that deviate from the norm. This makes unsupervised learning ideal for analyzing data where labeling is impractical or too expensive.

Key Characteristics of Unsupervised Learning:

- **No Labeled Data:** The model works entirely with unlabeled data, discovering patterns and relationships without pre-existing labels.
- **Exploratory Nature:** Often used for exploring data to find hidden patterns, making it useful for clustering or dimensionality reduction tasks.
- **Pattern Discovery:** Focuses on identifying groups of similar data points or uncovering anomalies that differ from the rest of the dataset.

2.2.3 Reinforcement Learning

Reinforcement learning involves training an agent to make a sequence of decisions by interacting with an environment. The agent performs actions, receives feedback in the form of rewards or penalties, and adjusts its strategy to maximize cumulative rewards over time. This approach is inspired by behavioral psychology and is widely used in areas such as robotics, game playing, and autonomous vehicles.

Key Characteristics of Reinforcement Learning:

- **Trial and Error:** The agent learns by interacting with the environment, exploring different actions, and observing the outcomes.

- **Feedback Mechanism:** Rewards or penalties guide the learning process, reinforcing behaviors that lead to positive outcomes.
- **Delayed Rewards:** The agent must learn to make decisions that may not yield immediate rewards but are beneficial in the long run.

2.2.4 Semi-Supervised Learning

Semi-supervised learning bridges the gap between supervised and unsupervised learning. It uses a small amount of labeled data and a large amount of unlabeled data. This is particularly useful in situations where obtaining labeled data is costly or time-consuming, but there is an abundance of unlabeled data available.

By combining both labeled and unlabeled data, semi-supervised learning allows models to learn more effectively, leveraging the labeled data to guide the learning process while still extracting useful information from the larger set of unlabeled examples.

Key Characteristics of Semi-Supervised Learning:

- **Limited Labeled Data:** The model is trained on a small labeled dataset combined with a much larger set of unlabeled data.
- **Efficient Use of Data:** This approach allows models to benefit from both labeled and unlabeled data, improving accuracy without needing extensive labeled datasets.

2.3 Definition of Insight

An **insight** [3] represents more than just observation; it is the product of deep analysis that transforms raw data into actionable knowledge. Insights reveal hidden patterns, correlations, or anomalies that provide clarity and direction, allowing organizations to make informed, strategic decisions. The ability to derive valuable insights is crucial for businesses across all industries, enabling them to leverage data as a competitive asset.

Insights empower decision-makers to go beyond surface-level information, uncovering relationships within data that can significantly impact performance and efficiency. For instance, in marketing, an insight might reveal that a particular customer segment is more likely to engage with personalized offers, leading to more effective campaign strategies. Similarly, in healthcare, insights from patient data can inform predictive models that anticipate health risks and suggest preventative measures. In the banking sector, insights derived from transaction data can help identify customers who are at a higher risk of defaulting on loans, enabling the bank to proactively adjust credit terms or offer financial counseling, thereby reducing the likelihood of default and improving customer retention.

2.3.1 Role of Insights in Modern Data-Driven Environments

In a data-rich world, **insights** act as the bridge between raw data and informed action. While data itself holds immense potential, it is through insights that this potential is realized. Insights allow organizations to:

- **Identify Trends:** By analyzing historical and real-time data, businesses can detect emerging trends that can inform product development, market strategies, or operational changes.
- **Understand Customer Behavior:** In industries like retail and marketing, insights derived from customer data (e.g., purchasing patterns, engagement metrics) help tailor services and offerings to meet consumer demands.
- **Optimize Operational Efficiency:** Insights enable companies to streamline processes by revealing inefficiencies or bottlenecks within operations, leading to better resource allocation and cost reductions.

For example, insights in the financial sector could involve the identification of patterns in transactional data that help in fraud detection, or in customer retention strategies where understanding user behavior enables better service delivery .

2.3.2 How Insights are Derived

The process of generating insights is systematic and involves several key steps, which ensure that the raw data is transformed into actionable conclusions:

1. **Data Acquisition:** The first stage involves gathering data from various sources, including structured (e.g., databases, spreadsheets) and unstructured (e.g., text, images, social media) formats. The diversity and quality of data sources are crucial in shaping the depth of insights generated .
2. **Data Processing:** Before analysis can occur, the data must be cleaned, organized, and prepared. This step includes handling missing values, normalizing data, and applying preprocessing techniques to ensure consistency and accuracy. Data preparation is fundamental in preventing errors during analysis that could distort insights.
3. **Contextual Understanding:** Analyzing data in its specific business context is essential. For example, in a retail environment, understanding customer behavior patterns through data analysis might provide insights into purchasing trends. Applying domain knowledge allows data analysts to identify meaningful patterns that directly address business challenges.
4. **Insight Generation:** The final step is deriving insights by interpreting the analyzed data. This involves using advanced tools like ML algorithms or predictive models to uncover deeper correlations or patterns that might not be immediately apparent. Insights at this stage provide actionable recommendations that can drive business strategies and decisions.

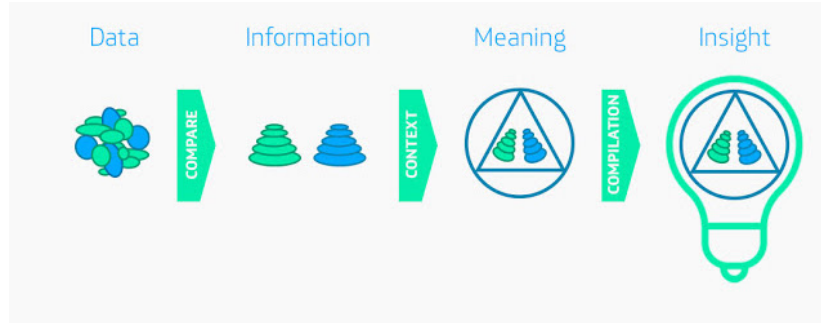


Figure 2.5: *Process of transforming data into insights*

2.3.3 Example of Insight Generation

To better understand the power of insights, consider a dataset from a bank containing customer transaction histories. A basic analysis might categorize customers by their spending habits, providing only a general overview of behavior. However, a deeper insight could reveal a correlation between customers’ credit card usage patterns and their likelihood of defaulting on a loan, showing that customers who frequently miss credit card payments are more likely to struggle with loan repayments.

This insight allows the bank to implement targeted risk mitigation strategies, such as offering financial counseling or adjusting credit terms for customers identified as being at higher risk of default. Figure 2.6 [19] illustrates a related insight by analyzing the correlation between educational level and the likelihood of loan default. The figure shows that customers with lower educational levels tend to have higher rates of default, highlighting how demographic factors can inform financial risk assessments.

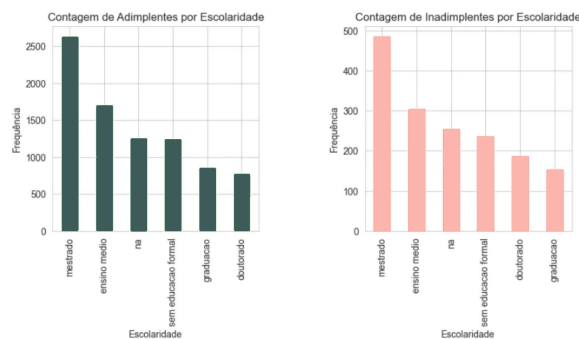


Figure 2.6: *Correlation between educational level and loan default rates*

The process of transforming banking data into actionable insights goes beyond basic data analysis. By understanding the context of customer behavior and applying advanced techniques such as ML or NLP, financial institutions can derive insights that inform strategic decisions. The ability to generate meaningful insights from data is a crucial asset in today’s competitive financial landscape, enabling banks to not only understand their current risk exposure but also predict future customer behavior and take proactive actions to minimize defaults and improve financial stability.

2.4 Data Pre-processing Techniques

Data pre-processing is an essential step in the ML pipeline, ensuring that raw data is cleaned and structured for optimal model performance. By handling missing values, detecting outliers, scaling features, and selecting relevant variables, pre-processing prepares the data for analysis and helps improve model accuracy and reliability. This section covers the key techniques used in data pre-processing, from handling null values to feature engineering, ensuring the data is ready for ML.

2.4.1 Handling Missing Data

Dealing with missing data is a common challenge in real-world datasets, and how it is handled can significantly impact the performance of ML models. Properly managing missing data ensures that analyses and models remain accurate and unbiased. The choice of method depends on the type and quantity of missing data, as well as the specific characteristics of the dataset. Below are several widely-used techniques for handling missing data [27].

2.4.1.1 Removing Missing Values

One of the simplest ways to handle missing data is by removing rows or columns that contain missing values. This approach, known as **listwise deletion**, is particularly useful when the amount of missing data is small and randomly distributed, as it may not significantly affect the dataset's overall representativeness.

In a survey dataset, if only 1-2% of the responses are missing, you might choose to remove those rows entirely. However, if missing data is concentrated in a key column (e.g., a critical feature like age or income), this method could lead to biased analysis if the missing values aren't random.

The advantages include ease of implementation and the assurance of clean data without missing entries. However, the disadvantages are that if many rows or columns have missing values, this method can lead to significant data loss. Additionally, it can introduce bias if the missing data is not randomly distributed.

2.4.1.2 Mean, Median, and Mode Imputation

Imputation replaces missing values with estimates derived from the available data [61]. Three common imputation methods are:

- **Mean Imputation:** Replaces missing values with the mean of the column. This method is best suited for continuous, normally distributed data, but it can be skewed by outliers.

$$\text{Mean} = \frac{\sum_{i=1}^n x_i}{n} \quad (2.1)$$

- **Median Imputation:** Uses the median (the middle value) of the data, making it more robust for datasets with outliers or skewed distributions.

$$\text{Median} = \begin{cases} x_{(\frac{n+1}{2})} & \text{if } n \text{ is odd} \\ \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2} & \text{if } n \text{ is even} \end{cases} \quad (2.2)$$

- **Mode Imputation:** For categorical data, missing values are replaced with the mode, the most frequent category. While effective, this method can distort the dataset if the missing data is non-random.

Suppose you're analyzing a dataset on customer demographics and purchasing behavior for a retail company. The dataset includes customer age (continuous data), income (continuous data), and preferred shopping day (categorical data).

- For **age**: Some customers did not provide their age. You could apply **mean imputation**, filling in missing ages with the average age of all customers. However, if the age distribution is skewed (e.g., more older customers), you might prefer **median imputation** to better represent the central tendency without being affected by outliers.
- For **income**: If there are missing values in the income field, and the data is highly skewed (e.g., a few very high-income individuals), **median imputation** would provide a more accurate estimate than the mean.
- For **preferred shopping day**: If some customers did not indicate their favorite day to shop, you could apply **mode imputation**, filling in the missing values with the most frequently selected day (e.g., "Saturday").

This approach works across different data types, ensuring that each missing value is imputed in a way that best represents the overall dataset.

2.4.1.3 Forward/Backward Fill

For time-series data, maintaining the sequence of data points is crucial. **Forward Fill (FFILL)** and **Backward fill (BFILL)** methods [44] propagate the last observed value forward or the next observed value backward, respectively.

- **FFILL:** Replaces missing values with the previous value in the sequence.
- **BFILL:** Fills missing values with the next available value.

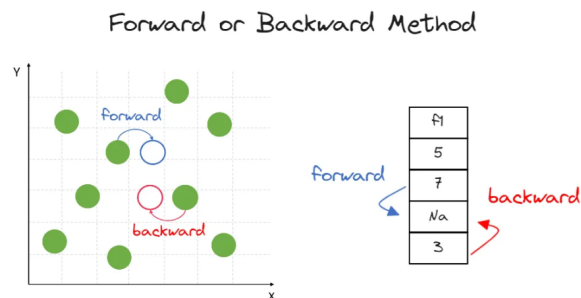


Figure 2.7: *Forward and Backward Fill Imputation*

While these techniques are useful in maintaining the sequence of data points, they can lead to misleading results if there are large gaps between observed data points or if the missing data follows non-random

patterns. FFILL and BFILL are best applied in scenarios where the values change gradually over time, as abrupt changes or trends might not be accurately captured.

For example, in stock market data, if a stock price is missing on a particular day, you might use the previous day's price (forward fill) to fill the gap. Similarly, if weather data for one hour is missing, you could fill it using the next recorded value (backward fill).

2.4.1.4 K-Nearest Neighbors (KNN) Imputation

KNN [17] imputation is a more advanced method where missing values are filled in by averaging the values of the k nearest neighbors (based on a distance metric, such as Euclidean distance) to impute missing values. It works by finding the nearest data points (in feature space) and using their values to fill in the gaps. KNN imputation works well in situations where there is a strong relationship between variables, allowing it to estimate missing values based on the similarity between data points.

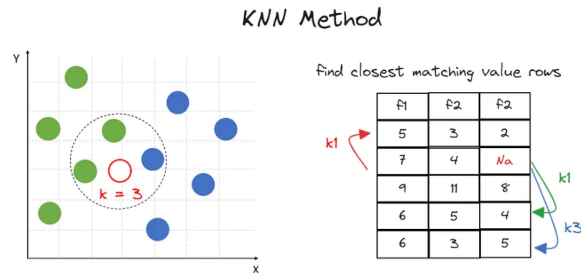


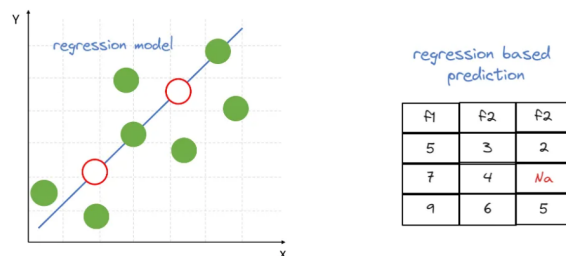
Figure 2.8: *KNN Imputation*

In a customer dataset with features like age, income, and spending score, if a customer's income is missing, you can predict it based on the income of the k most similar customers, taking into account their age and spending habits. This method is computationally expensive, but it often provides more accurate results than simple imputation methods.

This method has several advantages over simpler imputation techniques. By considering the relationships between data points, KNN provides more accurate estimations, especially in complex datasets. However, KNN imputation is computationally intensive and sensitive to the number of neighbors k and the choice of distance metric. Additionally, feature scaling is important in KNN imputation, as the method is sensitive to the magnitude of the features.

2.4.1.5 Regression Imputation

In **regression imputation**, missing values are predicted using a regression model, treating the missing data as the dependent variable and other available features as independent variables. The model is trained using the non-missing data to predict the missing values [29].

Figure 2.9: *Regression Imputation*

In a dataset of housing prices, if some values for "square footage" are missing, you can build a regression model using features like the number of bedrooms, lot size, and location to predict the missing values. This method is useful when there is a strong correlation between features.

This method is advantageous because it preserves the relationships between variables, providing more accurate imputations than simpler methods. However, it assumes that a linear or non-linear relationship exists between the features, which may not always be the case. Additionally, regression imputation requires careful validation of the model to ensure that the predictions are reliable.

2.4.1.6 Multiple Imputation

Multiple imputation accounts for uncertainty in missing data by creating several imputed datasets. Each dataset is analyzed independently, and the results are pooled to produce a final result [48]. The process involves:

1. Generating multiple datasets with randomly imputed values.
2. Analyzing each dataset separately.
3. Combining the results to reflect the uncertainty in the missing values.

An example of using multiple imputation is in a clinical trial dataset, where patient responses to a treatment are partially missing. In this case, multiple imputation can be used to generate several plausible imputed datasets. Each dataset is analyzed separately to reflect the uncertainty in the imputed values, which is particularly useful when the missing data is non-random and requires a more sophisticated approach.

Multiple imputation reduces bias and provides a more accurate reflection of the uncertainty associated with missing data. However, it is computationally intensive and more complex to implement, making it less suitable for very large datasets or real-time applications.

2.4.1.7 Expectation-Maximization (EM)

The Expectation-Maximization (EM) [45] algorithm is a sophisticated technique that estimates missing values by iteratively maximizing the likelihood of the observed data. EM alternates between two steps: the **Expectation Step (E-step)**, where the missing values are estimated based on the current parameters, and the **Maximization Step (M-step)**, where the parameters are updated to maximize the likelihood of the data. This process continues until the model converges on a solution.

A common application of the Expectation-Maximization (EM) algorithm is in marketing datasets where some customer purchase records are incomplete. The EM method helps estimate missing purchase amounts by modeling the underlying data distribution. This technique is particularly effective when the missing data aligns with a known distribution, like the normal or Poisson distribution.

EM is particularly well-suited for structured datasets where the underlying distribution is known or can be approximated. It provides statistically robust estimates but can be computationally expensive and sensitive to the initial parameter choices. Additionally, EM may converge to a local minimum, leading to suboptimal results in certain cases.

2.4.1.8 Autoencoders

Autoencoders, a type of artificial neural network, are powerful tools for handling missing data, especially in high-dimensional and complex datasets. An autoencoder learns to compress data into a lower-dimensional representation and then reconstructs the data from this compressed form. In the context of missing data imputation, autoencoders can learn the underlying structure of the data and use this knowledge to reconstruct the missing values [28].

This method is particularly effective for datasets with non-linear relationships between features, where traditional methods like mean or regression imputation may fail. However, autoencoders require substantial computational resources and large amounts of training data to perform effectively. The design and training of the neural network are also critical to its success.

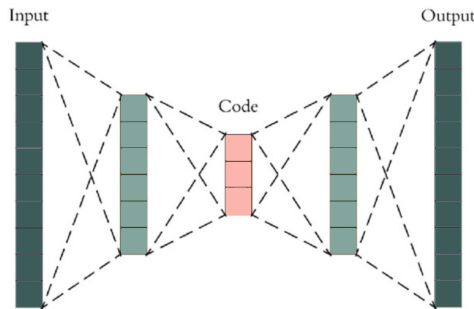


Figure 2.10: *Autoencoders Imputation*

In image data (e.g., partially corrupted photos), an autoencoder can reconstruct missing pixels by learning the underlying structure of the image. This method is powerful in tasks like image reconstruction or handling missing features in very large datasets with non-linear relationships.

2.4.2 Outlier Detection

Outliers are data points that significantly deviate from the majority of observations within a dataset. These unusual data points may arise due to measurement errors, anomalies, or rare events that require further investigation. Detecting and managing outliers is crucial for ensuring data quality, as outliers can skew the results of ML models and lead to unreliable predictions [58]. Proper handling of outliers helps improve model accuracy and robustness, particularly in sensitive fields such as finance, healthcare, and marketing, where outliers may indicate fraud, system failures, or unique opportunities.

Outliers can be categorized into several types based on the context and nature of the data:

- **Univariate outliers:** Outliers present in only one feature or variable.
- **Multivariate outliers:** Outliers involving abnormal relationships between two or more variables.
- **Artificial outliers:** Caused by data collection, input, or processing errors.
- **Natural outliers:** Valid, rare events or anomalies that deviate from normal patterns and can provide valuable insights.

Below, we explore various techniques for detecting outliers, each suited to different types of datasets and distributions.

2.4.2.1 Z-Score

The Z-Score method is a commonly used statistical approach for identifying outliers in datasets that follow a normal distribution. It represents how many standard deviations a particular data point deviates from the mean. Typically, a Z-Score greater than 3 or less than -3 is considered indicative of an outlier.

$$Z = \frac{x - \mu}{\sigma} \quad (2.3)$$

Where x is the data point, μ is the mean of the dataset, and σ is the standard deviation.

The Z-Score method is easy to implement and widely used due to its simplicity and straightforward interpretation. However, it assumes the data follows a normal distribution, making it less effective for datasets with skewed distributions or extreme outliers that can distort the mean and standard deviation [74]. Additionally, it may not perform well in datasets where outliers do not follow a normal pattern.

2.4.2.2 Interquartile Range (IQR)

The IQR method is a robust statistical technique that is less affected by extreme values, making it suitable for non-normal distributions. IQR is calculated as the difference between the first quartile (Q1) and the third quartile (Q3) of the data. Outliers are identified as data points that fall below $Q1 - 1.5 \times IQR$ or above $Q3 + 1.5 \times IQR$.

$$IQR = Q3 - Q1 \quad (2.4)$$

Outliers are identified as:

$$\text{Lower bound} = Q1 - 1.5 \times IQR \quad (2.5)$$

$$\text{Upper bound} = Q3 + 1.5 \times IQR \quad (2.6)$$

This method is widely used because it is not influenced by extreme outliers and works well with skewed data. However, the threshold of 1.5 times the IQR may not be appropriate for all datasets and might require adjustment based on the specific characteristics of the data. While effective, it may overlook extreme outliers that fall outside these predefined boundaries in some cases [7].

2.4.2.3 Box Plot

A Box Plot is a graphical representation of data distribution, highlighting key statistical properties such as the median, quartiles, and potential outliers. Outliers in a Box Plot are identified as data points that lie beyond the "whiskers," which are typically set at 1.5 times the IQR.

Box plots are useful for visualizing the spread and skewness of data, making them a quick and intuitive way to detect outliers. However, this method is limited in multivariate datasets where relationships between variables are not apparent from the distribution of individual features. While it provides a useful visual tool for univariate data, it may fail to identify complex outliers in multivariate settings [6].

2.4.2.4 Histogram

A Histogram allows for the visual detection of outliers by displaying the frequency distribution of the data. Outliers can be identified as unusually low or high bars compared to the majority of the dataset. Histograms are effective for detecting univariate outliers and provide a simple way to visually inspect data. However, like box plots, histograms are limited in detecting multivariate outliers, as they do not account for the relationships between different variables. This method works best when analyzing a single feature at a time and may miss outliers that only emerge when multiple variables are considered together [24].

2.4.2.5 Median Absolute Deviation (MAD)

The Median Absolute Deviation (MAD) is a robust alternative to the Z-Score for detecting outliers in skewed distributions. Instead of using the mean and standard deviation, MAD calculates the median of the absolute deviations from the median of the dataset.

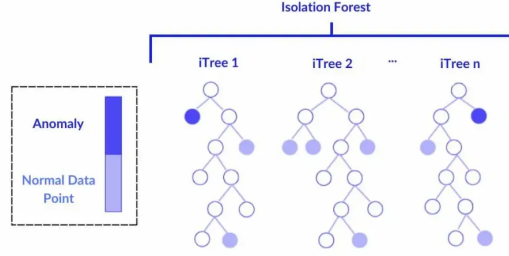
$$MAD = \text{median}(|X_i - \text{median}(X)|) \quad (2.7)$$

MAD works well in datasets with skewed distributions and is robust to extreme values, making it more reliable than the Z-Score in non-normally distributed data. However, MAD is more computationally complex and may be less familiar to users accustomed to simpler methods like Z-Score. Despite its effectiveness in handling outliers, it requires careful application in datasets with extreme variability [40].

2.4.2.6 Isolation Forest

Isolation Forest detects outliers by recursively partitioning the dataset. It builds a binary tree, and outliers are identified as points that are isolated with fewer splits than regular data points [26].

$$Isolation = \frac{\text{isolated points}}{\text{total points}} \quad (2.8)$$

Figure 2.11: *Isolation Forest method*

Isolation Forest is highly scalable and works well with large datasets, offering a robust approach to detecting anomalies. It is particularly useful for high-dimensional data. However, it requires careful tuning of parameters like the number of trees and maximum sample size to optimize performance, and the method can become computationally intensive when dealing with very large datasets. Furthermore, inappropriate parameter selection may lead to incorrect isolation of normal data points as outliers.

2.4.2.7 Local Outlier Factor (LOF)

The Local Outlier Factor (LOF) measures the local density deviation of a data point relative to its neighbors. Points with significantly lower densities than their neighbors are flagged as outliers.

$$LOF_k(p) = \frac{\sum (\frac{\text{local_reachability_density}(p)}{\text{local_reachability_density}(o)})}{k} \quad (2.9)$$

Where p is the point being evaluated, o represents neighboring points, and k is the number of neighbors. LOF is particularly effective for detecting local anomalies in densely clustered datasets where global methods (like Z-Score) might fail. It adapts well to complex and uneven data distributions. However, the method can be computationally expensive, especially for large datasets, and may require significant processing time to evaluate each data point against its neighbors. This sensitivity to parameters such as the number of neighbors also makes it harder to tune for optimal performance [37].

2.4.2.8 Mahalanobis Distance

Mahalanobis Distance detects multivariate outliers by measuring the distance between a point and the center of the data distribution, accounting for correlations between variables.

$$D_M(X) = \sqrt{(X - \mu)^T \sum^{-1} (X - \mu)} \quad (2.10)$$

Where X is the data point, μ is the mean vector, and \sum is the covariance matrix.

Mahalanobis Distance is particularly useful for multivariate data, especially when relationships between features are important. However, it requires an accurate estimation of the covariance matrix, which can be challenging with small or imbalanced datasets. Moreover, the method assumes that the underlying distribution is Gaussian, which might not always be the case [38].

2.4.2.9 Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

DBSCAN is a clustering algorithm that detects outliers as points that do not belong to any dense cluster [55]. It identifies outliers based on density and isolates them as noise points.

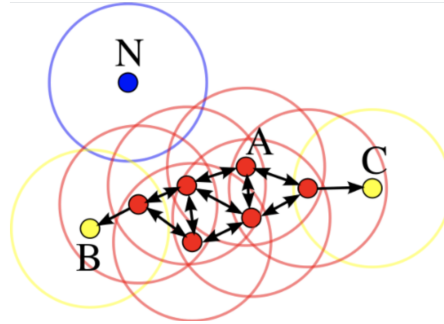


Figure 2.12: *Outlier detection using DBSCAN: points outside of dense clusters are considered outliers*

DBSCAN is effective in identifying both outliers and irregularly shaped clusters in datasets. It works particularly well for datasets with non-linear patterns, making it a versatile tool for anomaly detection. However, DBSCAN is sensitive to parameter selection (e.g., the distance threshold and minimum number of points required to form a cluster). Incorrect parameters may either miss outliers or overfit noise into clusters, making the method less effective if not properly tuned.

2.4.2.10 CleanLab

CleanLab[9] is a tool designed to detect label errors in datasets, especially in classification tasks where mislabeled data can act as hidden outliers. While CleanLab focuses on identifying mislabeled or noisy samples, rather than traditional outliers, it enhances the quality of supervised learning datasets by ensuring that labels are accurate, thus preventing misleading model performance. CleanLab uses algorithms that detect errors in the data labeling process and flags those samples that may have been mislabeled.

CleanLab can be highly beneficial for improving data integrity in scenarios where label quality is critical, such as healthcare or financial datasets, where accurate classification is essential. By identifying mislabeled data points, CleanLab helps mitigate the risk of poor model performance caused by noisy labels, ensuring that the ML model is trained on a clean and reliable dataset.

However, CleanLab does not directly detect statistical outliers in the sense of extreme or anomalous data points in feature space. Instead, it complements outlier detection by ensuring that classification labels are correct, reducing the noise that can obscure true outliers.

Advantages: Improves label quality in supervised learning datasets and enhances overall data integrity for better model performance.

Disadvantages: Focuses on detecting label errors rather than statistical outliers, making it more suitable for supervised learning tasks with labeled datasets rather than exploratory analysis in unsupervised settings.

2.4.3 Normalization and Standardization

Normalization and standardization are essential steps in data pre-processing that ensure features are on a similar scale, which is critical for many ML models. Although both techniques aim to scale features, they do so in different ways and are suited for different types of data and ML algorithms [46].

Normalization is the process of scaling data to fit within a specific range, typically $[0, 1]$. This is done by transforming the data so that the minimum value becomes 0 and the maximum value becomes 1. It is particularly useful for algorithms that rely on distance metrics, such as KNN and neural networks, where the relative scale of the features can significantly impact performance.

Standardization, on the other hand, transforms data so that it has a mean of 0 and a standard deviation of 1. This is achieved by subtracting the mean from each data point and dividing by the standard deviation. Standardization is especially useful when the data follows a normal distribution or when features have different units but the algorithm (e.g., Support Vector Machines, Logistic Regression) assumes normally distributed data.

2.4.3.1 Min-Max Normalization

Min-Max normalization scales the feature values to a specific range, often $[0, 1]$, using the following formula:

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (2.11)$$

This method preserves the relationships between the original values while ensuring all features are on the same scale. It's particularly useful for algorithms like neural networks, where the data should be bounded. However, it is sensitive to outliers, which can distort the scaling process if extreme values exist in the dataset [41].

2.4.3.2 Max-Abs Normalization

Max-Abs scaling transforms data by dividing each feature by its maximum absolute value:

$$X_{norm} = \frac{X}{X_{max}} \quad (2.12)$$

This method ensures that the values are between -1 and 1. Max-Abs normalization is effective when dealing with data that includes both positive and negative values. However, like Min-Max normalization, it is sensitive to outliers since a single extreme value can affect the scaling [39].

2.4.3.3 Robust Scaler

The Robust Scaler [63] is a variation of standardization that uses the median and IQR instead of the mean and standard deviation:

$$X_{robust} = \frac{X - Median}{IQR} \quad (2.13)$$

This method is robust to outliers since it focuses on the central data points and is unaffected by extreme values. It is especially useful when the data contains a significant number of outliers or has a skewed distribution. However, while it reduces the influence of outliers, it may not perform well on datasets that require precise scaling across the entire range of values.

2.4.3.4 L2 Normalization

L2, or vector normalization, [32] scales each data point so that the Euclidean norm (the square root of the sum of squared values) is 1:

$$X_{norm} = \frac{X}{\sqrt{\sum X_i^2}} \quad (2.14)$$

This method is commonly used in NLP (e.g., TF-IDF vectors) where the magnitude of the vector matters less than its direction. However, it can be sensitive to outliers, as large values can disproportionately affect the norm.

2.4.4 Vectorization

Vectorization is a critical step in pre-processing textual data for ML models. It converts text-based input into numerical vectors, which are the only format that most ML algorithms can process. By transforming text into numerical representations, vectorization enables algorithms to comprehend patterns and relationships between words and documents, allowing for tasks such as text classification, sentiment analysis, and topic modeling.

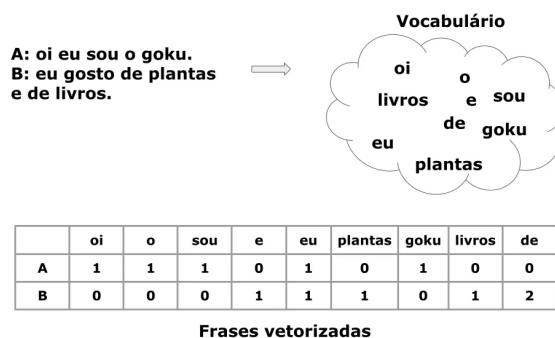
In NLP, vectorization is key to enabling computers to process and understand human language. Below, we will explore several vectorization techniques that play essential roles in NLP tasks.

2.4.4.1 Bag of Words (BoW)

The BoW method is one of the simplest and most commonly used techniques for converting text into a numerical format [65]. In this method, each document is represented as a vector of word counts, disregarding grammar and word order but maintaining the frequency of each word in the document.

The process can be broken down into four steps: tokenization, vocabulary construction, word count, and matrix construction.

1. **Tokenization:** The text is split into individual tokens, which are often words but could also be n-grams (sequences of n consecutive words).
2. **Vocabulary Construction:** A list of all unique words (or n-grams) across the dataset is compiled, forming the vocabulary.
3. **Word Count:** For each document, the occurrence of each word from the vocabulary is counted.
4. **Matrix Construction:** Finally, a matrix is created where each row represents a document, and each column corresponds to a word in the vocabulary. The values in the matrix represent the frequency of each word in each document.

Figure 2.13: *The creation of vectors*

BoW is simple and effective, but it has limitations. It ignores word order and semantic relationships between words, treating all words as independent entities. This limits its usefulness in understanding the deeper context of a document. For more complex tasks, models like Word2Vec or BERT are needed to capture semantic meaning.

2.4.4.2 TF-IDF (Term Frequency-Inverse Document Frequency)

The TF-IDF approach builds on the simplicity of BoW by weighing the importance of words in a document relative to their frequency across the entire corpus [62]. It balances two concepts: **Term Frequency (TF)**, which measures how frequently a word appears in a document, and **Inverse Document Frequency (IDF)**, which measures how rare or common a word is across all documents. This technique highlights words that are frequent in a particular document but not too common across many documents, thus assigning them higher importance.

- **Term Frequency (TF)**: Measures the frequency of a term in a document.

$$TF = \frac{\text{Number of times word "X" appears in a document}}{\text{Total number of words in the document}} \quad (2.15)$$

- **Inverse Document Frequency (IDF)**: Measures the rarity of the term across all documents.

$$IDF = \log\left(\frac{\text{Total number of documents}}{\text{Number of documents containing word "X"}}\right) \quad (2.16)$$

- **TF-IDF**: Combines both metrics to assign weight to each word.

$$TFIDF = TF * IDF \quad (2.17)$$

TF-IDF is particularly useful in applications like search engines and document classification, as it reduces the weight of common words (such as "and" or "the") while emphasizing less frequent, more informative terms.

2.4.4.3 Word2Vec

In BoW and TF-IDF, each word is treated as an individual entity, and semantics is completely ignored, unlike Word2Vec.

Word2Vec aims to capture semantic and syntactic relationships between words from the context in which they occur in large data sets [72]. Each word is represented as an n-dimensional vector, all words

are mapped onto this n-dimensional space in such a way that words with similar meanings exist in close proximity to each other.

Word2Vec vectors can be visualized and used to find semantic relationships and analyze the context in which words are used.

There are two main approaches: Skip-gram and CBOW (Continuous Bag of Words).

The Skip-gram method provides the neural network with a word, and it has to predict the context in which it is inserted, i.e. which words normally occur next to it. It predicts the probability of each word in the vocabulary being close to the word entered into the model.

CBOW is the inverse of skip-gram, i.e. instead of predicting words from the same context as the word, the vocabulary is inserted into the model to predict the current word.

Skip-gram works better for smaller data sets and better represents rare words, while CBOW trains faster and better represents frequent words.

2.4.4.4 GloVe (Global Vectors for Word Representation)

GloVe is another method for generating word embeddings, similar to Word2Vec, but it differs in its approach to capturing word relationships. Developed by Stanford, GloVe constructs a word co-occurrence matrix that records how frequently words appear together in a corpus [22]. The intuition behind GloVe is that words with similar meanings will have similar patterns of co-occurrence with other words.

The goal of GloVe is to find word vectors where the dot product between two vectors approximates the logarithm of the probability that the two words co-occur. By leveraging statistical information from the entire corpus, GloVe produces word vectors that capture both semantic similarity and word frequency.

GloVe tends to perform well on word analogy tasks (e.g., "king" is to "queen" as "man" is to "woman") and is widely used in NLP applications.

2.4.4.5 FastText

FastText, developed by Facebook AI Research, builds on the strengths of Word2Vec by considering subword information [67]. Instead of treating each word as a single unit, FastText breaks words into character n-grams, which makes it particularly useful for handling rare words or languages with complex morphology.

For example, the word "playing" might be broken down into "pla", "lay", "ayi", and so on. This allows FastText to generate word vectors even for words that were not present in the training corpus, improving performance on tasks involving rare or misspelled words.

FastText has become a popular choice for text classification tasks and has been integrated into various NLP systems due to its ability to generalize well across different languages and text domains.

2.4.5 Feature Selection

Feature selection is the process of selecting the most relevant features in a dataset for building a ML model [4]. It helps improve the model's accuracy, reduce overfitting, and decrease computation time.

In ML, irrelevant or redundant features can negatively impact the performance of models, making feature selection a vital step in data preprocessing.

There are three broad categories of feature selection techniques: **Filter Methods**, **Wrapper Methods**, and **Embedded Methods**. These approaches can be applied based on the complexity of the dataset, the size of the feature space, and the model's performance goals.

2.4.5.1 Filter Methods

Filter methods are based on general statistics and characteristics of the data. These methods are independent of the ML model and use statistical techniques to rank and select the most relevant features.

2.4.5.2 Correlation Coefficient

The correlation coefficient measures the linear relationship between two variables. In feature selection, it is used to identify features that are strongly correlated with the target variable [66]. Features with low correlation to the target are often removed.

- **Pearson's Correlation:** Measures linear correlation between continuous variables.

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (2.18)$$

- $r = 1$ implies perfect positive linear correlation.
- $r = -1$ implies perfect negative linear correlation.
- $r = 0$ implies no linear relationship.

- **Spearman's Rank Correlation:** Measures monotonic relationships and can be used for both continuous and ordinal data.

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (2.19)$$

- $\rho = 1$ indicates a perfect positive monotonic relationship.
- $\rho = -1$ indicates a perfect negative monotonic relationship.
- $\rho = 0$ suggests no monotonic relationship.

This method works best when the relationships between features and the target are linear. However, it may miss non-linear relationships between variables.

2.4.5.3 Chi-Square Test

The **Chi-Square Test** [8] is used for categorical data to test the independence between features and the target variable. It calculates the difference between the observed frequency and expected frequency of events. A higher chi-square value indicates that the feature has a stronger association with the target.

$$x^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad (2.20)$$

Where:

- O_i is the observed frequency.
- E_i is the expected frequency.

This method is particularly useful in classification problems, but it only works for categorical variables.

2.4.5.4 Information Gain (Entropy)

Information Gain measures the reduction in uncertainty (entropy) about the target variable given a particular feature [25]. It evaluates how much information a feature provides about the target. The higher the information gain, the more important the feature. Entropy ($H(X)$) is given by:

$$H(X) = - \sum_i P(x_i) \log P(x_i) \quad (2.21)$$

Where $P(x_i)$ is the probability of outcome x_i . Information gain (IG) for a feature X is the reduction in entropy:

$$IG = H(Y) - H(Y|X) \quad (2.22)$$

When $H(Y)$ is the entropy of the target variable, and $H(Y|X)$ is the entropy after observing the feature X .

2.4.5.5 Variance Thresholding

Variance Thresholding [68] removes features with low variance across samples, as features that don't vary much are unlikely to be useful in prediction tasks. It works well when you need to remove constant or near-constant features from the dataset. The method compares the variance of each feature to a predefined threshold and discards features with variance below the threshold.

$$Variance = \frac{\sum (x_i - \bar{x})^2}{N} \quad (2.23)$$

Where x_i represents each data point, \bar{x} is the mean, and n is the total number of samples.

This method is simple to implement but doesn't account for feature-target relationships.

2.4.5.6 Wrapper Methods

Wrapper methods use a ML model to evaluate feature subsets. The model is trained on different subsets of features, and the performance of each subset is evaluated. Wrapper methods are more computationally expensive but typically yield better results as they evaluate feature importance based on model performance.

Recursive Feature Elimination (RFE) RFE recursively removes the least important feature based on model performance and re-trains the model on the remaining features. The process is repeated until the desired number of features is reached. RFE ranks features by their importance scores and eliminates those that contribute the least to the model [50].

1. Train the model on all features.
2. Rank features by their importance.

3. Remove the least important feature and re-train the model.
4. Repeat steps 2-3 until the optimal number of features is reached.

RFE is widely used with algorithms that provide feature importance scores, such as decision trees, linear models, or SVMs.

Forward Selection

Forward Selection starts with an empty model and iteratively adds features that improve the model's performance [20]. At each step, the feature that most improves the performance (based on a scoring metric like AIC, BIC, or cross-validation score) is added to the model.

1. Start with no features.
2. Train the model and evaluate performance with each feature added.
3. Add the feature that improves model performance the most.
4. Repeat until no significant improvement is achieved.

This method works well but can be computationally expensive for large datasets.

Backward Elimination

In contrast to Forward Selection, Backward Elimination starts with all features and iteratively removes the least important feature. The feature whose removal improves model performance (or least deteriorates it) is eliminated.

1. Start with all features.
2. Train the model and evaluate performance with each feature removed.
3. Remove the feature whose removal improves performance.
4. Repeat until further removal of features no longer improves performance.

Backward Elimination works well with datasets that have a smaller number of features, as it evaluates the contribution of each feature.

2.4.5.7 Embedded Methods

Embedded methods perform feature selection during the model training process. These methods are more efficient than wrapper methods and often provide better generalization by considering feature interaction while training the model.

Lasso (L1 Regularization) Lasso regression [33] adds an $L1$ penalty to the loss function, which forces the coefficients of less important features to shrink to zero, effectively performing feature selection. The Lasso objective function is:

$$\min_{\omega} \left(\sum_{i=1}^n (y_i - X_i \cdot w)^2 + \lambda \sum_{j=1}^p |w_j| \right) \quad (2.24)$$

Where:

- y_i is the target variable,
- X_i is the feature matrix,
- w_j are the feature weights,
- λ is the regularization parameter.

This method is useful for high-dimensional data and tends to select a small subset of important features. However, it may struggle when multiple features are highly correlated.

Ridge Regression (L2)

Ridge regression adds an L2 penalty to the linear regression loss function, which shrinks the coefficients of less important features but doesn't reduce them to zero [33]. The objective function is:

$$\min_{\omega} \left(\sum_{i=1}^n (y_i - X_i \cdot w)^2 + \lambda \sum_{j=1}^p w_j^2 \right) \quad (2.25)$$

Ridge regression is useful when there are many small but important features, as it shrinks their coefficients without completely eliminating them.

Tree-Based Feature Selection

Decision tree-based models (e.g., Random Forests, Gradient Boosting) [64] naturally rank features by their importance. Features that are frequently used to make splits in the decision trees and contribute significantly to reducing the impurity or error are considered important.

Tree-based methods are powerful as they can handle both numerical and categorical data, capture non-linear relationships, and deal with large datasets.

2.4.6 Dimensionality Reduction

Dimensionality reduction is the process of reducing the number of independent variables in a dataset. This is useful for several reasons, including the elimination of redundancy and the simplification of analysis. Reducing the dimensionality can help to remove noise from the data and prevent overfitting, thus improving the performance and efficiency of ML models.

2.4.6.1 Principal Component Analysis (PCA)

PCA is a mathematical procedure that performs dimensionality reduction and is widely used in Data Science. The objective is to simplify and preserve the most important information in a high-dimensional dataset by reducing the number of variables while maintaining the essential characteristics of the data. This simplification can drastically reduce the time required to train algorithms on large datasets.

PCA highlights the similarities and differences in the existing data by identifying patterns [49]. The process starts by calculating the **covariance matrix**, which measures how two variables vary together. A positive covariance means that the variables tend to increase or decrease together, while a negative covariance indicates that one variable tends to increase when the other decreases.

Next, **eigenvectors** and **eigenvalues** of the covariance matrix are calculated. The eigenvectors represent the directions of maximum variance in the data, while each eigenvector is associated with an

eigenvalue that represents the amount of variance explained by that direction. The eigenvectors are then ordered based on the eigenvalues, from largest to smallest.

The original data is projected onto new feature spaces defined by the selected eigenvectors. This projection creates new variables called "principal components." These new variables are linear combinations of the original variables, with weights determined by the eigenvectors.

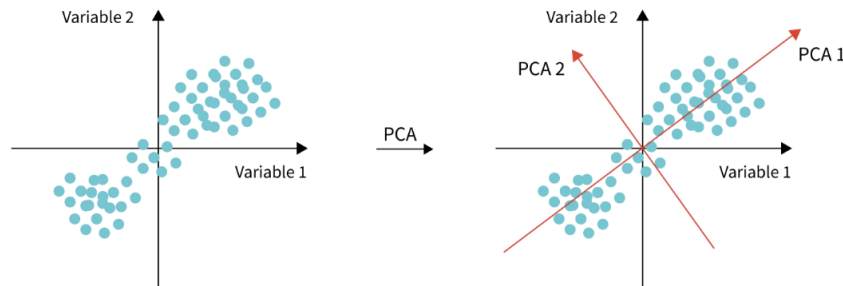


Figure 2.14: *Example of PCA*

It's important to note that PCA is a **linear** technique and assumes that the relationships between variables are linear. In scenarios where there are non-linear relationships, other non-linear dimensionality reduction techniques may be more appropriate.

2.4.6.2 t-Distributed Stochastic Neighbor Embedding (t-SNE)

t-SNE is a dimensionality reduction technique that is particularly useful for **visualization** [18]. It attempts to map data points from a high-dimensional space to a low-dimensional space (typically 2D or 3D) in a way that preserves the proximity relationships between them, as shown in Figure 2.15. It is particularly effective at preserving local structures and revealing clusters in data.

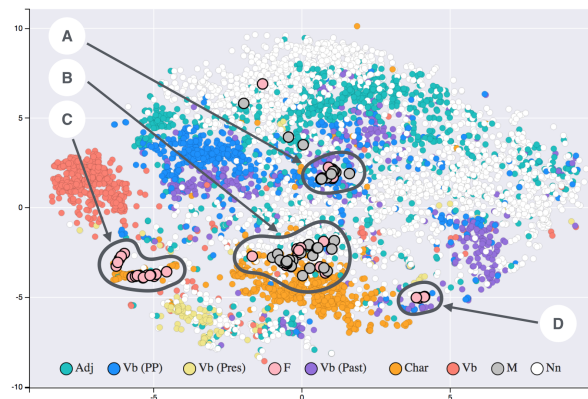


Figure 2.15: *Example of t-SNE technique*

The method uses a probabilistic approach to model the similarity between pairs of points in both high- and low-dimensional spaces, following four basic steps:

1. **Similarity Measurement (High Dimensionality)**: Calculate the conditional probabilities

$P(i|j)$, representing the probability of choosing point j as a neighbor of point i in high-dimensional space. This probability is typically computed using a Gaussian distribution based on the similarities between points.

2. **Similarity Measurement (Low Dimensionality):** Perform a similar process for the low-dimensional space, but using a Student t-distribution with a small degree of freedom to handle dense clusters better.
3. **Minimization of Kullback-Leibler Divergence:** Minimize the **Kullback-Leibler divergence** between the probability distributions in the high- and low-dimensional spaces.
4. **Stochastic Gradient Descent:** Use stochastic gradient descent to iteratively adjust the positions of the points in the low-dimensional space to minimize the KL divergence. The gradient is calculated with respect to the positions of the points in the low-dimensional space.

2.4.6.3 Linear Discriminant Analysis (LDA)

Linear Discriminant Analysis (LDA) is a supervised dimensionality reduction technique used primarily for classification tasks [36]. It finds a linear combination of features that best separates the different classes in the data. LDA works by maximizing the **between-class variance** while minimizing the **within-class variance**, making it particularly useful for multi-class problems.

$$\text{Maximizing } \frac{S_B}{S_W} \tag{2.26}$$

Where S_B is the between-class scatter matrix, and S_W is the within-class scatter matrix.

LDA is especially useful when the goal is to reduce dimensionality while preserving class discriminatory information.

2.5 Natural Language Processing and Topic Modeling

NLP is a subfield of AI and linguistics that focuses on the interaction between computers and human languages. The goal of NLP is to enable machines to understand, interpret, and generate human language in a meaningful way.

NLP techniques help in tasks like:

- **Tokenization:** Tokenization is a fundamental step in text processing, enabling further analysis and modeling by transforming text into a structured format.
- **Named Entity Recognition (NER):** NER helps in extracting and categorizing important information from text, which is useful for information retrieval, knowledge management and data organization.
- **Sentiment Analysis:** Sentiment analysis is widely used in social media monitoring, customer feedback analysis and brand reputation management to gauge public opinion and emotional tone.
- **Text Summarization:** Automatically generating concise summaries from large text documents.
- **Machine Translation:** Translating text from one language to another.

2.5.0.1 Latent Dirichlet Allocation (LDA)

LDA is a generative probabilistic model used for topic modeling. It assumes that documents are mixtures of topics, and topics are distributions over words. LDA identifies latent topics within a collection of text by assigning probabilities to each word for each topic.

LDA uses a bag-of-words approach, which ignores word order and context. While effective for many applications, LDA struggles with capturing the semantics and relationships between words in context, which is why transformer-based models like BERT and BERTopic offer improved performance.

2.5.0.2 BERTopic

BERTopic is a sophisticated NLP technique built on top of **Bidirectional Encoder Representations from Transformers** and designed to perform **topic modeling** by clustering text documents into semantically similar groups. It excels in capturing context and relationships between words, leveraging pre-trained transformer models to extract meaningful topics from large text corpora. This makes BERTopic especially valuable for discovering hidden structures in unstructured data such as customer reviews, social media posts, and other text-heavy datasets.

Before diving into BERTopic, it's important to understand the foundation it builds upon: **BERT**. Developed by Google, BERT uses transformer architecture to learn **bidirectional representations** of words in context. Unlike traditional NLP models, which process input data in a single direction (either left-to-right or right-to-left), BERT considers both preceding and succeeding words in a sentence, making it capable of fully understanding the context in which a word appears.

BERT was trained on a massive corpus of text, including over **800 million words from BooksCorpus** and **2.5 billion words from Wikipedia**, with the model designed to perform tasks like:

- **Masked Language Modeling (MLM)**: Predicting missing words in a sentence.
- **Next Sentence Prediction (NSP)**: Predicting the relationship between pairs of sentences.

This bidirectional approach allows BERT to generate more context-aware embeddings, making it highly effective for various NLP tasks.

BERTopic integrates BERT embeddings with **dimensionality reduction** and **clustering algorithms** to identify and group similar documents. The workflow of BERTopic involves:

1. **Embedding Extraction**: BERTopic starts by using a pre-trained BERT model to extract embeddings for each document in the corpus.
2. **Dimensionality Reduction**: BERTopic applies **UMAP (Uniform Manifold Approximation and Projection)**, a non-linear dimensionality reduction technique. UMAP reduces the dimensionality while preserving the local and global structure of the embeddings, making them more manageable for clustering.
3. **Clustering**: The reduced embeddings are then passed through **HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise)**. This density-based clustering algorithm groups similar documents into clusters (topics) without needing to specify the number of clusters beforehand.

4. **Similarity Calculation:** Using the embeddings, BERTopic calculates document similarities. Documents that are semantically close to each other are more likely to belong to the same cluster, forming the basis for topic identification.
5. **Topic Representation:** After clustering, BERTopic represents each topic by identifying the most representative words in the documents belonging to that cluster.

BERTopic is particularly useful for generating meaningful topics from unstructured text, such as customer reviews or social media data.

2.6 Interpretability and Explainability

As organizations increasingly rely on ML algorithms to make critical decisions, understanding how these algorithms function becomes essential. Two key concepts in this regard are **interpretability** and **explainability**.

Interpretability refers to the ability to understand how a ML model arrives at its decisions. It ensures that the decision-making process of AI models is transparent, enabling users to trust and verify the output of intelligent systems. Meanwhile, **explainability** focuses on providing clear, intuitive explanations of specific decisions made by a model, allowing users to understand why the model produced certain results in particular instances. Essentially, explainability emphasizes how the model's decisions can be justified and communicated in an understandable manner.

These concepts are vital for ensuring that models function as intended, reducing the risk of errors, and fostering trust among users. With explainability, users can verify that the model is operating fairly, without bias, and making decisions based on sound reasoning.

There are several techniques that improve both interpretability and explainability, including:

- **Visualization Methods:** Graphical representations of data and models help clarify how the model works by visually showing relationships and patterns in the data.
- **Model Decomposition:** Breaking down a complex model into simpler components can help users understand how each part contributes to the final decision.
- **Example-Based Explanations:** Providing concrete examples of how the model behaves under specific conditions helps illustrate its decision-making process.
- **Feature Importance:** By identifying which features (input variables) are most influential in the model's decisions, users can better understand what drives predictions.

Several tools and methods have emerged in recent years to enhance explainability and interpretability, such as **SHAP (SHapley Additive exPlanations)** and **LIME (Local Interpretable Model-agnostic Explanations)**, which offer insights into how individual model predictions are made. These tools break down predictions for complex models, such as neural networks or gradient-boosted trees, providing interpretable feedback for users.

By improving the transparency of ML models, organizations can ensure that AI-driven decisions align with their ethical and operational standards. In sensitive domains, such as healthcare or finance, explainability is not just a best practice but often a regulatory requirement to ensure fairness, accountability, and compliance with legal frameworks.

2.7 Clustering Methods

Clustering methods are unsupervised learning techniques used to group data points into more compact sets or *clusters* without predefined labels or additional supervision. The goal of clustering is to uncover hidden patterns and structures in the data that partition the dataset into meaningful clusters [53]. These methods are critical for exploratory data analysis, where relationships and patterns within the dataset are not known in advance.

In Figure 2.16 the left panel shows a dataset before clustering, while the right panel illustrates how the dataset is divided into three distinct clusters after applying a clustering method.

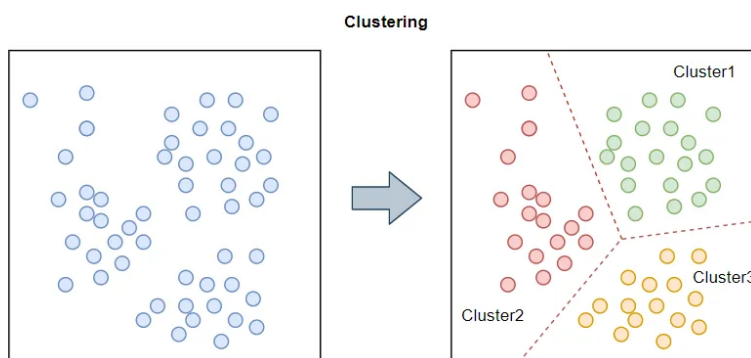
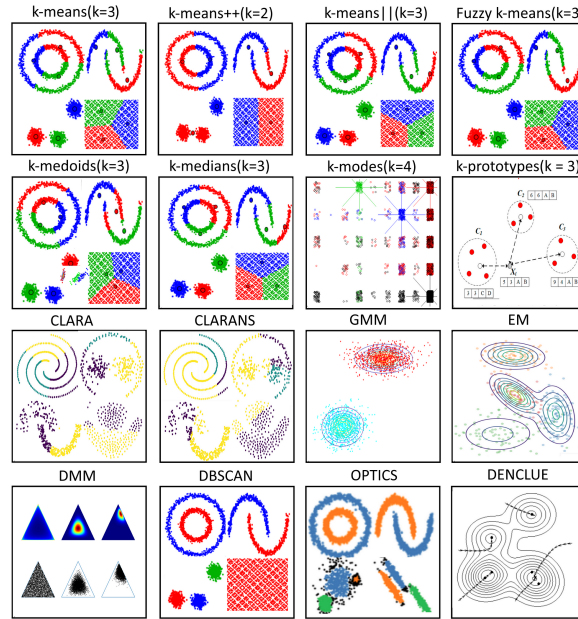


Figure 2.16: *Clustering Process*

Clustering methods differ from supervised learning approaches like classification and regression because they do not rely on labeled data. Instead, they autonomously identify relationships, similarities, and groupings within the dataset.

Clustering algorithms are broadly categorized into five types:

1. **Centroid-based Clustering:** Clusters are defined by their centroids, and data points are assigned to the nearest centroid based on proximity.
2. **Density-based Clustering:** Identifies regions of high density and clusters them, while points in lower-density regions may be treated as outliers.
3. **Distribution-based Clustering:** Assumes that data follows statistical distributions, such as Gaussian distributions, and groups points based on the probability of belonging to a distribution.
4. **Hierarchical Clustering:** Builds a tree of clusters by iteratively merging or splitting groups of data points.
5. **Other Methods:** Includes approaches like Affinity Propagation and Mean Shift, which do not fit neatly into the other categories.

Figure 2.17: *Examples of Different Clustering Algorithms*

Let's explore specific clustering methods within these categories.

2.7.1 Centroid-based Clustering

Centroid-based clustering organizes data into clusters by assigning each data point to the nearest centroid. The centroids represent the center of each cluster, and the goal is to minimize the distance between points and their assigned centroids. These algorithms aim to partition the data into non-hierarchical clusters. Popular algorithms like **K-Means** and **Mini-Batch K-Means** are known for their computational efficiency, especially with large datasets. Despite their sensitivity to initial conditions and outliers, centroid-based methods remain widely used due to their simplicity and effectiveness across various domains. In this section, we will delve into the most common centroid-based clustering algorithms: K-Means and Mini-Batch K-Means.

2.7.1.1 K-Means

K-Means [55] is one of the most popular and straightforward clustering methods. This algorithm is based on the concept of centroids, where each data point is assigned to the nearest centroid based on similarity. Various measures can be used to determine this similarity or distance between points, such as Euclidean and Manhattan distances. It operates by initializing k centroids randomly from the dataset, as depicted in Figure 2.18.

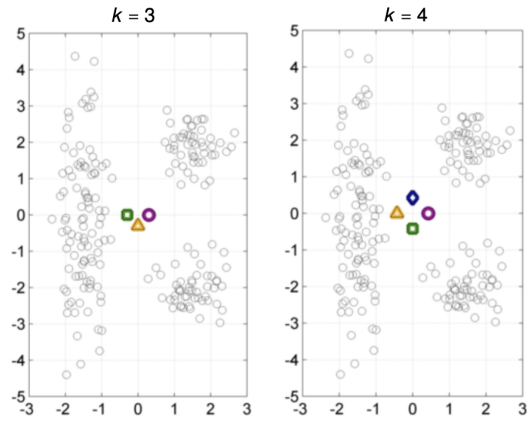


Figure 2.18: *Initial Data and Centroids*

Each data point is then assigned to the nearest centroid based on similarity, typically calculated using metrics like Euclidean or Manhattan distances, as shown in Figure 2.19

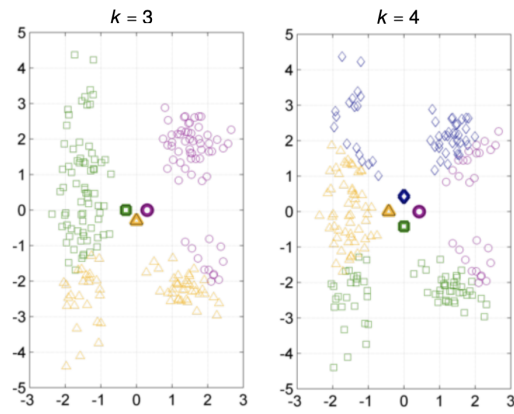


Figure 2.19: *Assigning Points to Clusters*

The algorithm then updates the centroids by calculating the average of all points assigned to each cluster, as shown in Figure 2.20. This process repeats until the centroids no longer change significantly.

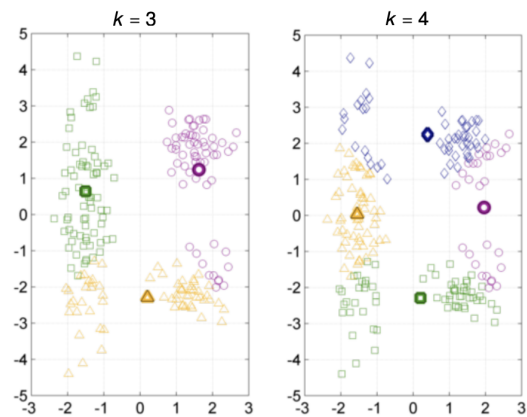


Figure 2.20: *Updating Centroids*

K-Means is favored for its simplicity and efficiency, making it widely used in various applications. However, it does have limitations.

- **Predefined k :** The number of clusters k must be specified before running the algorithm, which may require domain knowledge or trial and error.
- **Sensitivity to Outliers:** Outliers can heavily influence centroids, leading to distorted cluster shapes.
- **Assumption of Spherical Clusters:** K-Means assumes that clusters are spherical and of equal size, which might not be true for many real-world datasets.

In conclusion, while K-Means provides a straightforward approach to clustering, careful consideration of its assumptions and handling of outliers is essential to ensure accurate and meaningful clustering results for diverse datasets.

2.7.1.2 Mini Batch K-Means

Mini-Batch K-Means is a variant of the traditional K-Means clustering algorithm designed for improved efficiency, particularly with large datasets [69]. Unlike standard K-Means, which updates centroids using the entire dataset at each iteration, Mini-Batch K-Means operates on randomly sampled mini-batches of the data. This approach offers significant advantages in terms of computational efficiency and memory usage.

At the beginning of the algorithm, like K-Means, Mini-Batch K-Means initializes k centroids randomly or using a specific initialization method. However, instead of updating centroids based on all data points, it processes a subset (mini-batch) of the data in each iteration. The size of the mini-batch, often determined by a parameter called the *batch size*, influences the algorithm's performance and convergence speed.

During each iteration, Mini-Batch K-Means calculates the distances between the data points in the mini-batch and the current centroids, assigning each point to the nearest centroid. It then updates the centroids by computing the mean of all data points assigned to each centroid within the mini-batch.

Figure 2.21 visually illustrates the difference between standard K-Means and Mini-Batch K-Means. The left side shows the traditional approach where all data points are used in each iteration, while the right side demonstrates Mini-Batch K-Means processing a smaller subset of data.

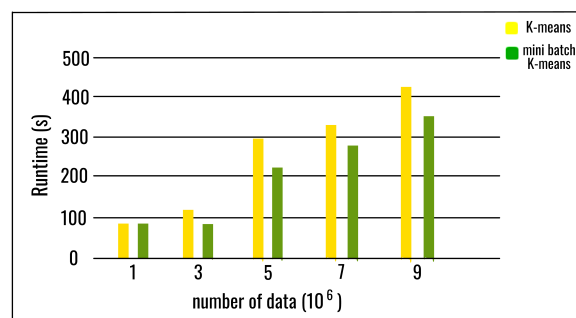


Figure 2.21: Comparison of standard K-Means and Mini-Batch K-Means

Due to its reliance on mini-batches, Mini-Batch K-Means may produce slightly different clustering results compared to standard K-Means, especially when the batch size is relatively small. However, these differences are typically minor and outweighed by the algorithm's advantages in terms of speed and scalability.

Overall, Mini-Batch K-Means is particularly beneficial for applications dealing with large datasets or real-time data streams where computational resources are limited. By processing data in smaller batches, it strikes a balance between accuracy and efficiency, making it a practical choice for many clustering tasks.

2.7.2 Density-based Clustering

Density-based clustering identifies clusters by detecting areas of high data density, separated by regions of lower density. Unlike centroid-based methods, density-based clustering can find clusters of arbitrary shapes and sizes, making it particularly useful for datasets where clusters are not well-separated or follow non-spherical distributions. Points in sparse regions are often treated as noise or outliers, making these methods robust to noise. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** is a widely recognized density-based algorithm that excels in detecting irregularly shaped clusters. This method is especially advantageous when the number of clusters is unknown and when noise is present in the data.

2.7.2.1 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

The DBSCAN (Density-Based Spatial Clustering of Applications with Noise) [21] algorithm is notable for its density-based approach to clustering. It can identify clusters in datasets with varying densities and irregular shapes without the need to specify the number of clusters beforehand. Clusters are defined as areas of high density within the data domain.

DBSCAN requires two main parameters:

- *epsilon*: The maximum radius within which two points are considered to be part of the same cluster;
- *min_samples*: The minimum number of points required within a region to ensure a desired density level.

The choice of parameters is crucial for the algorithm's performance. If *epsilon* is set too low, many points will be classified as noise because they won't have enough neighbors within the specified radius. Conversely, if *epsilon* is too high, the algorithm may produce overly generalized clusters. For *min_samples*, higher values are preferable for noisy datasets, as they help to ensure that clusters are formed from genuinely dense regions.

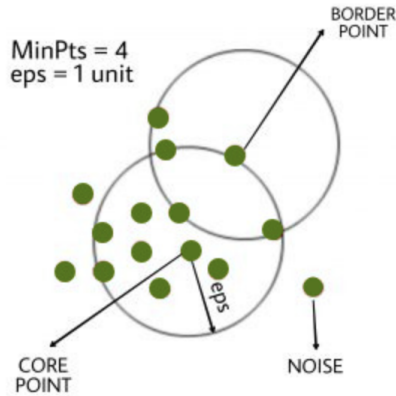


Figure 2.22: *Types of points in the DBSCAN algorithm*

In DBSCAN, data points are categorized into three types:

- *Core point*: A point that has at least $min_samples$ within its $epsilon$ radius;
- *Border point*: A point that has fewer than $min_samples$ within its $epsilon$ radius but is within the $epsilon$ radius of a core point;
- *Outlier*: A point that is neither a core point nor a border point. These points are considered noise and are not assigned to any cluster.

The DBSCAN algorithm starts by selecting a random point in the dataset and checking how many points are within its $epsilon$ radius. If the number of points within this neighborhood is greater than or equal to $min_samples$, the point is classified as a core point.

From a core point, the algorithm expands the cluster by including all points within the $epsilon$ distance from the core point. This process is performed recursively, including the neighbors of the newly added points into the cluster, as illustrated in Figure 2.23.

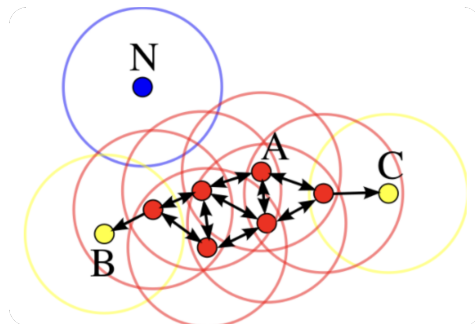


Figure 2.23: *Example of the DBSCAN algorithm in action*

DBSCAN is highly effective at detecting noise and identifying outliers, often being used specifically for these purposes in many applications. Its ability to handle data with varying densities and its robustness against noise make it a powerful clustering method for complex datasets.

1. **Initialization:** Treat each data point as a single cluster;
2. **Similarity Calculation:** Compute the similarity (or distance) between all pairs of clusters;
3. **Cluster Merging:** Merge the two closest clusters and update the similarity matrix;
4. **Iteration:** Repeat the similarity calculation and merging steps until only one cluster remains;
5. **Dendrogram Construction:** Throughout the process, build a dendrogram to visualize the merging steps and the hierarchical structure of the data.

Agglomerative Hierarchical Clustering is appreciated for its simplicity and its ability to capture the nested structure of data. However, it can be computationally intensive for large datasets due to the repeated similarity calculations and merges. Despite this, it remains a popular choice for exploratory data analysis and understanding data relationships.

2.7.4.2 Divisive Hierarchical

Divisive Hierarchical Clustering is another hierarchical clustering technique that takes an opposite approach to Agglomerative Hierarchical Clustering. Instead of starting with each data point as its own cluster, Divisive Hierarchical Clustering begins with all data points in a single cluster and progressively splits it into smaller subclusters until a stopping criterion is met.

To illustrate the difference between Agglomerative and Divisive approaches, refer to Figure 2.25.



Figure 2.25: *Comparison of Agglomerative and Divisive Clustering*

The Divisive algorithm starts by considering the entire dataset as a single cluster and progressively splits it into smaller subclusters. The goal is to separate data points into more distinct clusters. The division technique can vary and may be based on density, distance, or other similarity metrics. This splitting step is applied recursively to each resulting subcluster, meaning each subcluster is further divided into even smaller clusters.

The process continues until a stopping criterion is met. This criterion can be based on the desired number of clusters or a metric that indicates the quality of the division. The choice of the splitting metric significantly influences the clustering results.

During the process, a dendrogram can be constructed, providing a graphical representation of the cluster hierarchy and showing how clusters were divided throughout the iterations.

Divisive Hierarchical Clustering is particularly useful for analyzing the hierarchical structure of data and can offer valuable insights. However, it may be more sensitive to the choice of the splitting metric and the specified number of clusters.

Algorithm Steps:

1. **Initialization:** Start with all data points in a single cluster.
2. **Division:** Split the overall cluster into smaller subclusters based on a chosen metric.
3. **Recursion:** Apply the division step recursively to each subcluster.
4. **Stopping Criterion:** Continue the process until the desired number of clusters is achieved or another stopping criterion is met.
5. **Dendrogram Construction:** Build a dendrogram to visualize the hierarchical structure and the division process.

This approach is valuable for exploratory data analysis, offering a detailed view of the data's hierarchical structure. However, careful consideration is needed in selecting the division metric and determining the number of clusters.

2.7.4.3 BIRCH (Balance Iterative Reducing and Clustering using Hierarchies)

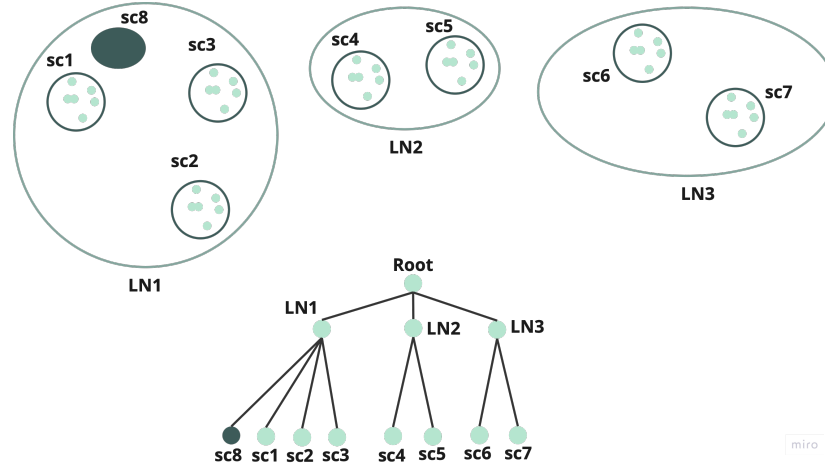
BIRCH is a hierarchical clustering algorithm developed to handle large datasets efficiently. It is particularly effective in scenarios with limited memory, operating incrementally and leveraging data summarization techniques.

The two most important parameters in the BIRCH algorithm are the *Branching_Factor* and the *Threshold*. The *Branching_Factor* defines the maximum number of children a node can have in the tree, influencing the tree's width. The *Threshold* sets the limit of data points a node can store before it needs to be split or merged.

The algorithm begins by creating a tree structure called the Clustering Feature Tree (CF Tree), which is a hierarchical representation of the clusters.

To insert a point into the tree, the distance between the point and the existing points in the node is calculated. If the distance is within the *Threshold*, the point is added to the node. If the distance exceeds the threshold, the node is split, or the point is merged with an existing point. When point insertion is complete, clusters are formed from the nodes of the tree, with each node representing a cluster.

To enhance the accuracy of the clusters, a more traditional clustering algorithm (such as K-means) can be applied to the points within each node.

Figure 2.26: *BIRCH Clustering Process*

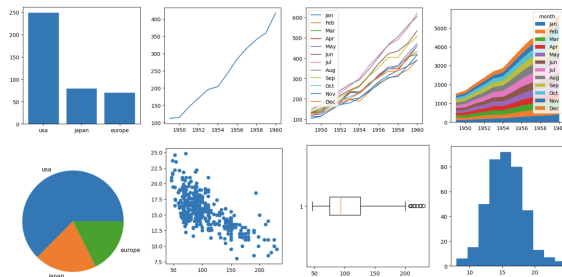
BIRCH is highly efficient in managing large datasets due to its incremental nature and data summarization capabilities. It is especially useful when working with data that cannot fit entirely into memory, as it constructs and refines clusters dynamically.

2.8 Visualization Libraries

Data visualization plays a key role in understanding distributions, relationships, and potential anomalies within datasets. Several visualization libraries can help provide these insights by creating intuitive and informative visual representations. Let's review some of the most prominent Python libraries for data visualization:

2.8.0.1 Matplotlib

Matplotlib is one of the oldest and most popular Python libraries for data visualization. It offers extensive flexibility, allowing users to create a wide variety of 2D and 3D visualizations, ranging from simple line plots to more complex figures like histograms, bar charts, scatter plots, and pie charts. Matplotlib's versatility makes it highly valuable for static visualizations, especially when detailed customization is required [56]. Its wide adoption ensures strong community support and ample resources for troubleshooting.

Figure 2.27: *Graph visualization using the Matplotlib library*

2.8.0.2 Plotly

Plotly is renowned for its ability to create interactive and dynamic visualizations. It allows users to create both 2D and 3D graphs that can be embedded into web applications. Plotly’s interactive features make it ideal for visualizing large datasets, offering drag-and-zoom functionality, tooltips, and real-time updates. Additionally, it supports a range of visualizations, including bar charts, scatter plots, heatmaps, and choropleth maps.

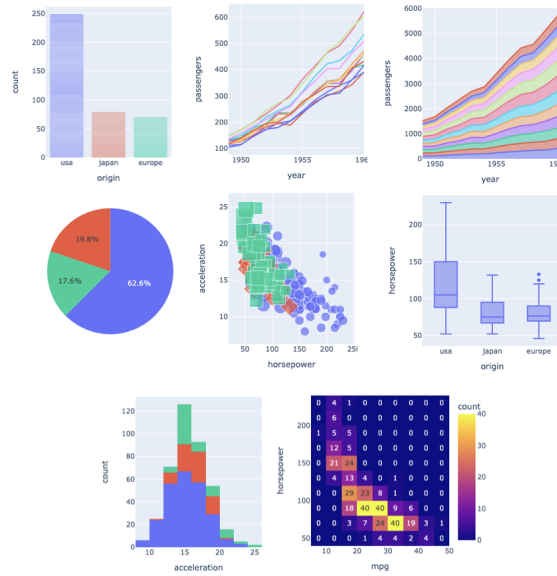


Figure 2.28: *Graph visualization using the Plotly library*

2.8.0.3 Bokeh

Bokeh is another interactive visualization library that focuses on enabling high-performance visualizations within web browsers. With support for a wide range of chart types, Bokeh allows developers to create interactive, aesthetically pleasing plots with minimal effort. This library is particularly useful for visualizing time-series data and large datasets, making it a go-to for interactive dashboards or data-driven applications.

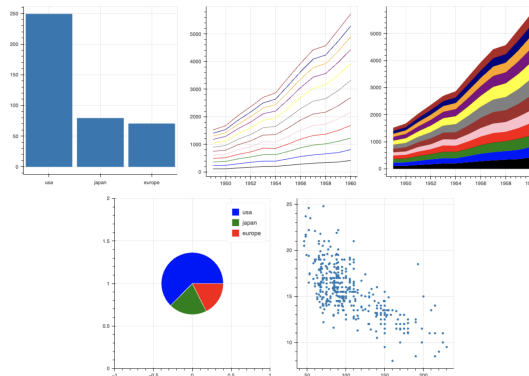


Figure 2.29: *Graph visualization using the Bokeh library*

2.8.0.4 Altair

Altair is a declarative visualization library that simplifies the process of creating visually appealing and meaningful graphs. It emphasizes a concise and human-readable syntax, making it easy to create basic and complex visualizations. Altair is particularly well-suited for exploratory data analysis, offering a smooth experience when working with bar charts, scatter plots, and histograms.

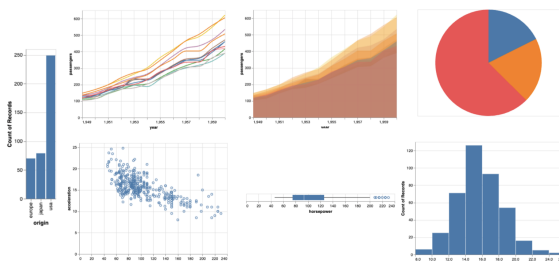


Figure 2.30: *Graph visualization using the Altair library*

These libraries provide diverse options for visualizing data, ranging from static plots to dynamic, interactive graphs. The choice of the right tool depends on the specific needs of the project, such as interactivity, customization, or simplicity.

2.9 Visualization of Model Results

Visualizing model [47] results is critical in ML for various reasons:

1. **Classifier Evaluation:** To assess classifier performance, visual tools such as confusion matrices, ROC curves, and precision-recall curves are used.
2. **Model Interpretability:** Visualization enhances understanding of model decisions, improving the transparency of complex models.
3. **Detecting Overfitting or Underfitting:** Learning curves visually show how well a model generalizes based on training and test data, helping identify overfitting or underfitting.
4. **Hyperparameter Tuning:** Visualization tools help optimize hyperparameters by showcasing their impact on model performance.
5. **Result Presentation:** Clear and informative visualizations are essential for effectively communicating project outcomes.

To visualize the results of the algorithm, there are a number of libraries that can be used. Scikit-learn offers various functions including the possibility of displaying ROC curves and confusion matrices which are very useful.

Chapter 3

Data Description

In the 2.3 section, the importance of using a reliable and versatile dataset to extract meaningful insights was highlighted [1]. This chapter provides an in-depth description of the datasets utilized in this project, including details on their sources, structure, and characteristics. Additionally, the chapter will present an analysis of the various data types and their distributions, which is essential for ensuring accurate data preprocessing and effective model development.

3.1 Data Source and Nature

The datasets used in this research are sourced from Kaggle, a widely-used platform that offers open-source datasets [31]. These datasets provided the foundation for model development and testing throughout the research process.

For documentation purposes, publicly available datasets from Kaggle will be used. These datasets, while anonymized and not identical to the proprietary data, closely mirror the structure and features needed for testing the models and documenting the research outcomes. The use of these anonymized public datasets allows for transparency in the research while maintaining the confidentiality of the proprietary data.

3.2 Data Structure and Analysis

This section examines the structure and analysis of the three main datasets used in this research: the *Retail Transactional Dataset*, the *Credit Card Customer Dataset*, and the *Investment Preferences Dataset*. These datasets allowed for comprehensive testing and validation throughout the project.

3.2.1 Structure and Analysis of the Retail Transactional Dataset

The Retail Transactional Dataset [51] provides detailed information about customer purchases and retail transaction behavior. The dataset is structured to give a comprehensive view of customer demographics, transaction specifics, and product details. This data is valuable for analyzing customer purchasing patterns, preferences, and loyalty behavior, making it an ideal source for developing data-driven business strategies.

The dataset contains **302 010 rows** and **30 columns**, each representing an individual customer transaction. The attributes within the dataset include customer demographic details, transaction

history, and product specifics, allowing for a comprehensive analysis of retail operations.

Column	Type
Transaction_ID	numeric
Customer_ID	numeric
Name	string
Email	string
Phone	numeric
Address	string
City	categorical
State	categorical
Zipcode	numeric
Country	categorical
Age	numeric
Gender	categorical
Income	categorical
Customer_Segment	categorical
Date	string
Year	numeric
Month	string
Time	datetime
Total_Purchases	numeric
Amount	numeric
Total_Amount	numeric
Product_Category	string
Product_Brand	string
Product_Type	string
Feedback	categorical
Shipping_Method	categorical
Payment_Method	string
Order_Status	categorical
Ratings	numeric
products	string

Figure 3.1: *Distribution of Data Types in the Retail Transaction Dataset*

The dataset includes the following key categories of information:

Customer Information:

- *Customer ID*: Unique identifier for each customer.
- *Age, Gender, Income*: Demographic attributes useful for understanding the target audience.

Transaction Details:

- *Last Purchase Date, Total Purchases, Amount Spent*: Capturing the transaction history and the overall monetary engagement of customers.

Product Information:

- *Product Category, Product Brand, Product Type*: Information about the products purchased, useful for analyzing product trends and preferences.

Feedback:

- *Customer Feedback*: Ratings or reviews given by customers for the products purchased.

Logistics:

- *Shipping Method, Payment Method, Order Status*: Important for analyzing logistics and payment preferences.

This rich dataset is invaluable for performing data analysis, including customer segmentation, trend analysis, and predictive modeling to enhance business decision-making. The next section will explore specific patterns observed within the data, focusing on customer purchasing behavior and product preferences.

3.2.2 Structure and Analysis of the Credit Card Dataset

The Credit Card Dataset [10] focuses on customer segmentation and is designed to help develop marketing strategies by analyzing the behavior of active credit card holders. This dataset provides valuable insights into customer spending habits, credit usage, and payment behaviors over a six-month period, making it an ideal source for understanding consumer behavior.

The dataset contains **8,950 instances** and includes **18 features** that capture key aspects of customer credit card usage. These features provide detailed behavioral data, allowing businesses to identify patterns and tailor marketing efforts.

Column	Type
CUST_ID	string
BALANCE	numeric
BALANCE_FREQUENCY	numeric
PURCHASES	numeric
ONEOFF_PURCHASES	numeric
INSTALLMENTS_PURCHASES	numeric
CASH_ADVANCE	numeric
PURCHASES_FREQUENCY	numeric
ONEOFF_PURCHASES_FREQUENCY	numeric
PURCHASES_INSTALLMENTS_FREQUENCY	numeric
CASH_ADVANCE_FREQUENCY	numeric
CASH_ADVANCE_TRX	numeric
PURCHASES_TRX	numeric
CREDIT_LIMIT	numeric
PAYMENTS	numeric
MINIMUM_PAYMENTS	numeric
PRC_FULL_PAYMENT	numeric
TENURE	numeric

Figure 3.2: *Distribution of Data Types in the Credit Card Dataset*

Among the most important features of the dataset are:

Customer information:

- *Cust ID*: A unique identifier for each customer, ensuring anonymity.
- *Tenure*: The length of time the customer has held the credit card.

Account balance information:

- *Balance*: The remaining balance available in the customer’s account for further purchases.

Purchasing behavior:

- *Purchases*: The total amount spent by the customer on purchases.
- *Oneoff Purchases*: The highest amount spent by the customer in a single transaction.

Credit Information:

- **Credit Limit**: The maximum credit available to the customer, a key factor in assessing credit risk.
- **Payments**: The total amount of payments made by the customer, indicating the ability to manage debt and maintain credit standing.

These key features provide critical insights for developing models that predict customer behavior, assess credit risk, and design personalized marketing strategies. The attributes selected are central to understanding how customers interact with their credit cards, how much they spend, and their credit management behavior.

3.2.3 Structure and Analysis of the Investment Preferences Dataset

In addition to the two aforementioned datasets, this research also integrates a dataset consisting of **40 records** and **24 columns** that captures various aspects of investment preferences and financial goals of individuals. This dataset provides additional insights into how individuals approach savings and investments across different financial avenues.

The dataset includes the following columns:

Column	Type
gender	string
age	numeric
Investment_Avenues	boolean
Mutual_Funds	numeric
Equity_Market	numeric
Debentures	numeric
Government_Bonds	numeric
Fixed_Deposits	numeric
PPF	numeric
Gold	numeric
Stock_Market	boolean
Factor	string
Objective	string
Purpose	string
Duration	string
Invest_Monitor	string
Expect	string
Avenue	string
What are your savings objectives?	string
Reason_Equity	string
Reason_Mutual	string
Reason_Bonds	string
Reason_FD	string
Source	string

Figure 3.3: *Distribution of Data Types in the Investment Preferences Dataset*

Demographics:

- *gender*: Represents the gender of the respondent.
- *age*: Indicates the age of the respondent.

Investment Preferences:

- *Investment Avenues*: Indicates whether the respondent invests in multiple avenues or not.
- *Mutual Funds, Equity Market, Debentures, Government Bonds, Fixed Deposits, PPF, Gold*: Capture the amounts or proportions invested in each of these respective avenues.
- *Stock Market*: Captures whether the respondent invests in the stock market or not.

Investment Motivations:

- *Factor*: Describes key factors influencing investment decisions.
- *Objective, Purpose*: Describes the investment objectives and purposes.
- *Duration*: Indicates the investment duration.
- *Expect*: Captures the respondent's expectations from the investment.

Additional Insights:

- *Reasons for Investing (Reason_Equity, Reason_Mutual, Reason_Bonds, Reason_FD)*: Describe the reasons behind the respondents' choice of particular investment avenues.
- *Invest Monitor*: Describes how respondents monitor their investments.
- *Savings Objectives*: Captures overall savings goals.
- *Source*: Describes the source of information that guides investment decisions.

This dataset provides a valuable addition to the overall analysis, offering unique insights into individual investment behaviors and preferences. By combining demographic information with detailed investment behavior data, it allows for the development of sophisticated models that can predict investment trends and aid in the design of targeted financial products.

In the following sections, the patterns and insights drawn from this dataset will be explored further, focusing on individual investment behavior and its implications for financial modeling and strategy development.

Chapter 4

Implementation

In this chapter, we will explore the technical implementation of the methodologies and tools used throughout the research project. The following sections detail the utilization of Apache Spark, Databricks, and the custom functions developed to process, clean, and analyze the datasets. These methods play a key role in efficiently managing large-scale data, applying ML techniques, and generating insights.

4.1 Apache Spark

Apache Spark [5] is an open-source framework designed for distributed data processing. Known for its speed and flexibility, Spark enables in-memory computing, significantly accelerating data processing compared to traditional disk-based methods. It supports multiple programming languages, including Python, Scala, Java, and R, making it adaptable for various use cases.

In this project, Spark was instrumental in handling large datasets, facilitating efficient parallel processing across clusters. The following features of Spark were used:

- **Spark SQL:** This library was utilized for structured data processing via SQL queries, streamlining complex data manipulations.
- **MLlib:** The ML library provided scalable implementations of algorithms such as clustering and outlier detection [43].
- **GraphX:** Although not central to this project, GraphX supports graph processing and is beneficial in scenarios involving network data.

The use of Spark enhanced scalability and performance, allowing us to handle millions of data points efficiently.

4.2 Databricks

Databricks [12] is a cloud-based platform built on Apache Spark, designed to simplify big data processing and enable collaborative data science. Its integration with cloud computing makes it an ideal tool for handling large-scale datasets. Key features include:

- **Interactive Notebooks:** Support for multiple languages (Python, SQL, Scala) and real-time collaboration among team members.
- **Cluster Management:** Databricks automatically manages Spark clusters, optimizing resource allocation and task execution.
- **MLflow:** Integrated for managing ML lifecycles, from experimentation to deployment [42].

Databricks significantly streamlined the workflow, enabling easy access to computational resources and seamless cloud integration.

4.3 Functions and Methodologies

To extract meaningful insights from large datasets, a set of custom functions was developed. These functions automate key aspects of the data preprocessing and analysis pipeline, enabling efficient processing and the application of ML techniques. Below, we present an overview of these functions, explaining their objectives and methodologies.

4.3.1 File Handling and Conversion

Data in various formats, such as CSV, JSON, Excel, and Parquet, is common in data science projects. Efficiently handling and converting these files into a standardized format is critical for smooth analysis. The following custom functions streamline the process of identifying file types, converting them into Spark DataFrames, and handling multiple formats in distributed environments like Apache Spark.

4.3.1.1 get_file_type

This function is used to identify the file type based on its extension, helping the system automatically recognize different formats such as CSV, JSON, Excel, and Parquet.

- **Objective:** Automatically determine the file type for easy processing.
- **How It Works:** The function splits the file name to extract its extension and uses a mapping dictionary to return the appropriate `FileType` enum. Unsupported formats raise an error.

```

1 def get_file_type(file_name):
2     extract extension from file_name
3
4     if extension matches known types:
5         return corresponding file_type
6     else:
7         raise error

```

Listing 4.1: Pseudocode for function `get_file_type()`

4.3.1.2 to_spark

This function converts various data formats into a Spark DataFrame. Whether it's a CSV, JSON, Excel file, or Pandas DataFrame, this function ensures that the data is loaded efficiently into Spark for processing.

- **Objective:** Transform different file formats into Spark.
- **How It Works:** Depending on the file type, it uses the appropriate Spark read function to load the data into memory. This function handles common formats like CSV, JSON, Excel, Parquet, and more.

```

1 def to_spark(data, file_type):
2
3     if file_type is CSV:
4         return data as CSV
5     elif file_type is JSON:
6         return data as JSON
7     elif file_type is EXCEL:
8         return data as EXCEL
9     elif file_type is PARQUET:
10        return data as PARQUET
11    elif file_type is ORC:
12        return data as ORC
13    else:
14        return data as SPARK

```

Listing 4.2: Pseudocode for function to_spark()

4.3.2 Data Type Identification and Conversion

In any data analysis pipeline, correctly identifying and converting data types is critical for ensuring accurate processing and analysis. Data often comes in raw formats, with columns incorrectly typed as strings or other general types. This section focuses on two key functions that help automate the identification of column types and the conversion of those columns into appropriate data types.

4.3.2.1 get_type_of_column

This function automatically identifies the type of data contained in a column by analyzing its values. It is particularly useful in determining whether a column should be categorized as numeric, boolean, datetime, or categorical, based on the data's uniqueness and distribution.

- **Objective:** Identify the correct data type of a column in a DataFrame.
- **How It Works:** The function checks the data type of the column, then analyzes the values by calculating metrics like uniqueness ratio (to distinguish between categorical and string types) or directly testing if the data can be interpreted as numeric, boolean, or datetime. Columns are classified based on predefined thresholds and data characteristics.

```

1 def get_type_of_column(column, threshold):
2     if column contains boolean values:
3         return 'boolean'
4     elif column contains numeric values:
5         return 'numeric'
6     elif column contains date or time values:
7         return 'datetime'
8     else:

```

```

9     calculate uniqueness_ratio
10    if uniqueness_ratio < threshold:
11        return 'categorical'
12    else:
13        return 'string'

```

Listing 4.3: Pseudocode for function `get_type_of_column()`

4.3.2.2 `convert_column_types`

This function automatically converts columns of string type into more specific data types (such as boolean, numeric, or datetime) based on their contents.

- **Objective:** Automatically identify and convert string columns into more specific types to improve data accuracy and consistency.
- **How It Works:** It samples unique values from a column and determines if the values represent booleans, numbers, or dates. If identified, the column is cast into the corresponding Spark data type.

```

1  def convert_column_types(dataframe, column):
2
3      if column contains booleans:
4          convert to boolean
5      elif column contains numbers:
6          convert to numeric
7      elif column contains dates:
8          convert to datetime
9
10     return updated dataframe

```

Listing 4.4: Pseudocode for function `convert_column_types()`

4.3.3 Handling Null Values

In real-world datasets, missing values are a common issue that can significantly impact the accuracy of machine learning models. Proper handling of these missing values is essential to ensure robust and reliable results. Various strategies, such as mean, median, or mode imputation, as well as more advanced techniques like forward fill and KNN imputation, can be applied depending on the type of data and the specific problem.

4.3.3.1 `handle_missing_data`

This function handles missing values in a specified column using various imputation strategies like mean, median, mode, or forward fill.

- **Objective:** Fill missing data efficiently using different strategies based on the data type and user preference.
- **How It Works:** The function allows the user to specify a method for imputation, such as using the mean, median, or mode, or applying forward/backward fill.

```
1
2 def handle_missing_data(dataframe, method):
3
4     if method id 'mean':
5         fill missing values with mean
6     elif method id 'median':
7         fill missing values with median
8     elif method id 'mode':
9         fill missing values with mode
10    elif method id 'default':
11        fill missing values with default values
12    elif method id 'forward fill':
13        apply forward fill
14    elif method id 'backward fill':
15        apply backward fill
16    elif method id 'min':
17        fill missing values with min
18    elif method id 'max':
19        fill missing values with max
20    elif method id 'interpolate':
21        fill missing values with interpolate
22    elif method id 'custom':
23        fill missing values with custom values
24    elif method id 'knn':
25        fill missing values with knn
26
27    return updated dataframe
```

Listing 4.5: Pseudocode for function handle_missing_data()

4.3.4 Vectorization

In many datasets, categorical or textual data cannot be directly used by ML algorithms without being transformed into a numerical representation. This process, known as vectorization, converts string or categorical columns into a numerical format, making them compatible with machine learning models. For data that is numeric or boolean, no further transformation is needed since boolean values are converted into integers. However, for non-numeric and non-boolean columns, this vectorization process is crucial.

4.3.4.1 vectorize_column

The vectorize_column function transforms non-numeric and non-boolean data into numerical vectors, which can be used in machine learning models. It supports a variety of vectorization techniques such as one-hot encoding, bag of words, TF-IDF, Word2Vec, and more. This function is specifically useful for categorical or textual data, ensuring it is properly formatted for model training.

- **Objective:** Convert non-numeric or non-boolean data (e.g., text or categorical values) into a numerical format that machine learning models can interpret.

- **How It Works:** Depending on the method chosen (such as one-hot encoding, TF-IDF, Word2Vec, etc.), the function tokenizes and processes the input column, applying the corresponding vectorization algorithm. The output is a DataFrame where the original string or categorical data is replaced with numerical vectors.

```

1 def vectorize_column(dataframe, method):
2
3     if method is 'one-hot encoding':
4         apply one-hot encoding
5     elif method is 'label':
6         apply label transformation
7     elif method is 'binary':
8         apply binary transformation
9     elif method is 'Bag-of-words':
10        apply Bag-of-words transformation
11    elif method is 'TF-IDF':
12        apply TF-IDF transformation
13    elif method is 'Word2Vec':
14        apply Word2Vec transformation
15    elif method is 'Hashing':
16        apply Hashing transformation
17    elif method is 'N-Gram':
18        apply N-Gram transformation
19
20    return vectorized dataframe

```

Listing 4.6: Pseudocode for function `vectorize_column()`

4.3.5 Outlier Detection

Outliers are data points that significantly differ from the majority of the dataset and can potentially distort statistical analyses and ML model performance. Detecting and handling outliers is a crucial step in the data preprocessing pipeline to ensure the accuracy and reliability of the model's predictions. Various methods can be applied to identify outliers based on the nature of the data, ranging from statistical techniques like Z-score to advanced machine learning approaches such as Isolation Forest.

4.3.5.1 detect_outliers

This function detects outliers in a specified column using different methods such as Z-score, IQR, Isolation Forest, DBSCAN, and others.

- **Objective:** Identify and handle outliers in a dataset using various statistical and machine learning-based approaches.
- **How It Works:** Depending on the chosen method, the function applies the corresponding technique to detect outliers. It then filters the data to exclude outliers or marks them for further analysis.

```

1 def detect_outliers(dataframe, method):
2

```

```

3   if method is 'Z-score':
4       calculate Z-scores and filter outliers
5   elif method is 'Box-Plot':
6       calculate Box-Plot and filter outliers
7   elif method is 'IQR':
8       calculate IQR and filter outliers
9   elif method is 'Histogram':
10      calculate Histogram and filter outliers
11  elif method is 'Isolation-Forest':
12      apply Isolation-Forest and filter outliers
13  elif method is 'Lof':
14      calculate LocalOutlierFactor and filter outliers
15  elif method is 'Mahalanobis':
16      calculate Mahalanobis distance and filter outliers
17  elif method is 'CleanLab':
18      apply CleanLab and filter outliers
19  elif method is 'DBSCAN':
20      apply DBSCAN and filter outliers
21  elif method is 'KMeans':
22      apply KMeans and filter outliers
23  elif method is 'MAD':
24      calculate Median Absolute Deviation (MAD) and filter outliers
25  elif method is 'Regression':
26      apply Regression analysis and filter outliers
27
28  return dataframe without outliers

```

Listing 4.7: Pseudocode for function detect_outliers()

4.3.6 Correlation

Correlation analysis helps understanding relationships between variables within a dataset. By measuring the strength and direction of these relationships, we can determine how changes in one variable might influence another. In data science and machine learning, identifying highly correlated features helps in feature selection, reducing multicollinearity, and improving model performance.

4.3.7 Feature Selection

Feature selection is a critical process in machine learning where the most relevant variables are selected from a dataset to improve model performance and reduce computational cost. By eliminating irrelevant or redundant features, we ensure that the model focuses on the most informative variables, thus enhancing its accuracy, reducing overfitting, and speeding up training.

4.3.7.1 select_features

This function performs feature selection by applying various algorithms, including correlation analysis, Chi-Square, Random Forest importance, Lasso regression, and PCA. It helps reduce the feature space, making the dataset more manageable for machine learning tasks.

- **Objective:** Select the most important features from the dataset to improve model performance and reduce dimensionality.

- **How It Works:** The function offers multiple feature selection methods, from filtering based on correlations to using supervised models like Random Forest or PCA. It outputs the selected features for further modeling.

```

1 def select_features(dataframe, method):
2
3     if method is 'correlation':
4         select features based on correlation threshold
5     elif method is 'select_k_best':
6         apply select_k_best method
7     elif method is 'random_forest':
8         apply Random Forest importance for feature selection
9     elif method is 'gbt':
10        apply Gradient Boosted Trees for feature selection
11    elif method is 'lasso':
12        apply Lasso regression for feature selection
13    elif method is 'logistic_regression':
14        apply Logistic Regression for feature selection
15    elif method is 'pca':
16        apply Principal Component Analysis (PCA) for feature selection
17    elif method is 'variance_threshold':
18        select features based on variance threshold
19
20    return selected features

```

Listing 4.8: Pseudocode for function `select_features()`

4.3.8 Clustering

Clustering is a fundamental unsupervised machine learning technique used to group data points into clusters based on their similarity. It helps in discovering hidden patterns, segmenting data into meaningful groups, and providing insights that are crucial for data exploration. In this section, we present the `apply_clustering_methods` function, which allows the application of various clustering algorithms such as K-Means, DBSCAN, Hierarchical Clustering, and others. Each algorithm has its own advantages depending on the dataset and the problem being solved.

However, before applying clustering algorithms like K-Means, one crucial step is determining the optimal number of clusters, k , that best fits the data. This choice can significantly influence the performance and interpretation of the clustering results.

4.3.8.1 find_optimal_k

The `find_optimal_k` function addresses this need by helping to identify the optimal number of clusters, k , using various methods such as the Elbow Method, Silhouette Score, and others. Selecting the appropriate k ensures that the clustering results are meaningful and well-structured, avoiding underfitting or overfitting the data.

- **Objective:** Find the optimal number of clusters (k) for clustering algorithms.
- **How It Works:** The function evaluates multiple values of k using different metrics. It computes the optimal k by iterating over a range of cluster numbers and applying one of several available

methods. The function then returns the k that optimizes the selected metric, ensuring better clustering quality.

```

1 def find_optimal_k(df, feature_columns, k_min, k_max, method):
2
3     assembler = data to vector
4     VectorAssembler para criar coluna de features assembler =
5
6     k_values = list of k values from k_min to k_max
7
8
9     if method is 'elbow':
10        for each k in k_values:
11            apply K-Means algorithm with k clusters
12            calculate WSSSE (training cost)
13
14        optimal_k = k value corresponding to the lowest WSSSE
15
16    elif method is 'silhouette':
17        for each k in k_values:
18            apply K-Means algorithm with k clusters
19            calculate Silhouette score
20
21        optimal_k = k value corresponding to the highest Silhouette score
22
23    elif method is 'calinski-harabasz':
24        for each k in k_values:
25            apply K-Means algorithm with k clusters
26            calculate Calinski-Harabasz score
27
28        optimal_k = k value corresponding to the highest Calinski-Harabasz score
29
30    elif method is 'gap':
31        for each k in k_values:
32            calculate Gap statistic and standard deviation for k
33
34        optimal_k = k value corresponding to the largest Gap statistic
35
36    elif method is 'canopy':
37        for each k in k_values:
38            apply Bisecting K-Means with k clusters
39            calculate training cost for Canopy method
40
41        optimal_k = k value corresponding to the lowest training cost
42
43    elif method is 'davies-bouldin':
44        for each k in k_values:
45            apply K-Means algorithm with k clusters
46            generate cluster predictions
47            calculate cluster centers
48            compute Davies-Bouldin Index
49
50        optimal_k = k value corresponding to the lowest Davies-Bouldin Index
51
52

```

```
53 return optimal_k
```

Listing 4.9: Pseudocode for function find_optimal_k()

4.3.8.2 apply_clustering_methods

This function applies a range of clustering methods to group data points in a dataset. Depending on the method chosen, the function adjusts its parameters to fit the best possible model for clustering.

- **Objective:** Group data points into clusters using various clustering algorithms.
- **How It Works:** The function supports multiple clustering algorithms like K-Means, DBSCAN, Agglomerative Clustering, and others. It allows tuning of parameters such as the number of clusters (k for K-Means) or the neighborhood distance (eps for DBSCAN). After clustering, it outputs the cluster labels for each data point.

```
1 def apply_clustering_methods(dataframe, method):
2
3     if method is 'K-Means':
4         apply K-Means clustering
5     elif method is 'Affinity':
6         apply Affinity Propagation clustering
7     elif method is 'Mean_Shift':
8         apply Mean Shift clustering
9     elif method is 'Spectral':
10        apply Spectral clustering
11    elif method is 'Hierarchical':
12        apply Hierarchical clustering
13    elif method is 'DBSCAN':
14        apply DBSCAN clustering
15    elif method is 'HDBSCAN':
16        apply HDBSCAN clustering
17    elif method is 'OPTICS':
18        apply OPTICS clustering
19    elif method is 'BIRCH':
20        apply BIRCH clustering
21
22    return cluster labels
```

Listing 4.10: Pseudocode for function apply_clustering_methods()

4.3.9 Visualization Methods

Visualization is an integral part of analyzing and understanding the dataset. Several functions were developed for generating different types of visualizations, including histograms, box plots, and bar plots.

This function plots a histogram to visualize the distribution of numerical data.

Visualization is an essential aspect of data analysis and machine learning. It helps to understand the structure of the data, explore relationships, and present results in a way that is interpretable and actionable. In this section, we focus on various visualization techniques used for analyzing data distributions, visualizing clustering results, and presenting the outcomes of topic modeling via BERTopic.

4.3.9.1 Data Visualization for Analysis

Visualizing data distributions is an essential step in understanding the nature of numerical, categorical, and boolean attributes within the dataset. At this stage of exploratory data analysis, charts such as histograms, box plots, and bar plots are used to identify patterns, outliers, and trends within the data.

4.3.9.2 Clustering Visualization

After performing clustering, visualizing the results helps to better understand the groupings of data points. Scatter plots are commonly used for 2D representations of clusters, allowing for a visual interpretation of how well the clustering algorithm has segmented the data.

4.3.9.3 BERTopic Visualization

For topic modeling using BERTopic, visualizations are crucial for understanding the identified topics and their relationships. Graphs such as bar charts display the number of documents per topic, while 2D scatter plots show how topics are distributed across the dataset, often using dimensionality reduction techniques like PCA.

Chapter 5

Testing and Validation

The purpose of this chapter is to rigorously validate the performance and effectiveness of the proposed clustering algorithms by conducting a comprehensive series of tests. These tests were designed to simulate diverse and intricate real-world scenarios, reflecting the complexities and challenges that clustering algorithms may encounter in practical applications. Although only a subset of tests is documented here, the complete testing scope encompassed a broad range of experiments, probing specific algorithmic behaviors across different configurations, including hyperparameter tuning, dataset setups, preprocessing techniques, and clustering scenarios.

5.1 Testing Scope and Challenges

One of the primary challenges during the testing phase was managing the significant computational cost required to execute an extensive series of tests on large and complex datasets. The combination of high data volume and computationally intensive algorithms necessitated strategic planning, including prioritization of test configurations that provided the most meaningful insights into the algorithms' behavior.

Another critical decision was to avoid modifications to the datasets, such as sampling or dimensionality reduction, to preserve their intrinsic structure. Any such alterations could have introduced biases or obscured important patterns, potentially distorting the clustering results. This ensured that the clustering outcomes faithfully represented the inherent complexity of the data, minimizing the risk of distortions or oversights.

To validate the results under varying conditions, cross-validation strategies were employed, such as testing different data splits, initialization parameters, and synthetic datasets. These efforts ensured that the findings were robust and not overly sensitive to specific configurations or random initialization.

5.2 Data Preprocessing

Data preprocessing is pivotal in preparing data for clustering algorithms, as the quality of input data directly influences algorithm performance and reliability. This phase involved critical tasks such as addressing missing data, detecting and removing outliers, normalizing features, and selecting relevant variables for clustering. Each task aimed to enhance data quality, providing the algorithms with optimal inputs. For illustration, the *Credit Card Customer* dataset is used to detail preprocessing steps, with

similar methods applied across all datasets to maintain consistency.

5.2.1 Handling Missing Data

An initial evaluation identified **314 missing values** across **8,950 entries and 18 columns**. To address this, multiple strategies were tested, including data removal, imputation techniques (mean, median, and mode), and predictive modeling. Through these methods, all missing values were effectively resolved, leaving a complete and consistent dataset.

- **Before processing:** 314 missing values.
- **After processing:** 0 missing values.

The comprehensive resolution of missing data issues resulted in a complete and robust dataset, providing a solid foundation for accurate and reliable clustering analysis.

5.2.2 Outlier Detection and Removal

Outliers, or anomalous data points, are values that deviate significantly from the expected range of a dataset. Outliers can significantly distort clustering outcomes, especially in distance-based algorithms like K-means, where extreme values can shift centroids and obscure meaningful patterns. By employing multiple detection techniques, including Box Plot and Interquartile Range (IQR), we ensured that the impact of these anomalies was minimized, resulting in more accurate and interpretable clusters.

In this section, the results of two specific methods—Box Plot and Interquartile Range (IQR)—are presented due to their straightforward visual clarity. These methods were selected to illustrate the outlier removal process effectively.

5.2.2.1 Box Plot Method

The Box Plot method provides a clear visualization of outliers and allows for quick identification of extreme values in individual variables. Figure 5.1 shows the Box Plots for the "BALANCE" and "PURCHASES" variables before and after outlier removal. The initial box plots (top row) reveal a substantial number of outliers represented as points outside the whiskers, indicating extreme values that could have impacted clustering accuracy.

After the removal of these outliers, the updated box plots (bottom row) show a more refined distribution, confirming the effective elimination of extreme points. This process allowed the dataset to reflect a more accurate and representative range of values, enhancing the clustering quality by reducing the influence of atypical data points.

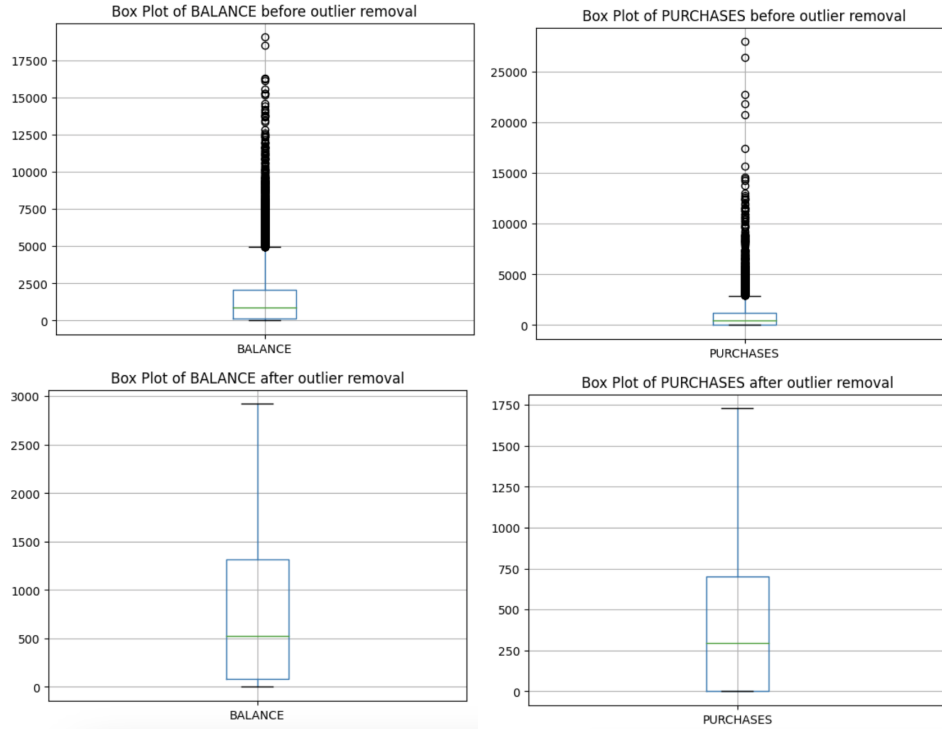


Figure 5.1: *Outlier Detection and Removal Process using Box Plot Example*

5.2.2.2 IQR Method

In addition to the Box Plot approach, the Interquartile Range (IQR) method was also applied to detect and remove outliers. This method focuses on values that fall beyond a specified range (usually 1.5 times the interquartile range), which helps in isolating extreme data points. Figure 5.2 shows histograms of the "BALANCE" and "PURCHASES" variables before and after applying the IQR-based outlier removal.

Before IQR Outlier Removal: The top histogram in each pair illustrates the original distribution, with many extreme values contributing to a skewed appearance. **After IQR Outlier Removal:** The bottom histogram shows the distribution after outlier removal, demonstrating a significant reduction in skewness and a more balanced frequency of values. This refined data distribution helped improve the clustering results by mitigating the impact of extreme, non-representative values.

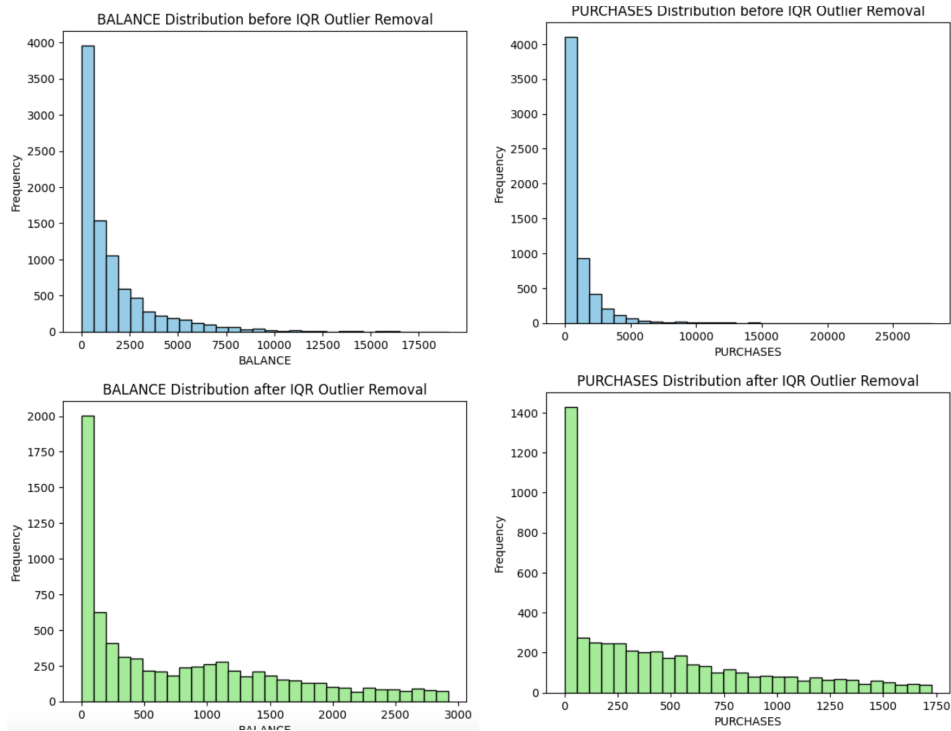


Figure 5.2: *Histograms of BALANCE and PURCHASES Variables Before and After IQR-Based Outlier Removal*

5.2.2.3 Additional Methods

Various outlier detection methods were tested, including Z-score analysis, isolation forests, and clustering-based detection. However, the Box Plot and IQR methods were chosen for representation here due to their straightforward visualization and ease of interpretation. These methods effectively illustrate the impact of outlier removal, making them the most suitable for conveying the results in a clear and accessible manner.

5.3 Feature Selection

Once the dataset was cleaned, feature normalization was applied to ensure consistency in the scale of all variables. This step is crucial, particularly for distance-based clustering algorithms like K-means, where differences in feature magnitudes can skew the clustering process.

For feature selection, we explored multiple techniques, notably *SelectKBest* and *Variance Threshold*, which each offer unique advantages. *SelectKBest* focuses on selecting the highest scoring features based on statistical tests, while *Variance Threshold* helps by eliminating features with low variance, which contribute little to distinguishing data points.

Interestingly, despite the different criteria, both methods consistently identified similar subsets of relevant features, which strengthened our confidence in their selection. This consistency suggests that these features are robust and consistently contribute to improving the clustering outcomes. To further emphasize this, Table 5.1 presents a comparison between the initial set of features and the final selected subset.

Feature	Method Variance	Method SelectKBest
CUST_ID	-	-
BALANCE	✓	✓
BALANCE_FREQUENCY	✓	✓
PURCHASES	✓	✓
ONEOFF_PURCHASES	✓	✓
INSTALLMENTS_PURCHASES	✓	-
CASH_ADVANCE	✓	✓
PURCHASES_FREQUENCY	✓	✓
ONEOFF_PURCHASES_FREQUENCY	✓	✓
PURCHASES_INSTALLMENTS_FREQUENCY	✓	-
CASH_ADVANCE_FREQUENCY	✓	✓
CASH_ADVANCE_TRX	✓	-
PURCHASES_TRX	✓	✓
CREDIT_LIMIT	✓	✓
PAYMENTS	✓	✓
MINIMUM_PAYMENTS	✓	✓
PRC_FULL_PAYMENT	✓	-
TENURE	-	-

Table 5.1: Comparison of Selected Features Between Variance Threshold and SelectKBest Methods

To illustrate how different methods impacted feature selection:

- **Correlation-Based Selection** resulted in a set of 16 features, including notable variables like "BALANCE", "PURCHASES", and "MINIMUM_PAYMENTS", while excluding others like "TENURE" and "CUST_ID" due to their lower relevance.
- **SelectKBest Analysis** focused on identifying the top contributing features for clustering, leading to a smaller subset of 12 features, including key financial metrics like "BALANCE", "PURCHASES", and "CREDIT_LIMIT". This reduction further emphasized variables that are critical to understanding customer behavior, thereby simplifying the dataset without significant loss of information.

Normalization was performed using the Min-Max scaling method, which rescaled all features to a range between 0 and 1. This ensures that no single feature disproportionately affects the clustering due to its magnitude.

By focusing on these selected subsets of features, we were able to optimize the clustering process by ensuring that only the most impactful variables were included while minimizing noise from less relevant data. This careful selection of features allowed for improved computational efficiency and enhanced interpretability of the clustering results.

The results obtained using other feature selection methods were largely similar to those achieved with the Variance Threshold and SelectKBest approaches. This consistency across different methods gave us further confidence in the robustness of the selected features.

Ultimately, we chose to proceed with the features selected by the Variance Threshold method, as it allowed us to maintain a high level of data integrity while ensuring we did not lose significant

information. This approach provided a balanced trade-off, retaining the key features necessary for effective clustering without unnecessarily complicating the model with excessive or redundant variables.

5.4 Cluster Analysis Using Various Methods

This section delves into the application of multiple clustering algorithms to segment customers into meaningful groups based on their transactional and financial behaviors. By employing a range of clustering methods, we aimed to ensure robust and comprehensive insights into the data, uncovering distinct patterns and actionable segments.

A critical aspect of clustering analysis is determining the optimal parameters for each method, such as the number of clusters (K). This ensures that the results are both interpretable and reflective of the inherent data structure. For each method, we evaluated its effectiveness in capturing meaningful clusters using relevant metrics and visualizations.

5.4.1 Optimal Parameter Determination

Determining the optimal number of clusters, K , is a crucial step in clustering analysis as it directly impacts the interpretability and quality of the resulting clusters. For this analysis, we applied the Elbow Method to observe where the Within-Cluster Sum of Squared Errors (WSSE) begins to level off, indicating a suitable balance between the number of clusters and clustering quality.

As shown in Figure 5.3, the "elbow" point appears to be between $K = 4$ and $K = 5$, as both values lead to a significant reduction in WSSE. While $K = 5$ offers a slightly lower WSSE, the reduction is marginal when compared to $K = 4$. Choosing $K = 4$ thus provides a simpler model that captures most of the data's inherent structure without adding unnecessary complexity.

To confirm the selection of $K = 4$ as the optimal choice, additional metrics such as the Silhouette Score and Gap Statistic were considered as secondary checks. These metrics supported $K = 4$ as an effective compromise between interpretability and clustering quality. Therefore, based on this analysis, $K = 4$ was selected for the final clustering model.

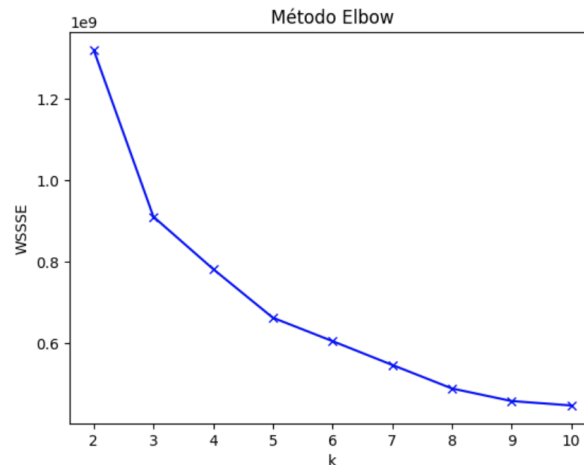


Figure 5.3: *Optimal K Determination using the Elbow Method*

5.4.2 Clustering Methods and Results

The objective of this analysis is to segment credit card customers based on their spending behavior and credit usage, identifying groups with similar characteristics in an unsupervised manner. This segmentation aims to provide valuable insights to support marketing strategies, risk assessment, and customer relationship management.

To achieve this, we applied various clustering algorithms, exploring different approaches to capture distinct patterns in the data. The methods include centroid-based techniques like K-means, as well as hierarchical and density-based algorithms, ensuring a comprehensive analysis.

Each method's effectiveness was evaluated based on its ability to segment customers into meaningful groups with similar behaviors. Key features such as **PURCHASES**, **ONEOFF_PURCHASES**, and **PAYMENTS** were selected for their relevance in representing customer financial behavior. The following sections detail the results of each method, accompanied by visualizations that illustrate the identified clusters.

5.4.2.1 K-means Clustering Results

Using K-means with $K = 4$, the algorithm generated four distinct clusters. The clusters reveal patterns in customer spending and payment behavior, offering actionable insights into different customer profiles. Figure 5.4 illustrates these clusters across key financial metrics.

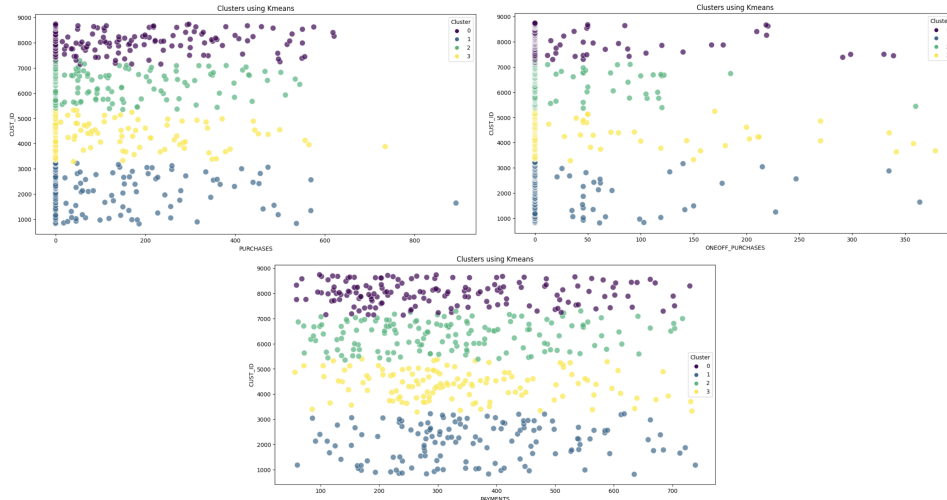
The scatter plots below highlight the key dimensions of the analysis:

- **Purchases:** The first plot demonstrates the total purchase behavior across clusters. Customers are grouped into segments with significantly high purchase volumes and others with lower spending tendencies.
- **One-off Purchases:** The second plot focuses on one-off purchasing patterns, revealing clear distinctions between clusters with frequent, high-value transactions and those with fewer or smaller one-off purchases.
- **Payments:** The third plot captures payment behavior, with clusters highlighting customers who exhibit regular payment patterns compared to those with sporadic or lower payments.

These visualizations emphasize the separation between clusters, validating the effectiveness of the K-means algorithm in identifying meaningful customer segments. For instance:

- One cluster groups high-value customers with frequent purchases and consistent payments, indicating a premium segment.
- Another cluster represents customers with minimal spending and sporadic payments, suggesting a low-risk but low-engagement group.

The results underscore how K-means can assist in tailoring financial services, such as personalized offers for high-value customers or strategies to engage low-spending customers.

Figure 5.4: *Clusters generated by K-means*

5.4.2.2 Mean Shift Clustering Results

The Mean Shift algorithm was used to identify clusters by locating the maxima of a density function, allowing the discovery of customer groups without needing to specify the number of clusters in advance. This approach is particularly useful when patterns are not clearly defined, and the number of clusters is unknown.

The Mean Shift algorithm was applied to segment customers based on behavioral patterns without predefining the number of clusters. One critical parameter in Mean Shift is the bandwidth, which controls the smoothness of the clustering process. In this analysis, the bandwidth was determined using a quantile-based approach, and the relationship between the quantile value and the resulting number of clusters is illustrated in Figure 5.5.

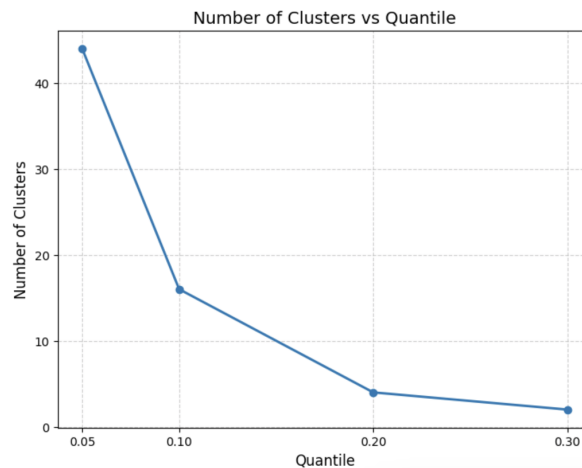
Figure 5.5: *Quantile vs. Number of Clusters*

Figure 5.5 demonstrates how the quantile parameter influences the number of clusters generated:

- **Lower Quantile Values:** At smaller quantiles (e.g., 0.05), the algorithm detects a higher number of clusters. This is because smaller bandwidths focus on finer details, splitting the data into more

granular groups.

- **Higher Quantile Values:** As the quantile increases (e.g., 0.2 or 0.3), the number of clusters decreases. This occurs because larger bandwidths merge nearby density peaks, resulting in broader and fewer clusters.

For this analysis, a quantile value of 0.2 was selected. This choice represents a balance between granularity and generalization, yielding a manageable number of clusters that preserve meaningful customer segmentation while avoiding over-segmentation.

This relationship highlights the trade-off between **precision** and **generalization** in the clustering process:

- Lower quantiles can capture subtle variations in customer behavior, which is useful for detailed segmentation.
- Higher quantiles yield broader, more general groupings, which are ideal for identifying overarching patterns.

Once the quantile value of 0.2 was applied, the algorithm identified clusters based on the relationships between key features. The results are visualized in Figure 5.6.

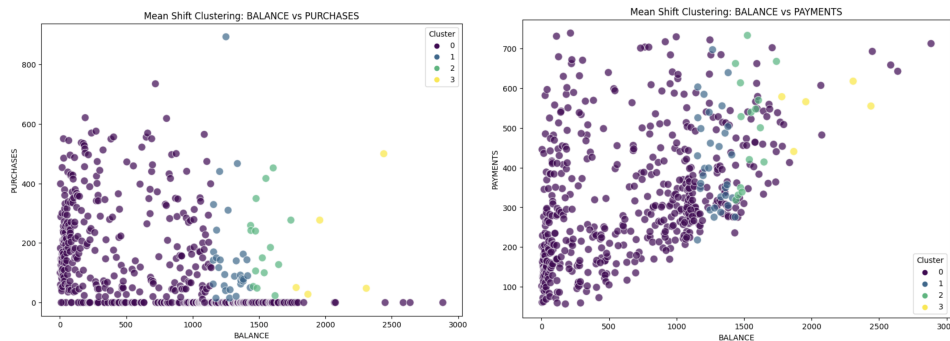


Figure 5.6: *Clusters from the Mean Shift algorithm*

1. Clustering of **BALANCE** vs. **PAYMENTS**:

- Most customers are concentrated in a large cluster with low balances and payments (Cluster 0, purple). This represents a dominant group of low-risk, low-engagement customers.
- Smaller clusters (e.g., green and yellow) highlight high-balance and high-payment customers, reflecting premium or active clients.

2. Clustering of **BALANCE** vs. **PURCHASES**:

- Similarly, the dominant purple cluster shows customers with low balances and minimal purchases.
- Distinct outlier clusters (e.g., yellow and green) capture customers with higher balances and purchase activity, emphasizing unique spending profiles.

5.4.2.3 Birch Clustering Results

The Birch (Balanced Iterative Reducing and Clustering using Hierarchies) algorithm was employed to segment customers into groups based on their spending and credit usage behavior. This method is particularly well-suited for large datasets and efficiently handles hierarchical tendencies. For this analysis, the algorithm was configured to produce **4 clusters**, aligning with the number of clusters used in the K-means analysis.



Figure 5.7: *Clusters from the Birch algorithm*

The scatter plot in Figure 5.7 illustrates the clusters identified by Birch, showing the relationship between **BALANCE** and **PURCHASES**. The results reveal four distinct clusters, each representing unique customer segments. The largest cluster, Cluster 0 (purple), corresponds to customers with low balances and minimal purchasing activity, forming the densest and most prominent group. This cluster likely represents the majority of the customer base with limited financial engagement.

Cluster 1 (blue) captures customers with moderate balances and purchasing behavior. This group includes individuals whose spending patterns suggest a balanced, mid-tier customer segment. In contrast, Cluster 2 (green) represents customers with slightly higher balances but moderate purchasing activity, suggesting a group with greater financial potential that could be targeted with specialized offers. Lastly, Cluster 3 (yellow) stands out as a smaller segment of high-value customers with both high balances and purchasing levels. This cluster reflects a premium customer group characterized by significant financial activity and engagement.

The hierarchical nature of Birch clustering ensures that the identified clusters reflect natural groupings within the data, capturing not only the majority behavior but also smaller, distinct segments. The presence of these smaller clusters, particularly Clusters 2 and 3, highlights opportunities to tailor marketing strategies to premium customers while also identifying low-engagement groups that could benefit from reactivation efforts.

5.4.2.4 Agglomerative Clustering Results

Agglomerative Clustering, a hierarchical clustering method, was used to group customers into distinct clusters based on their financial behavior. This method constructs a hierarchy of nested clusters, providing insight into relationships within the data.

The dendrogram in Figure 5.8 illustrates the hierarchical structure of the data. Each vertical line represents a cluster merging step, while the height of the merge indicates the distance or dissimilarity between clusters. By cutting the dendrogram at a threshold of 25 (indicated by the red dashed line), we identified 4 clusters, as shown by the distinct colored branches below the threshold. This threshold was selected to balance granularity and interpretability, ensuring the most meaningful customer groupings.

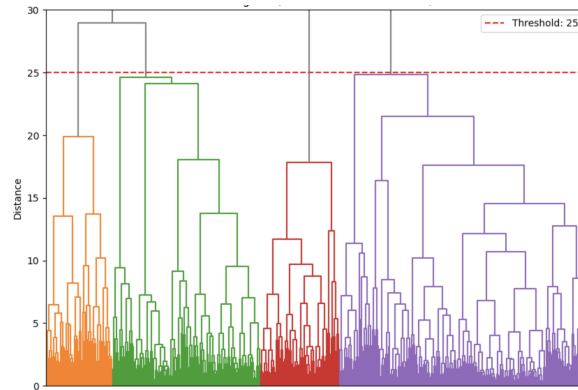


Figure 5.8: *Dendrogram Showing Hierarchical Clustering and Optimal Number of Clusters*

This approach effectively highlights the clusters, with the branches below the threshold visually grouped into four distinct colors, representing the identified customer segments.

The identified clusters were visualized using scatter plots across different feature combinations, highlighting customer behavior:

1. BALANCE vs. PURCHASES:

This plot reveals four distinct clusters.

- **Cluster 0 (purple):** Represents customers with low balances and low purchasing activity, forming the largest and most densely populated group.
- **Cluster 1 (blue):** Corresponds to customers with moderate balances and purchases, indicating a balanced spending profile.
- **Cluster 2 (green):** Captures customers with higher balances but moderate purchasing activity, suggesting an intermediate level of financial engagement.
- **Cluster 3 (yellow):** Highlights a smaller group of customers with high balances and high purchases, representing a premium customer segment.

2. INSTALLMENTS_PURCHASES vs. MINIMUM_PAYMENTS:

The scatter plot shows how the clusters differ in terms of installment purchases and minimum payments.

- **Cluster 0 (purple):** Includes customers with consistently low minimum payments and installment purchases.
- **Cluster 1 (blue):** Features customers with slightly higher installment purchases and moderate minimum payments.

- **Cluster 2 (green):** Highlights customers with a more diverse range of installment purchases and payment behaviors.
- **Cluster 3 (yellow):** Represents a smaller group with high installment purchases and relatively higher minimum payments, further reinforcing their premium status.

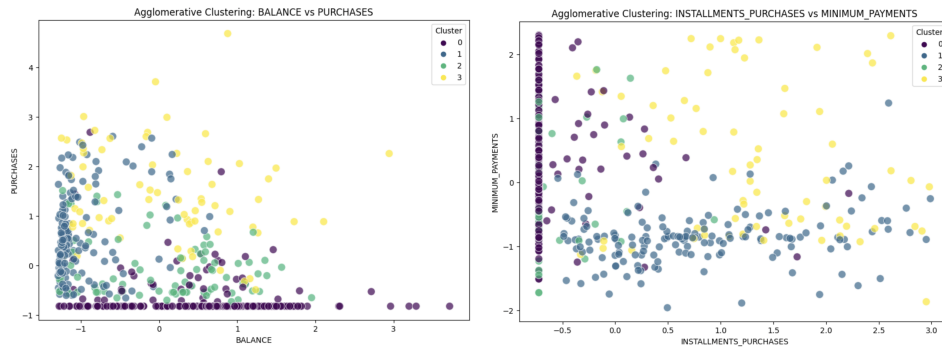


Figure 5.9: *Clusters from the Agglomerative Clustering algorithm*

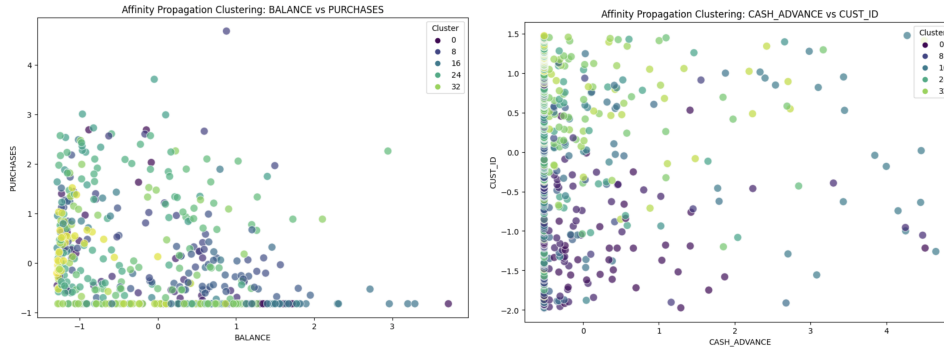
5.4.2.5 Affinity Propagation Clustering Results

Affinity Propagation Clustering was applied to segment customers dynamically, leveraging its ability to determine the number of clusters based on data similarity without predefined input. The results of the clustering process revealed a meaningful segmentation of the dataset, capturing diverse customer behaviors.

The clustering process was influenced by key parameters, particularly the damping factor and preference values, which determine how the algorithm balances stability and cluster formation. Based on the sensitivity analysis, the results showed that Affinity Propagation identified multiple clusters that align well with the underlying data structure.

The clusters identified by Affinity Propagation reveal meaningful distinctions in customer behavior. The plot of **BALANCE** versus **PURCHASES** shows that customers with low balances and purchases dominate the dataset, forming the majority of the clusters. This segment represents low-engagement customers, likely requiring reactivation strategies or targeted campaigns to increase activity. On the other hand, distinct clusters emerge for customers with higher balances and purchasing activity, highlighting high-value segments with significant financial engagement. These clusters could be targeted with tailored offers, premium services, or loyalty programs.

When analyzing the relationship between **CASH_ADVANCE** and **CUST_ID**, the clustering reveals diverse usage patterns. Most clusters consist of customers with low reliance on cash advances, indicating more traditional financial behavior. However, some clusters represent outliers who frequently rely on cash advances. These customers may depend on revolving credit or require short-term liquidity solutions, presenting an opportunity to design specific products or financial support services to meet their needs.

Figure 5.10: *Affinity Propagation Clustering*

5.5 BERTopic Analysis and Results

The integration of BERTopic within this analytics pipeline yielded significant insights, particularly in the realm of large-scale textual data analysis. By leveraging its state-of-the-art topic modeling capabilities, BERTopic proved instrumental in uncovering latent patterns, clustering information, and transforming unstructured data into actionable insights. This analysis focused on datasets detailed in earlier chapters, including the Retail Transactional Dataset, characterized by a significant volume of string-based information ideally suited for this approach. The following section explores the results in depth, highlighting their relevance, practical interpretation, and potential implications for achieving the project’s overarching objectives.

5.5.1 Results and Insights

5.5.1.1 Customer Names and Demographics

The analysis of the "Name" column identified 6,940 unique topics, revealing patterns in common names and their occurrences. Clusters such as “Tiffany James” and “Donna Thomas” highlighted frequently used names, reflecting cultural or geographic naming conventions. However, some clusters lacked coherence, including empty or mixed entries grouped under a noise topic (-1), suggesting potential preprocessing gaps. These results provide actionable insights into customer profiling, allowing businesses to design culturally tailored outreach strategies.

5.5.1.2 Email Patterns

From the "Email" column, BERTopic identified over 10,000 topics, effectively clustering emails by structural patterns and naming conventions. Recurring themes included personal and organizational domains, such as “Jacob54@gmail.com,” providing insights into customer affiliations and preferences. Sparse or placeholder data formed distinct noise topics, indicating opportunities for data refinement. These patterns are highly useful for segmenting customers and optimizing targeted email marketing campaigns.

5.5.1.3 Address-Based Clustering

The clustering of address data revealed 2,614 topics, organized around recurring terms such as “Johnson Lock” and “Michael Squares.” This structured grouping uncovered geographic trends and regional address conventions. A significant proportion of data grouped under the noise topic suggested that

some entries lacked standardization. Nonetheless, the results provide a strong foundation for regional market analysis and logistical planning.

5.5.1.4 Product Categories and Brands

In the analysis of product-related data, BERTopic extracted meaningful topics across both product categories and brands. The "Product Category" column revealed dominant clusters like "Home Decor," "Clothing," and "Books," while brand analysis highlighted key players such as "Samsung," "IKEA," and "Pepsi." These results validate the model's ability to identify market trends and customer preferences. Sparse or unclassified data in both columns formed noise topics, but the overall findings align with known consumer behaviors, enabling better inventory management and brand-specific promotions.

5.5.1.5 Feedback and Sentiment Analysis

BERTopic's clustering of feedback data offered valuable insights into customer sentiment. The analysis identified clusters corresponding to common ratings, such as "Excellent," "Good," and "Average," highlighting areas of satisfaction and dissatisfaction. Clusters associated with "Bad" feedback provided actionable information for targeted quality improvements. This clustering demonstrates how BERTopic can effectively capture sentiment trends, enabling businesses to prioritize enhancements that align with customer expectations.

5.5.1.6 Behavioral Insights: Payment and Order Patterns

The analysis of behavioral data, including "Payment Method" and "Order Status," uncovered meaningful trends. Payment data clustered around methods like "Debit Card," "Credit Card," and "PayPal," reflecting consumer preferences for digital transactions. Similarly, order statuses such as "Pending," "Shipped," and "Delivered" formed distinct and interpretable clusters, offering insights into operational performance. These findings allow businesses to streamline payment options and optimize order fulfillment workflows.

5.5.1.7 Temporal and Seasonal Trends

The clustering of "Date" and "Month" data revealed significant seasonal patterns. For example, clusters associated with specific months, such as "December" or "July," suggested peaks in activity or sales. This information is critical for planning promotions and resource allocation during high-demand periods.

5.5.2 Interpretation of Results

The results obtained through BERTopic highlight its robustness in clustering complex datasets and revealing non-obvious relationships. The coherent clusters across customer, product, and behavioral data validate the model's ability to uncover actionable insights. However, the presence of noise topics in certain columns emphasizes the importance of data preprocessing and standardization.

The interpretability of BERTopic's clusters, coupled with the interactive visualizations such as word clouds and scatter plots, makes it an accessible tool for stakeholders across technical and non-technical backgrounds. These visualizations aid in identifying key trends and relationships, such as customer preferences, regional dynamics, and seasonal behaviors.

The insights derived from BERTopic offer a multitude of practical applications:

- **Customer Segmentation:** Clusters of names, emails, and feedback enable the creation of detailed customer personas for targeted marketing.
- **Operational Optimization:** Order and payment data clusters highlight potential bottlenecks, allowing for improved resource allocation.
- **Market Trends:** Product and brand clustering provide valuable information for inventory management and competitive positioning.
- **Sentiment Analysis:** Feedback clusters offer a direct line to customer satisfaction metrics, guiding quality improvements.

By integrating these insights into strategic decision-making processes, businesses can achieve enhanced efficiency, customer satisfaction, and market responsiveness.

5.6 Summary

This chapter presented a comprehensive validation of clustering algorithms, ensuring their robustness and effectiveness in segmenting customers based on transactional and financial behaviors. The validation process was structured across three major phases: **data preprocessing**, **feature selection**, and **clustering analysis**, each contributing to the overall quality and interpretability of the results.

Data Preprocessing played a pivotal role in ensuring the integrity and reliability of the datasets. Missing data was meticulously addressed using imputation techniques, while outlier detection and removal were executed through Box Plot and IQR methods. These steps refined the dataset, reducing the impact of anomalies and ensuring that the results reflected genuine patterns inherent in the data.

In the **feature selection** phase, various methods, including Variance Threshold and SelectKBest, were employed to isolate the most relevant variables. This process resulted in a streamlined dataset with minimal noise, enabling the clustering algorithms to focus on key features like `{BALANCE}`, `{PURCHASES}`, and `{MINIMUM_PAYMENTS}`. Normalization ensured consistency in variable scales, further enhancing algorithm performance and result interpretability.

The **clustering analysis** demonstrated the diverse strengths of each algorithm:

- **K-means** provided clear and interpretable clusters, highlighting high-value customers with significant spending and payment activity, as well as low-risk, low-engagement groups.
- **Mean Shift** excelled in adapting to data density, uncovering natural groupings without requiring a predefined cluster count. Its flexibility made it ideal for capturing nuanced patterns in customer behavior.
- **Birch** showcased its efficiency in handling large datasets, clustering hierarchically and highlighting both dominant and smaller customer groups, such as premium segments.
- **Agglomerative Clustering** revealed nested relationships among customers, enabling insights into hierarchical groupings that could inform personalized marketing and loyalty strategies.
- **Affinity Propagation** dynamically segmented customers by identifying exemplars, effectively capturing diverse behaviors such as high reliance on cash advances or consistent spending patterns.

Each method provided unique insights into customer behavior, from low-engagement users requiring reactivation strategies to high-value customers suited for tailored offers and loyalty programs. These findings highlight the algorithms' adaptability to real-world data complexities and their potential for applications in customer relationship management, targeted financial services, and risk assessment.

In conclusion, the rigorous testing and validation process confirmed the algorithms' effectiveness in addressing the challenges posed by large-scale and complex datasets. The analysis not only revealed meaningful customer segments but also demonstrated the importance of methodical preprocessing and feature selection in achieving accurate clustering results. Future work could extend this analysis to dynamic or time-series data, enabling the exploration of evolving behavioral patterns over time and further enhancing the applicability of clustering for predictive and real-time analytics.

Chapter 6

Conclusions and Future Work

The development of this AI-powered analytics pipeline marks a substantial step forward in processing and extracting insights from large-scale and unstructured data. By integrating ML models, NLP, and advanced clustering techniques, this project has demonstrated how businesses can efficiently transform raw data into actionable intelligence. Despite its achievements, this work is far from complete. It emphasizes the need for continuous learning, refinement, and collaborative efforts to address challenges and ensure sustainable progress.

One of the standout contributions of this pipeline is the implementation of the BERTopic model for topic modeling. This approach effectively extracted meaningful topics from extensive unstructured textual data, allowing the identification of trends in customer behavior, risk management, and anomaly detection. The pipeline successfully demonstrated the potential of clustering techniques to reveal deeper insights into structured and unstructured data types. However, we acknowledge that these results represent just an initial exploration of the vast possibilities in this domain.

6.1 Key Insights and Challenges

A major insight from this project was the critical role of data preprocessing, feature selection, and algorithm validation in achieving meaningful clustering results. Metrics such as the Silhouette Score and Davies-Bouldin Index proved invaluable in quantifying the quality of the clusters, while preprocessing techniques, including outlier removal and normalization, directly impacted the accuracy and interpretability of the results.

Nonetheless, significant challenges were encountered, particularly with the computational demands of clustering large datasets. These limitations necessitated strategic compromises, such as focusing on smaller datasets during validation to allow for in-depth analysis under manageable conditions. While this approach enabled a clearer understanding of the system's functionality, it also underscored the need for future work to optimize scalability and adapt to increasingly complex real-world datasets.

Most importantly, we recognize that the insights derived from this project are just one step in a broader journey. The clustering results, while promising, reveal opportunities for refinement and optimization, particularly as the pipeline evolves to address larger, more diverse datasets.

6.2 Future Work

Building on the findings and challenges of this project, several avenues for further exploration and improvement emerge. These include both technical advancements and practical applications:

1. **Enhanced Computational Efficiency:**

Future efforts should prioritize optimizing computational efficiency to better handle large-scale datasets. Investigating alternative data processing frameworks, leveraging distributed computing environments, and refining outlier handling methods could significantly reduce processing times. This will be crucial for scaling the pipeline while maintaining performance.

2. **Improving Clustering Accuracy and Scalability:**

While the clustering algorithms demonstrated strong performance, further study is needed to refine their accuracy and scalability. Advanced approaches, such as deep clustering, adaptive hierarchical clustering, or ensemble methods, should be explored. Additionally, improving feature selection through automated techniques or domain-specific knowledge could enhance segmentation quality in diverse datasets.

3. **Integration of Temporal and Real-Time Data:**

Expanding the pipeline to handle time-series data and real-time processing would enable the analysis of dynamic trends, such as evolving customer behaviors or real-time risk assessments. These capabilities would be especially valuable in industries like finance, healthcare, and e-commerce, where timely insights are critical.

4. **Validation Across Diverse Domains:**

While the pipeline has shown promise in customer segmentation, applying it across other domains—such as manufacturing, supply chain management, or public health—would provide a broader validation of its adaptability and utility. These cross-domain experiments could uncover new challenges and opportunities for refinement.

5. **Continuous Learning and Adaptation:**

Implementing mechanisms for continuous learning and adaptation would allow the pipeline to evolve as new data becomes available. This approach could reduce the need for frequent retraining while ensuring the system remains relevant and accurate in dynamic environments.

6. **Focus on Ethical and Secure Data Handling:**

As the pipeline scales to include sensitive data, integrating advanced encryption and privacy-preserving methods will be critical. Ensuring compliance with data protection regulations, such as GDPR, will not only safeguard user information but also establish trust with stakeholders.

7. **Collaborative Efforts for Improvement:**

Lastly, this project recognizes that achieving truly impactful results requires collaboration across disciplines and industries. By sharing findings and integrating feedback from diverse perspectives, we aim to enhance the pipeline's robustness and relevance.

6.3 Final Reflection

In conclusion, this project represents an important step forward in leveraging AI for advanced data analysis and decision-making. The pipeline's ability to process and extract actionable insights from large-scale datasets underscores its potential to transform decision-making processes across various industries. However, we remain acutely aware of its limitations and the need for ongoing improvement.

The insights gained from this project provide a strong foundation, but they also reveal the vast scope for further exploration. We recognize that many challenges remain, particularly in ensuring scalability, optimizing clustering techniques, and integrating real-time and dynamic data capabilities. Addressing these challenges will require continuous learning, humility, and a commitment to refining the methods and tools developed here.

Ultimately, we hope that this work will serve as a stepping stone for future advancements, inspiring further research and collaboration. By embracing these challenges with an open and iterative mindset, we believe this pipeline can evolve into a powerful tool for extracting value from data, supporting strategic decision-making, and driving innovation across industries.

Bibliography

- [1] *A Comprehensive Guide to Data Preprocessing in Machine Learning*. Towards Data Science. 2022. URL: <https://towardsdatascience.com/data-preprocessing-in-machine-learning-986e99c590f6> (cit. on p. 53).
- [2] *AI Analytics: What It is, Why It Matters, Use Cases*. Qlik.com. 2023. URL: <https://www.qlik.com/us/augmented-analytics/ai-analytics> (cit. on p. 5).
- [3] *An introduction to Machine Learning: "What are Insights?"*. Telefónica Tech. 2019. URL: <https://telefonicatech.com/en/blog/what-are-insights> (cit. on p. 16).
- [4] *An Overview of Feature Selection Techniques in Machine Learning*. Towards Data Science. 2022. URL: <https://towardsdatascience.com/feature-selection-techniques-in-machine-learning-with-python-f24e7da3f36e> (cit. on p. 31).
- [5] *Apache Spark: Unified Analytics Engine for Big Data*. Apache Software Foundation. 2023. URL: <https://spark.apache.org/docs/latest/> (cit. on p. 59).
- [6] *Box Plot for Outlier Detection*. Towards Data Science. 2020. URL: <https://towardsdatascience.com/understanding-boxplots-5e2df7bcbd51> (cit. on p. 25).
- [7] D. S. Central. *Using Interquartile Range (IQR) for Outlier Detection*. Data Science Central. 2021. URL: <https://www.datasciencecentral.com/iqr-outliers> (cit. on p. 24).
- [8] *Chi-Square Test for Feature Selection in Machine Learning*. Geeks for Geeks. 2022. URL: <https://www.geeksforgeeks.org/feature-selection-using-chi-square-test/> (cit. on p. 32).
- [9] *Cleanlab Open-Source*. Telefónica Tech. 2019. URL: <https://pypi.org/project/cleanlab/> (cit. on pp. 9, 10, 27).
- [10] *Credit Card Dataset for Clustering*. Kaggle. 2021. URL: <https://www.kaggle.com/datasets/arjunbhasin2013/ccdata> (cit. on p. 55).
- [11] F. in Data Science. *Natural Language Processing in the Context of Big Data Analytics*. Frontiers. 2021. URL: <https://www.frontiersin.org/articles/10.3389/fdata.2021.666703/full> (cit. on p. 2).
- [12] *Databricks: Unified Data Analytics Platform*. Databricks. 2023. URL: <https://databricks.com/product/unified-data-analytics-platform> (cit. on p. 59).
- [13] Dataversity. *The Importance of Real-Time Analytics for Modern Businesses*. Dataversity. 2022. URL: <https://www.dataversity.net/importance-real-time-analytics/> (cit. on p. 4).

- [14] J. Doe. *Introduction to Machine Learning: Concepts and Algorithms*. Springer, 2022 (cit. on p. 12).
- [15] Dremio. *Vectorization in NLP*. Dremio. 2023. URL: <https://www.dremio.com/wiki/vectorization-in-nlp/> (cit. on p. 6).
- [16] Elsevier. *Big Data: The Data Deluge, Challenges and Solutions*. ScienceDirect. 2019. URL: <https://www.sciencedirect.com/science/article/pii/S1877050919303408> (cit. on p. 1).
- [17] M. Experts. *K-Nearest Neighbors Imputation for Missing Data*. TechTarget. 2023. URL: <https://www.techtarget.com/knn-imputation> (cit. on p. 21).
- [18] V. Experts. *t-SNE for Visualizing High-Dimensional Data*. AI Research Journal. 2022. URL: <https://www.airesearchjournal.com/tsne> (cit. on p. 36).
- [19] *Explorando as causas da inadimplência: uma análise dos dados de clientes em uma instituição financeira*. Medium. 2023. URL: <https://samequefarias.medium.com/explorando-as-causas-da-inadimplência-uma-análise-dos-dados-de-clientes-em-uma-instituição-7c63313fe463> (cit. on p. 18).
- [20] *Feature Selection Techniques: Forward and Backward Selection*. Towards Data Science. 2022. URL: <https://towardsdatascience.com/forward-and-backward-feature-selection-in-machine-learning-bd557dd0a1d5> (cit. on p. 34).
- [21] GeeksforGeeks. *DBSCAN Clustering in ML: Density-Based Clustering*. GeeksforGeeks. 2023. URL: <https://www.geeksforgeeks.org/dbscan-clustering-in-ml-density-based-clustering/> (cit. on p. 44).
- [22] *GloVe: Global Vectors for Word Representation*. Stanford NLP Group. 2021. URL: <https://nlp.stanford.edu/projects/glove/> (cit. on p. 31).
- [23] M. Grootendorst. *BERTopic: A Comprehensive Guide to Topic Modeling with Transformers*. Maarten Grootendorst (GitHub). 2022. URL: <https://github.com/MaartenGr/BERTopic> (cit. on pp. 5, 9).
- [24] *Histograms for Detecting Outliers*. Data Science Society. 2021. URL: <https://www.datasciencesociety.net/detecting-outliers-using-histograms/> (cit. on p. 25).
- [25] *Information Gain and Entropy for Feature Selection*. Analytics Vidhya. 2023. URL: <https://www.analyticsvidhya.com/blog/2017/03/an-introduction-to-information-gain-and-entropy/> (cit. on p. 33).
- [26] A. Innovations. *Isolation Forest for Outlier Detection*. AI Innovations Journal. 2022. URL: <https://www.aiinnovations.com/isolation-forest> (cit. on p. 25).
- [27] D. S. Insights. *Introduction to Data Preprocessing for Machine Learning*. Data Science Central. 2022. URL: <https://www.datasciencecentral.com/preprocessing> (cit. on p. 19).
- [28] A. R. Journal. *Autoencoders for Missing Data Imputation*. AI Research. 2021. URL: <https://www.airesearchjournal.com/autoencoder-imputation> (cit. on p. 23).
- [29] M. L. Journal. *Regression Imputation in Machine Learning*. MLJ. 2022. URL: <https://www.mljournal.com/regression-imputation> (cit. on p. 21).

-
- [30] S. Journals. *Challenges and Opportunities in NLP-Based Analytics: A Comprehensive Review*. Sage. 2020. URL: <https://journals.sagepub.com/doi/full/10.1177/2158244020931591> (cit. on p. 2).
- [31] *Kaggle Datasets for Data Science Projects*. Kaggle. 2023. URL: <https://www.kaggle.com/datasets> (cit. on p. 53).
- [32] *L2 Normalization in Machine Learning*. Scikit-learn Documentation. 2023. URL: <https://scikit-learn.org/stable/modules/preprocessing.html#normalization> (cit. on p. 29).
- [33] *Lasso and Ridge Regression for Feature Selection*. Analytics Vidhya. 2021. URL: <https://www.analyticsvidhya.com/blog/2020/03/feature-selection-techniques-in-machine-learning/> (cit. on pp. 34, 35).
- [34] C. Lee. *Applications of Machine Learning in Finance and Healthcare*. IEEE Transactions on Neural Networks. 2021, pp. 988–997 (cit. on p. 13).
- [35] A. D. Library. *Artificial Intelligence and Machine Learning for Big Data: Concepts, Algorithms, and Applications*. ACM. 2021. URL: <https://dl.acm.org/doi/10.1145/3343031.3356127> (cit. on pp. 1, 2).
- [36] *Linear Discriminant Analysis (LDA) for Dimensionality Reduction*. Towards Data Science. 2022. URL: <https://towardsdatascience.com/linear-discriminant-analysis-in-python-b0e7e8ed8da5> (cit. on p. 37).
- [37] *Local Outlier Factor (LOF) for Outlier Detection*. Scikit-learn Documentation. 2023. URL: <https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.LocalOutlierFactor.html> (cit. on p. 26).
- [38] *Mahalanobis Distance: Multivariate Outlier Detection*. Analytics Vidhya. 2021. URL: <https://www.analyticsvidhya.com/blog/2021/05/anomaly-detection-using-mahalanobis-distance/> (cit. on p. 26).
- [39] *Max-Abs Normalization Explained*. Medium. 2021. URL: https://medium.com/@maxabs_normalization (cit. on p. 28).
- [40] *Median Absolute Deviation for Outlier Detection*. Geeks for Geeks. 2022. URL: <https://www.geeksforgeeks.org/median-absolute-deviation/> (cit. on p. 25).
- [41] *Min-Max Normalization in Machine Learning*. Geeks for Geeks. 2023. URL: <https://www.geeksforgeeks.org/ml-min-max-scaler/> (cit. on p. 28).
- [42] *MLflow: An Open Source Platform for the Machine Learning Lifecycle*. MLflow Project. 2023. URL: <https://mlflow.org/> (cit. on p. 60).
- [43] *MLlib: Scalable Machine Learning on Spark*. Apache Spark. 2023. URL: <https://spark.apache.org/mllib/> (cit. on p. 59).
- [44] S. Mudadla. *What is Forward and Backward filling of Imputation in Exploratory Data Analysis (EDA)?* Medium. 2024. URL: <https://medium.com/@sujathamudadla1213/what-is-forward-and-backward-filling-of-imputation-in-exploratory-data-analysis-eda-74ec33aca4bd> (cit. on p. 20).
- [45] K. Mukherjee. *Missing value treatment using EM algorithm — A comparative study*. Medium. 2021. URL: <https://kushalmukherjee.medium.com/missing-value-treatment-using-em-a-comparative-study-e4b0d6c9da61> (cit. on p. 22).

- [46] *Normalization vs Standardization in Data Science*. Towards Data Science. 2022. URL: <https://towardsdatascience.com/normalization-vs-standardization-7ff5e57b6e03> (cit. on p. 28).
- [47] OpenText. *What are Intelligent Insights?* OpenText. 2023. URL: <https://blogs.opentext.com/what-are-intelligent-insights/> (cit. on pp. 5, 51).
- [48] E. A. S. Peng Li and D. B. Allison. *Multiple Imputation: A Flexible Tool for Handling Missing Data*. PubMed Central. 2016. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4638176/> (cit. on p. 22).
- [49] *Principal Component Analysis (PCA) for Dimensionality Reduction*. Towards Data Science. 2020. URL: <https://towardsdatascience.com/principal-component-analysis-for-dimensionality-reduction-115a3d157bad> (cit. on p. 35).
- [50] *Recursive Feature Elimination (RFE) for Feature Selection*. Medium. 2023. URL: <https://towardsdatascience.com/recursive-feature-elimination-for-feature-selection-661cf7c99f7d> (cit. on p. 33).
- [51] *Retail Transactional Dataset*. Kaggle. 2021. URL: <https://www.kaggle.com/datasets/bhavikjikadara/retail-transactional-dataset> (cit. on p. 53).
- [52] H. B. Review. *The Big Data Future*. Harvard Business Review. 2023. URL: <https://hbr.org/the-big-data-future> (cit. on pp. 9, 10).
- [53] T. D. Science. *An Introduction to Clustering in Machine Learning*. Towards Data Science. 2021. URL: <https://towardsdatascience.com/an-introduction-to-clustering-in-machine-learning> (cit. on p. 40).
- [54] T. D. Science. *Advanced Data Visualization Techniques for Machine Learning Insights*. Towards Data Science. 2023. URL: <https://towardsdatascience.com/advanced-data-visualization-techniques-machine-learning-insights/> (cit. on pp. 5, 9).
- [55] T. D. Science. *Understanding K-Means Clustering in Machine Learning*. Towards Data Science. 2023. URL: <https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1> (cit. on pp. 9, 27, 41).
- [56] Seaborn. *Seaborn Heatmap Documentation*. Seaborn. 2023. URL: <https://seaborn.pydata.org/generated/seaborn.heatmap.html> (cit. on pp. 11, 49).
- [57] A. Smith and B. Jones. “A Comprehensive Survey of Machine Learning Algorithms and Applications”. In: *Journal of Artificial Intelligence* 45.2 (2023), pp. 234–256 (cit. on p. 12).
- [58] O. D. Specialists. *Survey of Outlier Detection Techniques*. Outlier Research Journal. 2023. URL: <https://www.outlierdetection.com/survey> (cit. on p. 23).
- [59] SpringerLink. *Big Data and the Internet of Things: A Roadmap for Smart Environments*. Springer. 2016. URL: <https://link.springer.com/article/10.1007/s10796-016-9691-6> (cit. on p. 1).
- [60] SpringerLink. *Machine Learning for Big Data Analytics*. Springer. 2023. URL: https://link.springer.com/chapter/10.1007/978-3-031-55639-5_9 (cit. on pp. 10, 12).
- [61] I. Techniques. *Survey of Imputation Methods in Machine Learning*. Towards Data Science. 2023. URL: <https://towardsdatascience.com/imputation-survey> (cit. on p. 19).

- [62] *TF-IDF Explained: Weighing Terms by Importance*. Geeks for Geeks. 2022. URL: <https://www.geeksforgeeks.org/tfidf-for-word-vectorization/> (cit. on p. 30).
- [63] *The Robust Scaler in Machine Learning*. Analytics Vidhya. 2022. URL: <https://www.analyticsvidhya.com/blog/2022/05/a-complete-guide-to-data-scaling-in-machine-learning/> (cit. on p. 28).
- [64] *Tree-Based Feature Selection Techniques*. Geeks for Geeks. 2023. URL: <https://www.geeksforgeeks.org/feature-selection-techniques-in-machine-learning/> (cit. on p. 35).
- [65] *Understanding Bag of Words (BoW) for Text Representation*. Towards Data Science. 2020. URL: <https://towardsdatascience.com/understanding-bag-of-words-5828b5e26038> (cit. on p. 29).
- [66] *Understanding Correlation Coefficients for Feature Selection*. Statistics How To. 2021. URL: <https://www.statisticshowto.com/probability-and-statistics/correlation-coefficient-formula/> (cit. on p. 32).
- [67] *Understanding FastText for NLP*. Facebook AI Research. 2021. URL: <https://fasttext.cc/> (cit. on p. 31).
- [68] *Variance Thresholding for Feature Selection*. Scikit-learn Documentation. 2023. URL: https://scikit-learn.org/stable/modules/feature_selection.html#variance-threshold (cit. on p. 33).
- [69] A. Vidhya. *Clustering Algorithms in Machine Learning: K-Means, DBSCAN, and More*. Analytics Vidhya. 2022. URL: <https://www.analyticsvidhya.com/blog/2022/01/clustering-algorithms-ml/> (cit. on pp. 13, 43).
- [70] *What Are the 5 Vs of Big Data?* Coursera. 2024. URL: <https://www.coursera.org/articles/5-vs-of-big-data> (cit. on pp. 9, 11).
- [71] *What is machine learning (ML)?* IBM. URL: <https://www.ibm.com/topics/machine-learning> (cit. on pp. 9, 13).
- [72] *Word2Vec: The Use of Word Embeddings for NLP*. Medium. 2023. URL: <https://medium.com/@word2vec> (cit. on p. 30).
- [73] I. Xplore. *Big Data Challenges and Opportunities in the Heterogeneous World of Information Technology*. IEEE. 2014. URL: <https://ieeexplore.ieee.org/document/6690320> (cit. on p. 1).
- [74] *Z-Score Outlier Detection: Understanding How it Works*. Statistics How To. 2021. URL: <https://www.statisticshowto.com/probability-and-statistics/z-score/> (cit. on p. 24).
- [75] “*What is Data Science?*” AWS. Amazon. URL: <https://aws.amazon.com/what-is/data-science/> (cit. on p. 10).

@is@a@figure