

IFC-ECCBSO-Bank of Spain Workshop on "New insights from financial statements"

17 October 2024

Outlier detection optimization using machine learning for improving data quality¹

Salah Ahl Mbarek,
Bank Al-Maghrib (Central Bank of Morocco)

¹ This contribution was prepared for the workshop. The views expressed are those of the authors and do not necessarily reflect the views of the European Committee of Central Balance Sheet Data Offices (ECCBSO), the Bank of Spain, the BIS, the IFC or the other central banks and institutions represented at the event.



Topic

Outlier Detection Optimization using Machine Learning for improving Data Quality

Abstract

This paper explores the implementation of advanced machine learning techniques to optimize outlier detection for enhancing data quality within the financial sector at Bank Al Maghrib. The primary focus was on developing and integrating robust algorithms capable of identifying anomalies in extensive datasets of annual financial statements of Moroccan companies. We began by thoroughly examining existing methodologies and their limitations, which informed the selection of cutting-edge techniques such as Isolation Forest, Mahalanobis Distance, DBSCAN, HDBSCAN, and Adversarially Learned Anomaly Detection (ALAD), along with state-of-the-art reinforcement methods like ML-consensus achieved through Stacking and Meta-learning.

The selected algorithms were meticulously analyzed and integrated into an ensemble model to leverage their individual strengths, thereby improving overall anomaly detection accuracy. The implementation phase involved extensive preprocessing and feature engineering to ensure data integrity and relevance. Rigorous testing and validation procedures confirmed the model's efficacy and reliability, demonstrating significant improvements in outlier detection. These improvements contribute to enhanced financial data quality, ensuring more accurate and reliable insights for decision-making.

Additionally, the solution's deployment included a user-friendly interface for analysts to interact with the model, visualize results, and make informed decisions. Comprehensive documentation and training sessions ensured a smooth transition and knowledge transfer to the in-house team, guaranteeing the solution's sustainability.

The paper concludes with several promising avenues for future work, including algorithm refinement, real-time anomaly detection, scalability and performance optimization, integration with other financial systems, expansion to other domains, and continuous learning and adaptation. These perspectives aim to enhance the model's capabilities and ensure its continued relevance in addressing evolving challenges in data quality management.

By harnessing the power of advanced algorithms and ensuring their practical applicability within Bank Al Maghrib, this report lays a robust foundation for future advancements in the field. The insights and results not only contribute to the existing body of knowledge but also provide a practical framework for ongoing innovation and improvement in data quality management.

Keywords : Ensemble Model, Isolation Forest, Mahalanobis Distance, Generative Adversarial Networks (GANs), Adversarially Learned Anomaly Detection (ALAD), Meta-Learning.

Acronymes

	Acronym	Meaning
1	BKAM	Bank Al Maghrib
2	CSE	Casablanca Stock Exchange
3	ML	Machine Learning
4	DL	Deep Learning
5	SVM	Support Vector Machine
6	CRISP-DM	Cross-Industry Standard Process for Data Mining
7	EDA	Exploratory Data Analysis
8	GDPR	General Data Protection Regulation
9	DBSCAN	Density-Based Spatial Clustering of Applications with Noise
10	HDBSCAN	Hierarchical DBSCAN
11	ALAD	Adversarially Learned Anomaly Detection
12	UI	User Interface
13	UML	Unified Modeling Language
14	BDD	Behavioral-Driven Development
15	SEI	Software Engineering Institute

List of Figures

1.1	Logo of the Organization.....	14
1.2	Organizational Chart	15
1.3	Key Axes of the Data and Statistics Strategy	20
2.1	CRISP-DM Chart	31
2.2	Gantt Diagram.....	32
4.1	Use Case Diagram	65
4.2	Class Diagram.....	66
4.3	Sequence Diagram	67
4.4	Activity Diagram.....	68
4.5	Workflow.....	70
4.6	Execution Time Comparison	74
4.7	Accuracy Comparison.....	76
4.8	Precision Comparison.....	76
4.9	Recall Comparison.....	76
4.10	F1-Score Comparison.....	76
4.11	Execution Time Comparison	77
4.12	Accuracy Comparison.....	78
4.13	Precision Comparison.....	78
4.14	Recall Comparison.....	78
4.15	F1-Score Comparison.....	78
4.16	Area Under the Curve Variation for Models	79
4.17	Training Time Comparison of Deep Learning-Based Methods.....	80
4.18	Accuracy Comparison.....	80
4.19	Precision Comparison.....	80
4.20	Recall Comparison.....	81
4.21	F1-Score Comparison.....	81
4.22	AUC Comparison of Deep Learning-Based Methods.....	81
4.23	Loss Function Evolution of Autoencoders for varying number of data entries	82
4.24	Accuracy Variation by Model.....	85
4.25	Prediction Time Variation by Model.....	85
5.1	Development Tools.....	88
5.2	Programming Language.....	88
5.3	Application Development Tools.....	88
5.4	Model Creation Tools	88
5.5	Testing Tools: Jest and Cucumber	88
5.6	Version Control: Git and GitHub	89
5.7	Form Interface.....	90
5.8	App passing unitary and integration tests.....	92

5.9	App passing Behavioral Driven tests.....	93
5.10	User documentation detailing installation and usage instructions.	94
5.11	Developer documentation showcasing application architecture and codebase details.	95

List of Algorithms

1	Synthetic Dataset Generation	73
2	Stacking for Anomaly Detection.....	84

Contents

Abstract.....	3
Acronymes.....	4
General introduction	10
1 Context	13
1.0 Introduction to chapter	13
1.1 Presentation of the hosting organization.....	14
1.2 Framework of Monetary Policy.....	15
1.3 General context.....	19
1.4 Thematic.....	21
1.5 Conclusion.....	23
2 Problem Analysis.....	24
2.0 Introduction to chapter	24
2.1 Problem Statement	25
2.2 Critique of Existing Methods.....	26
2.3 Solution Proposal.....	27
2.4 Requirements Specification.....	28
2.5 Methodology of work – CRISP-DM	30
2.6 Project Planning.....	31
2.7 Conclusion.....	32
3 Theoretical Background.....	33
3.0 Introduction to chapter	33
3.1 Background History.....	35
3.2 Distance-Based Approaches	37
3.3 Density-Based Approaches	41
3.4 Machine Learning-Based Approaches.....	46
3.5 Deep Learning-Based Approaches	50
3.6 ML-Consensus.....	53
3.7 Data Synthesis	59
3.8 Conclusion.....	61
4 Solution Engineering.....	63
4.0 Introduction to chapter	63
4.1 Solution Architecture	64

4.2	Data Synthesis	71
4.3	Solution Analysis.....	73
4.4	Optimisations and/or Alternatives.....	83
4.5	Conclusion.....	85
5	Implementation of the Solution.....	87
5.0	Introduction to chapter	87
5.1	Used Technologies.....	88
5.2	The desktop App	89
5.3	Quality Assurance and Testing.....	91
5.4	Documentation	93
5.5	Conclusion.....	96
	Conclusions and Perspectives.....	98
	Bibliography	102

General introduction

Background

In the modern financial landscape, data accuracy and integrity are paramount, especially for institutions like the Central Bank of Morocco (Bank Al Maghrib) that rely on precise financial data to inform policy decisions and maintain economic stability. Financial data quality directly impacts the credibility and reliability of financial statements, risk assessments, and economic forecasts. One significant challenge in maintaining high data quality is the presence of outliers—data points that deviate significantly from other observations. These outliers can result from errors in data entry, fraud, or unusual but legitimate transactions, and can distort analyses, leading to incorrect conclusions and potentially harmful financial decisions.

In Morocco, the annual financial statements of companies provide critical information for economic analysis and policy-making. Ensuring the accuracy of these financial statements is essential for regulatory oversight, investor confidence, and overall economic health. The traditional process of identifying and correcting outliers in these financial statements is labor-intensive and subject to human error. Manual re-evaluation of financial figures is not only time-consuming but also prone to inconsistencies. As the volume and complexity of financial data continue to grow, there is a pressing need for an automated, robust, and accurate method for outlier detection.

Problem Statement

The Central Bank of Morocco faces significant challenges in maintaining the accuracy and integrity of financial data reported by Moroccan companies. The current manual methods for re-evaluating financial statements to identify and correct outliers are inefficient and error-prone. These methods cannot keep pace with the increasing volume and complexity of data, leading to potential inaccuracies in financial reporting and analysis. This thesis addresses the need for an automated approach to outlier detection that can efficiently and accurately identify anomalies in financial data, thereby improving the overall data quality.

Objectives

This thesis aims to develop and optimize machine learning models for the automated detection of outliers in the annual financial statements of Moroccan companies. The primary objectives of the study are:

- **To Review Existing Outlier Detection Techniques**

- Conduct a comprehensive review of current outlier detection methods, including statistical approaches and machine learning techniques.
- Identify the strengths and limitations of these methods in the context of financial data.
- **To Develop Machine Learning Models for Outlier Detection**
 - Implement and evaluate various machine learning algorithms such as Isolation Forest, and other relevant models.
 - Compare the performance of these models in terms of accuracy, efficiency, and scalability.
- **To Optimize the Performance of Machine Learning Models**
 - Apply optimization techniques such as hyperparameter tuning, ensemble methods, and feature engineering to enhance model performance.
 - Evaluate the impact of these optimizations on the accuracy and reliability of outlier detection.
- **To Integrate the Optimized Models into the Central Bank's Data Processing Pipeline**
 - Develop a practical framework or internal tool (Software) for implementing the optimized outlier detection models within the existing data infrastructure at the Central Bank of Morocco.
 - Ensure the solution is scalable, efficient, and capable of handling large volumes of financial data.

Importance of Study

Enhancing the process of outlier detection in financial data has far-reaching implications for the Central Bank of Morocco and the broader financial ecosystem. Improved data quality leads to better decision-making, more accurate financial analyses, and more effective regulatory oversight. By automating and perfecting the re-evaluation of financial data, this study contributes to the reliability and transparency of financial reporting, thereby supporting the overall stability and integrity of the Moroccan financial system.

Automated outlier detection will not only reduce the workload on analysts but also provide a consistent and reliable method for identifying anomalies. This, in turn, enhances the credibility of financial statements and boosts confidence among investors, regulators, and other stakeholders.

Structure of the Report

The thesis is structured as follows:

- **Chapter 1: Context** This chapter Provides an overview of the hosting organization, Bank Al Maghrib, and the monetary policy framework in Morocco.
- **Chapter 2: Problem Analysis** Discusses the problem statement, current methodologies, their limitations, and the proposed solution along with detailed specifications and work methodology.
- **Chapter 3: Theoretical Background** Reviews the existing literature and theories related to outlier detection and cutting-edge techniques in the field.

- **Chapter 4: Solution Engineering** Details the data collection and synthesis, preprocessing, model development, optimization techniques, and evaluation metrics.
- **Chapter 5: Implementation of the Solution** Describes the practical implementation of the models, including user interfaces, model training and testing, quality assurance, documentation, and deployment.

Finally, the thesis concludes with a summary of findings, implications, and suggestions for future work, followed by a comprehensive bibliography and annexes.

By the end of this study, we aim to provide a comprehensive solution for the automated detection of outliers in financial data, enhancing the Central Bank of Morocco's capability to ensure data quality and make informed decisions based on accurate financial information.

Chapter 1

Context

This chapter sets the foundation for the entire study by providing a comprehensive overview of Bank Al Maghrib. Understanding the organizational structure, functions, and missions of Bank Al Maghrib is crucial, as it is the central authority responsible for monetary policy and financial stability in Morocco. We will explore the general description of the institution, its roles, and its organizational hierarchy to understand how it operates and how decisions are made.

Furthermore, the chapter delves into the framework of monetary policy in Morocco, which includes an analysis of the financial market, the Moroccan financial system, and the monetary policy mechanisms. This detailed context is essential for comprehending the environment in which the study is conducted and highlights the importance of maintaining high data quality in financial reporting. By understanding the intricacies of monetary policy and the financial system, we can appreciate the critical role that accurate financial data plays in informing policy decisions and ensuring economic stability.

This chapter's content is directly linked to the problem analysis in the following chapter, where we will discuss the challenges associated with data quality and outlier detection in financial datasets. By establishing the context, we provide a background that underscores the relevance and necessity of our study, ensuring that the reader comprehends the foundational elements that drive the need for improved outlier detection methods.

1.1 Presentation of the hosting organization

1.1.1 Overview



Figure 1.1 – Logo of the Organization

Bank Al Maghrib, the Central Bank of Morocco, is a pivotal institution in the Moroccan financial system. Established in 1959, its primary mission is to ensure price stability and the soundness of the financial system. As a central authority, it regulates and oversees the banking sector, implements monetary policy, and manages the country's foreign exchange reserves. The bank also plays a critical role in formulating and executing policies that promote sustainable economic growth and financial inclusion.

1.1.2 Functions and Missions

Bank Al Maghrib's functions and missions are diverse and encompass several key areas:

Monetary Policy Implementation

The bank designs and implements monetary policy to maintain price stability and control inflation. It uses various instruments such as interest rate adjustments, open market operations, and reserve requirements to influence money supply and demand.

Financial Stability

Ensuring the stability of the financial system is a core mandate. The bank monitors and assesses systemic risks, supervises financial institutions, and enforces prudential regulations to prevent financial crises.

Banking Supervision and Regulation

Bank Al Maghrib oversees the banking sector to ensure it operates efficiently and safely. This involves licensing new banks, conducting regular inspections, and ensuring compliance with regulatory standards.

Foreign Exchange Management

The bank manages Morocco's foreign exchange reserves, conducts foreign exchange operations, and maintains the stability of the national currency, the Moroccan Dirham.

Economic Research and Statistics

Conducting economic research and providing reliable statistical data are essential functions. The bank analyzes economic trends, publishes reports, and offers insights to inform policy decisions and guide market participants.

Consumer Protection

Protecting the interests of consumers and promoting financial literacy are also important missions. The bank ensures transparency in financial services, addresses consumer complaints, and educates the public about financial products and services.

Payment Systems Oversight

Bank Al Maghrib oversees the payment systems to ensure their efficiency, security, and reliability. It regulates and monitors payment instruments and infrastructures to facilitate smooth and secure transactions.

1.1.3 Organizational Chart

The organizational chart of Bank Al-Maghrib provides an overview of the hierarchical structure and roles within the central bank. It illustrates the relationships and reporting lines among different departments and units, highlighting key positions and their responsibilities in the bank's operations and decision-making processes.

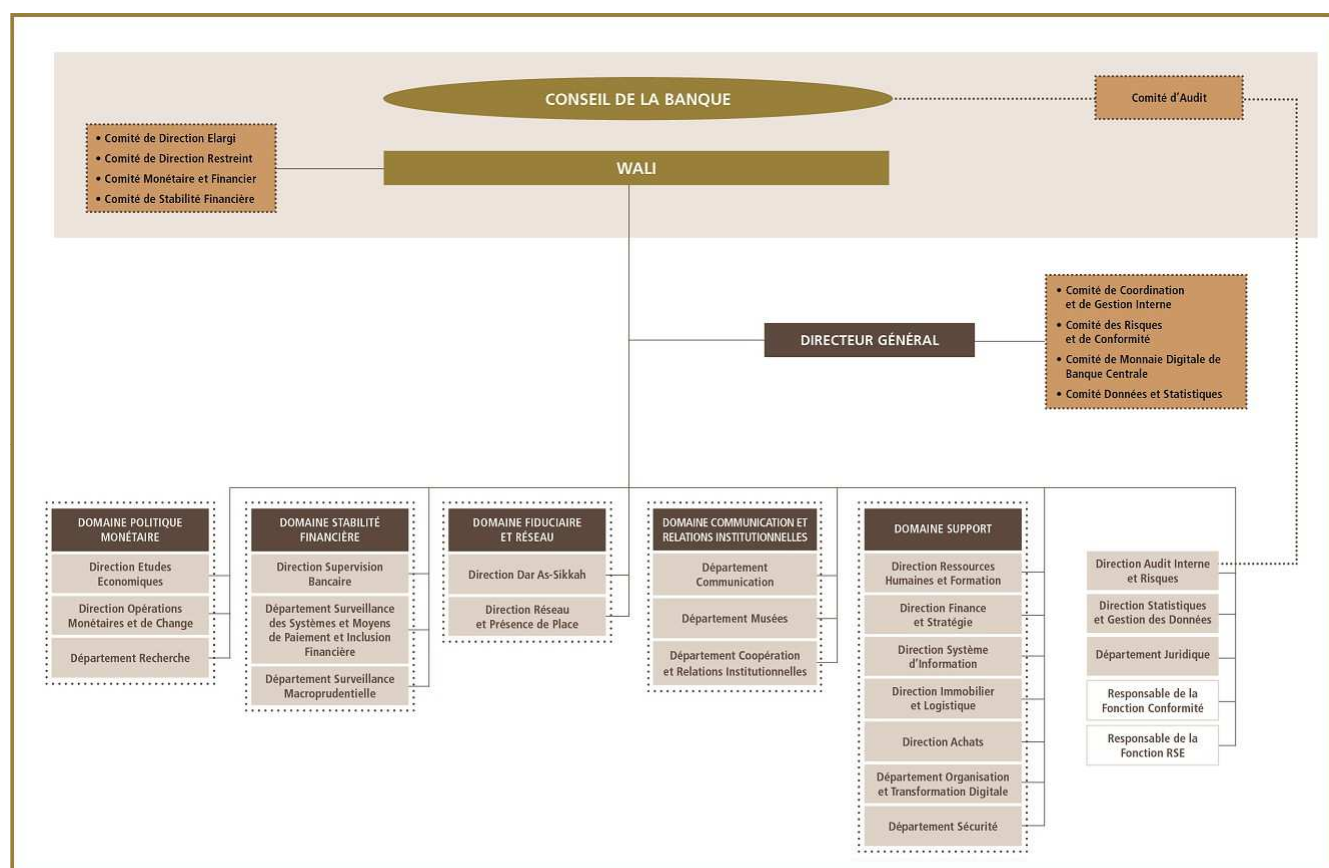


Figure 1.2 – Organizational Chart

This chart serves as a visual representation of the bank's organizational hierarchy, demonstrating how various departments collaborate and contribute to fulfilling the bank's mandates and objectives.

1.2 Framework of Monetary Policy

1.2.1 Financial Market

The financial market in Morocco plays a crucial role in the overall economic framework of the country. It encompasses a variety of sectors, including the stock market, bond market, and foreign exchange market. The primary stock exchange, the Casablanca Stock Exchange (CSE) – Bourse

de Casablanca, is one of the largest in Africa and serves as a vital platform for raising capital for companies. The bond market, meanwhile, provides a mechanism for the government and corporations to borrow funds through the issuance of debt securities.

The foreign exchange market in Morocco, regulated by Bank Al Maghrib, facilitates the exchange of the Moroccan Dirham with other currencies, thus supporting international trade and investment. This market is integral to maintaining the stability of the Dirham and managing the country's foreign exchange reserves. Moreover, the financial market in Morocco is characterized by a growing number of financial instruments and services, which are designed to meet the diverse needs of investors and borrowers.

Bank Al Maghrib plays a key role in overseeing and regulating these markets to ensure their stability, transparency, and efficiency. The central bank implements various policies and regulations to safeguard the interests of investors and maintain the overall health of the financial system.

1.2.2 Moroccan Financial System

The Moroccan financial system is characterized by a structured network of institutions including commercial banks, insurance companies, pension funds, and other financial intermediaries. Dominating this landscape are several prominent commercial banks offering a wide range of services such as deposits, loans, and payment processing. While the sector features some competition, it is primarily led by a few well-established domestic banks, supported by a regulatory framework overseen by Bank Al Maghrib.

Insurance companies and pension funds also play significant roles, providing critical services for risk management and retirement planning, enhancing financial security for individuals and businesses. Additionally, microfinance institutions and non-bank financial entities contribute to financial inclusion by serving underserved populations.

Bank Al Maghrib maintains the stability of the financial system through rigorous supervision and regulation. The central bank sets prudential standards, conducts regular inspections, and enforces compliance with regulatory requirements, thereby mitigating risks and ensuring a stable financial environment supportive of economic growth.

1.2.3 Monetary Policy

Monetary policy in Morocco is formulated and implemented by Bank Al Maghrib with the primary objective of maintaining price stability. The central bank employs a range of tools and instruments to achieve this goal, including interest rate adjustments, open market operations, and reserve requirements.

Interest rate policy is a cornerstone of monetary policy. By setting the key policy rate, Bank Al Maghrib influences short-term interest rates, which in turn affect borrowing and lending rates across the economy. This mechanism helps control inflation and support economic activity.

Open market operations involve the buying and selling of government securities in the open market. Through these operations, Bank Al Maghrib can manage the money supply and influence liquidity conditions in the banking system. This tool is used to smooth out fluctuations in interest rates and ensure that the banking system has adequate liquidity to meet the needs of the economy.

Reserve requirements are another important tool. By setting the minimum reserves that banks must hold, Bank Al Maghrib can influence the amount of money that banks can lend. This helps

regulate the money supply and control inflationary pressures.

In addition to these traditional tools, Bank Al Maghrib monitors various economic indicators and financial conditions to make informed policy decisions. The central bank's actions are guided by a commitment to transparency and accountability, with regular communication to the public and financial markets about policy decisions and their rationale.

1.2.4 Instruments of Monetary Policy

Bank Al Maghrib employs a range of monetary policy instruments to achieve its primary objective of price stability and to support the broader economic goals of the country. These instruments include:

Interest Rate Policy

The central bank sets the key policy rate, which influences short-term interest rates across the economy. Adjustments to this rate impact borrowing and lending rates, affecting economic activity and inflation. The policy rate serves as a benchmark for other interest rates within the banking system.

Open Market Operations

Bank Al Maghrib conducts open market operations by buying and selling government securities. These operations manage liquidity conditions in the banking system, influencing the money supply and short-term interest rates. By increasing or decreasing the amount of money in circulation, the central bank can smooth out fluctuations in interest rates and ensure sufficient liquidity for economic needs.

Reserve Requirements

The central bank sets minimum reserve requirements for commercial banks, determining the proportion of customer deposits that must be held as reserves. By adjusting these requirements, Bank Al Maghrib can influence the amount of money banks can lend, thereby controlling the money supply and exerting an impact on inflation and economic activity.

Standing Facilities

These facilities provide and absorb overnight liquidity to and from the banking system. The central bank offers lending facilities (e.g., overnight lending facilities) to banks facing short-term liquidity shortages, and deposit facilities to absorb excess liquidity. These facilities help maintain stability in the money market and ensure the smooth functioning of the banking system.

Foreign Exchange Interventions

To maintain the stability of the Moroccan Dirham, Bank Al Maghrib may intervene in the foreign exchange market. These interventions involve buying or selling foreign currencies to influence exchange rates and manage the country's foreign exchange reserves. The central bank's actions in

the foreign exchange market help mitigate excessive volatility and ensure a stable external value of the Dirham.

Communication and Forward Guidance

Transparency and effective communication are critical components of Bank Al Maghrib's monetary policy strategy. The central bank regularly communicates its policy decisions, economic outlook, and rationale for its actions to the public and financial markets. Forward guidance provides insights into the likely future path of monetary policy, helping to shape expectations and influence economic behavior.

1.2.5 Recent Monetary Policy Decisions and Their Impacts

In recent years, Bank Al Maghrib has made several key monetary policy decisions aimed at addressing the evolving economic challenges and maintaining price stability. These decisions and their impacts are summarized below:

Interest Rate Adjustments

In response to global economic uncertainties and domestic inflationary pressures, Bank Al Maghrib has adjusted the key policy rate several times. For instance, during periods of economic slowdown, the central bank has lowered the policy rate to stimulate borrowing and investment, supporting economic growth. Conversely, to counteract rising inflation, the policy rate has been increased to curb excessive demand and control price levels.

- **Impact:** Lower interest rates have facilitated increased access to credit for businesses and households, boosting consumption and investment. Higher rates have helped contain inflationary pressures and stabilize the economy.

Liquidity Management through Open Market Operations

Bank Al Maghrib has actively conducted open market operations to manage liquidity conditions in the banking system. By purchasing or selling government securities, the central bank has ensured adequate liquidity to support economic activity while preventing excessive money supply growth.

- **Impact:** Effective liquidity management has contributed to the stability of short-term interest rates, enhanced the functioning of the money market, and supported the banking system's ability to meet the economy's credit needs.

Foreign Exchange Interventions

To address external shocks and maintain exchange rate stability, Bank Al Maghrib has intervened in the foreign exchange market. These interventions have included buying or selling foreign currencies to manage exchange rate fluctuations and stabilize the Moroccan Dirham.

- **Impact:** Foreign exchange interventions have helped mitigate excessive volatility in the exchange rate, supported export competitiveness, and maintained confidence in the Dirham's stability.

Enhanced Communication and Forward Guidance

Bank Al Maghrib has strengthened its communication strategy by providing clear and transparent information on its policy decisions, economic assessments, and future policy intentions. This forward guidance has helped shape market expectations and enhance the effectiveness of monetary policy.

- **Impact:** Improved communication has reduced uncertainty, influenced economic agents' expectations, and facilitated better-informed decision-making by businesses and households.

Reserve Requirement Adjustments

In response to changing economic conditions, Bank Al Maghrib has adjusted the reserve requirements for commercial banks. These adjustments have been used to influence the amount of money banks can lend, thereby impacting credit conditions and money supply growth.

- **Impact:** Changes in reserve requirements have helped manage liquidity in the banking system, supported financial stability, and ensured that credit conditions align with the central bank's monetary policy objectives.

Overall, Bank Al Maghrib's recent monetary policy decisions have been instrumental in maintaining price stability, supporting economic growth, and ensuring the stability of the financial system. The central bank's proactive and adaptive approach has enabled it to navigate complex economic challenges and contribute to Morocco's economic development.

1.3 General context

The "Direction Statistiques et Gestion des Données" at Bank Al-Maghrib plays a pivotal role in enhancing the central bank's data collection and statistical production capabilities to support its policies and fulfill its mandates. Over recent years, significant advancements have been made in several key areas within the department:

1.3.1 Internal Developments

Bank Al-Maghrib has strengthened its information assets through various initiatives:

- **Enrichment of Data Sources:** Establishing frameworks for data exchange with key partners and expanding data collection to include macro-prudential frameworks, payment system surveillance, and financial inclusion.

- **Financial Information Centers:** Implementation of financial information centers covering data on monetary and exchange markets, bank accounts, check and bill defaults, as well as company balances and profiles.

- **New Indices:** Introduction of new indices such as the real estate asset price index and banking service price index.

- **Statistical Surveys:** Launch of new statistical surveys, including credit granting conditions and inflation expectations.

- **Modernization of Statistical Production:** Continuous upgrading of monetary statistics

production to ensure comprehensive coverage of institutional sectors in line with international standards.

1.3.2 Enhanced Data Utilization

The department has invested in operational systems to enhance data collection, dissemination tools, and diverse formats for data exchange. This includes:

- Deployment of statistical series warehouses and a decision support system.
- Establishment of a functional competence center to unify practices in data processing, calculation, indicator reporting, and statistical analysis, fostering knowledge sharing and enhancement.

1.3.3 Regulatory Framework and Centralized Data Management

Bank Al-Maghrib has implemented a governance framework for data governance, focusing on monetary policy, banking operations, and common interest services. Key initiatives include:

- Creation of a centralized entity for data management and statistical production, allowing business entities to focus on core activities.
- Consolidation of common data repositories and initial deployment of advanced data exploration and analysis capabilities.

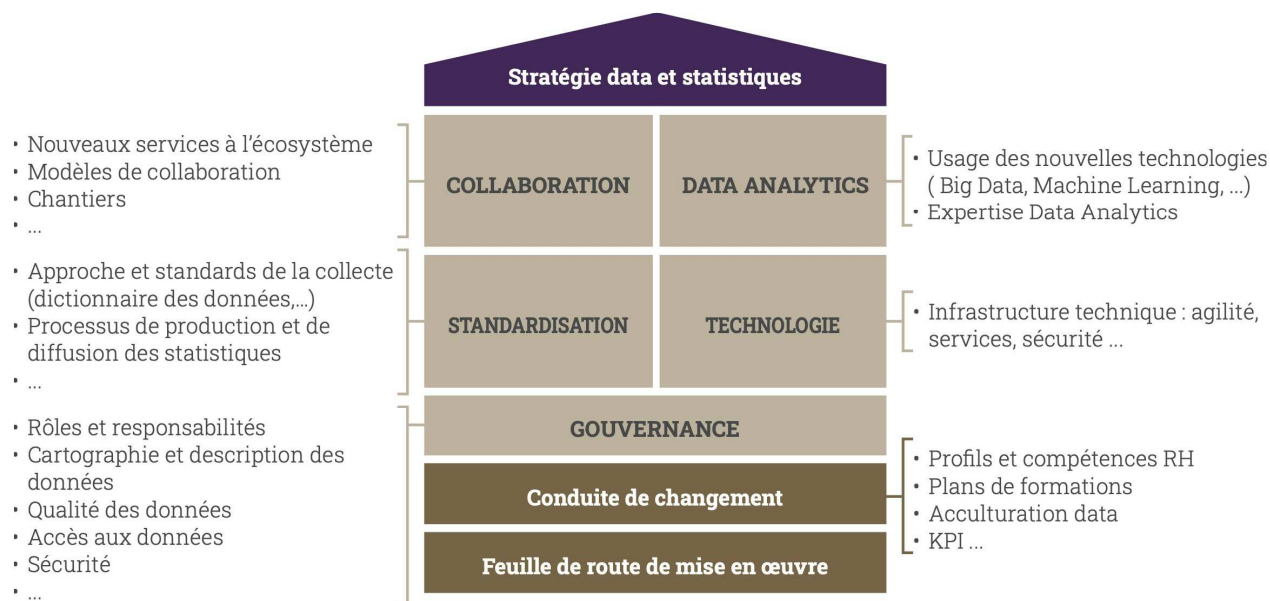


Figure 1.3 – Key Axes of the Data and Statistics Strategy

1.3.4 External Challenges and Technological Advancements

Externally, central banks operate in an environment characterized by:

- Rapid growth in new and diverse data sources following economic changes and increased reporting requirements post-financial crisis.

- Growing demand for granular data for integrated analytical needs, necessitating enhanced data access through sharing platforms.
- Emergence of new information technologies facilitating advanced data exploration and analysis, including alternative sources such as websites, Google searches, and multimedia content.

These challenges underscore the importance for central banks to deliver high-quality information efficiently, supporting international comparisons and individual behavior analysis. The focus on mastering specific technologies like Big Data has shifted data to the forefront of business strategy, emphasizing sharing over silos and prioritizing data quality management and information system resilience.

Bank Al-Maghrib's high maturity in data and statistical management positions it well to accelerate its transformation, aligning with international central banking practices. This transformation, anchored in a structured data and statistical strategy, aims to enhance operational efficiency, mitigate risks through data quality control, and improve governance for prioritizing data-related investments.

1.3.5 Strategy Implementation Approach

The implementation of Bank Al-Maghrib's data and statistical strategy involves collaborative efforts among internal stakeholders and key partners, whether subject to the bank's oversight or not. This approach encompasses:

- Assessment of internal and ecosystem landscapes, alongside benchmarking against international trends.
- Definition of target strategy focusing on data governance, process standardization, technological infrastructure, data analytics development, and ecosystem collaboration.
- Development of an implementation roadmap, estimation of financial and human resources, and change management plan to guide execution.

This comprehensive strategy aims to empower Bank Al-Maghrib with robust data capabilities, fostering innovation, resilience, and informed decision-making in support of its core mandates and strategic objectives.

1.4 Thematic

This section provides a detailed overview of the thematic focus, outlining the subject matter, the tasks involved, the constraints, objectives, and planning required to achieve the project's goals.

1.4.1 Presentation of the Subject

The project is centered on enhancing the process of detecting and managing outliers in financial data using advanced machine learning techniques. The goal is to optimize existing methodologies to improve data quality, ensuring accurate and reliable financial statistics and reports. This involves understanding and refining the current approaches used by Bank Al-Maghrib to identify anomalies in vast datasets.

1.4.2 Tasks Involved

The project involves several critical tasks:

Analysis and Understanding of the Existing Program

- Review and comprehend the existing data processing and anomaly detection systems.
- Engage in discussions with domain experts to gather insights into the current methodologies and their limitations.
- Identify potential areas for improvement in terms of efficiency, accuracy, and robustness of the current system.

Optimization of the System

- Implement improvements based on the analysis to enhance the system's performance and reliability.
- Refine processes to ensure better handling of data and more accurate detection of outliers.

Testing and Validation

- Conduct thorough testing of the optimized system to ensure it meets the desired performance criteria.
- Validate the system's accuracy in identifying outliers and its robustness in handling different data scenarios.

Documentation and Reporting

- Update the system documentation to reflect the improvements and changes made.
- Prepare comprehensive reports detailing the project progress, methodologies used, and outcomes achieved.

1.4.3 Constraints

The project faces several constraints, including:

- **Data Quality:** Ensuring the collected data is accurate, complete, and reliable.
- **Time Limitations:** Completing the project within the specified period.
- **Resource Availability:** Managing the computational and human resources required for the project.
- **Regulatory Compliance:** Adhering to legal and regulatory requirements related to data privacy and security.

1.4.4 Objectives to Achieve

The key objectives of the project include:

- **Enhancing Data Quality:** Improving the overall quality of financial data through better anomaly detection.
- **Efficiency:** Streamlining processes to reduce the time and effort required for data validation.

- **Scalability:** Developing solutions that can handle increasing amounts of data without compromising performance.
- **Knowledge Transfer:** Ensuring that the improved methodologies and processes can be effectively transferred to and utilized by Bank Al-Maghrib's team.

1.4.5 Planning

The project planning involves several key steps:

1. **Initial Assessment:** Conducting a comprehensive assessment of the current anomaly detection processes and identifying key areas for improvement.
2. **Design Phase:** Developing a detailed plan for the enhancements, including theoretical frameworks and methodologies.
3. **Implementation Phase:** Applying the planned improvements to the existing system.
4. **Testing Phase:** Rigorous testing of the enhanced system to ensure it meets all specified requirements.
5. **Deployment Phase:** Implementing the enhanced system in the operational environment.
6. **Training and Knowledge Transfer:** Conducting training sessions and providing documentation to ensure effective use and maintenance of the enhanced system.
7. **Evaluation and Feedback:** Continuously evaluating the performance of the enhanced system and making necessary adjustments based on feedback.

1.5 Conclusion

This chapter has provided an essential overview of Bank Al Maghrib, including its structure, functions, and the framework of monetary policy in Morocco. This contextual understanding is critical as it sets the stage for the problem analysis in the next chapter. By comprehending the operational environment and the significance of accurate financial data in the context of monetary policy, we are better positioned to address the challenges related to data quality.

The insights gained from this chapter highlight the importance of maintaining high data quality in financial reporting, which is a central theme in the subsequent chapters. In Chapter 2, we will delve into the specific problems associated with outlier detection in financial datasets, analyzing existing methods, their limitations, and proposing a solution to enhance data quality. The foundational knowledge provided in this chapter ensures that the reader is well-prepared to understand and engage with the problem analysis and the proposed methodologies in the following chapters.

Chapter 2

Problem Analysis

Building upon the foundational understanding established in Chapter 1, this chapter delves into the core problem statement addressed in this thesis: optimizing outlier detection methods to enhance the quality and reliability of financial data at Bank Al-Maghrib. By situating the problem within the broader context of financial data integrity and regulatory compliance, we aim to critically evaluate existing methodologies and propose advanced techniques that can overcome current limitations.

Financial institutions, particularly central banks like Bank Al-Maghrib, rely on accurate financial data for crucial decision-making processes. The integrity of this data is pivotal not only for regulatory compliance but also for ensuring economic stability and facilitating effective policy formulation. Anomalies in financial datasets, whether due to errors, fraud, or systemic risks, can significantly impact the reliability of analysis and subsequent decisions. Therefore, robust outlier detection mechanisms are essential to safeguarding data quality and maintaining trust in financial systems.

This chapter aims to provide a comprehensive analysis of the existing outlier detection methods employed at Bank Al-Maghrib. It will critically examine the strengths and limitations of these methods, identifying key challenges that hinder their effectiveness in a dynamic financial environment. By elucidating these challenges, we lay the groundwork for proposing enhancements and advanced techniques that promise to improve anomaly detection accuracy, scalability, and efficiency.

Chapter 1 established the foundational context of Bank Al-Maghrib's role and the significance of financial data integrity. Chapter 2 builds upon this context by focusing specifically on outlier detection within financial datasets. The insights gained from Chapter 2 will serve as a bridge to Chapter 3, where we delve into the theoretical underpinnings of outlier detection techniques. By understanding the specific challenges and requirements identified in this chapter, we can explore theoretical frameworks and methodologies that address these challenges comprehensively.

2.1 Problem Statement

2.1.1 General context

Bank Al-Maghrib, the Central Bank of Morocco, plays a pivotal role in maintaining financial stability and driving economic growth. A critical aspect of its mandate involves meticulous analysis of financial data to inform strategic monetary policy decisions. Given the intricate nature and vast volume of financial transactions, detecting anomalies within this data is paramount to safeguarding its integrity and ensuring accurate policy formulation. Anomalies may signify errors, fraudulent activities, or irregularities that, if undetected, could potentially distort policy decisions or compromise financial stability.

2.1.2 Problematic

This project focuses on optimizing outlier detection methods to enhance the quality of financial datasets at Bank Al-Maghrib. Current practices primarily rely on algorithms such as manual double-checking and Mahalanobis Distance for anomaly detection. While effective, these methods face challenges related to performance, accuracy, and scalability. The objective is to refine these existing methods to achieve more precise and efficient anomaly detection capabilities. By doing so, the project aims to elevate the overall quality and reliability of financial data analysis, thereby reinforcing the central bank's ability to make informed decisions and maintain financial stability.

2.1.3 Analysis of existing method

The current outlier detection methods utilized at Bank Al-Maghrib primarily involve the computation of Mahalanobis Distance (MD) (Mahalanobis, 1936) to identify anomalous values in financial data. This method plays a critical role in enhancing the reliability of financial data analysis by focusing efforts on potentially anomalous data points as in the identification of outliers in multivariate data, rather than manually reviewing every entry in annual financial statements (Rocke and Woodruff, 1996).

Mahalanobis Distance measures the distance of a point from the centroid of a distribution, considering the covariance structure of the data. This method plays a critical role in enhancing the reliability of financial data analysis by focusing efforts on potentially anomalous data points rather than manually reviewing every entry in annual financial statements of Moroccan companies (Bank Al Maghrib, 2023).

Mahalanobis Distance is computed using the formula:

$$D_M(x) = \sqrt{(x - \mu)^T \Sigma^{-1}(x - \mu)}$$

where x is the data point, μ is the mean vector, and Σ is the covariance matrix.

However, despite its benefits, MD faces challenges related to computational complexity and scalability. The computation involves matrix inversion and multiplication operations, which become computationally expensive as the dimensionality of data increases (Ren et al., 2021). This computational burden limits its scalability for large-scale financial datasets, where timely anomaly detection is crucial for effective decision-making and policy formulation (Xu et al., 2023).

Moreover, the assumption of multivariate normality underlying Mahalanobis Distance may not hold in real-world financial data, where distributions are often non-Gaussian and exhibit complex patterns. This limitation necessitates preprocessing steps to transform data or alternative outlier detection techniques when dealing with diverse data distributions (Kamenetsky Yadan, 2021).

Despite these challenges, Mahalanobis Distance remains integral to the anomaly detection framework at Bank Al-Maghrib, providing a foundational approach to flagging potentially irregular financial transactions and data entries. Future enhancements could explore hybrid approaches combining MD with machine learning techniques like deep learning-based models or ensemble methods to improve accuracy and scalability in detecting anomalies across diverse financial datasets (Chalapathy et al., 2019).

2.2 Critique of Existing Methods

The current outlier detection methods utilized at Bank Al-Maghrib, primarily employing Mahalanobis Distance, play a crucial role in safeguarding the integrity of financial data. Mahalanobis Distance (MD) calculates the distance of each data point from the distribution, facilitating the identification of anomalies in Moroccan companies' annual statements. This statistical approach is advantageous for its ability to account for correlations between variables, providing a robust measure of deviation from the norm.

However, despite its strengths, MD and Isolation Forest encounter significant challenges that impact their effectiveness under evolving data landscapes. One of the prominent challenges is the computational efficiency required for large-scale financial datasets. Isolation Forest, renowned for its capability to isolate anomalies efficiently, can suffer from computational inefficiencies as dataset sizes expand. This can lead to prolonged processing times and resource-intensive operations, affecting the system's responsiveness to dynamic market conditions and emerging financial risks (Liu et al., 2008).

Similarly, Mahalanobis Distance's reliance on covariance matrix calculations poses computational burdens, particularly in high-dimensional datasets where the inversion of matrices becomes computationally expensive. As the dimensionality of data increases, the scalability of MD diminishes, making it less suitable for real-time anomaly detection and proactive decision-making in financial monitoring (Ren et al., 2021).

Moreover, the complexity of implementing and maintaining these algorithms presents operational challenges. Isolation Forest and Mahalanobis Distance require expertise in parameter tuning and algorithmic optimization to achieve optimal performance. However, the lack of standardized guidelines and best practices for their application in financial settings can lead to inconsistent results and sub-optimal anomaly detection outcomes, impacting the reliability of financial data analysis (Park et al., 2018).

Furthermore, the assumptions underlying Isolation Forest and Mahalanobis Distance algorithms limit their applicability to diverse data distributions and anomaly patterns. Isolation Forest, designed for data with a uniform distribution, may struggle with datasets exhibiting complex data distributions or skewed anomalies. Similarly, Mahalanobis Distance assumes multivariate normality, restricting its effectiveness in detecting anomalies in non-Gaussian datasets without rigorous preprocessing and transformation steps (Kamenetsky Yadan, 2021).

To address these challenges and enhance the reliability of outlier detection systems at Bank Al-Maghrib, future advancements should explore scalable anomaly detection techniques. Deep learning-based approaches and ensemble methods offer promising alternatives capable of handling large-scale

financial datasets with diverse distributions and evolving anomaly profiles (Chalapathy et al., 2019).

By integrating these innovations, the central bank can bolster its data-driven decision-making processes, strengthen financial stability measures, and ensure the robustness of financial data analysis in Morocco.

2.3 Solution Proposal

To effectively address the challenges inherent in outlier detection within financial datasets, the proposed solution at Bank Al-Maghrib emphasizes integrating machine learning (ML) and deep learning (DL) techniques. The primary objective is to enhance the accuracy, efficiency, and scalability of anomaly detection systems, which are critical for maintaining the integrity of financial data and supporting informed policy decisions.

The proposed improvements include:

- Leveraging Machine Learning over traditional ways of Outlier detection
- Enhancing the codebase for better readability, maintainability, and scalability.
- Incorporating advanced data structures and caching mechanisms to boost performance.
- Implementing rigorous testing protocols to ensure the reliability and robustness of the optimized system.

The solution leverages ML models such as ensemble methods, support vector machines (SVM), and isolation forests to augment traditional statistical methods for anomaly detection (Chalapathy et al., 2019). These techniques excel in capturing complex patterns and dependencies within data, thereby improving overall detection accuracy by identifying subtle anomalies that may evade traditional statistical approaches (Chalapathy et al., 2019).

Deep Learning models, including deep neural networks (DNNs) and deep isolation forests, are deployed to address the non-linear relationships and high-dimensional complexities present in financial datasets (Xu et al., 2023). DNNs can learn hierarchical representations of data, enabling them to detect anomalies based on intricate features and relationships across multiple layers (Goodfellow et al., 2014). Deep isolation forests leverage unsupervised learning principles to efficiently isolate anomalies, making them suitable for real-time anomaly detection applications (Liu et al., 2008).

Improving the codebase is crucial for ensuring the system's maintainability, scalability, and performance. By adhering to best practices in software engineering and optimizing algorithms for faster execution, the system can efficiently handle large-scale financial datasets with minimal computational overhead (Pressman, 2014). This approach accelerates anomaly detection and reduces resource consumption, enhancing the system's cost-effectiveness and sustainability over time (Pressman, 2014).

Advanced data structures and caching mechanisms are integrated to expedite data retrieval and processing. Optimized data access strategies, such as in-memory caching and distributed computing frameworks, enable the system to achieve real-time anomaly detection capabilities (Pressman, 2014). These techniques are essential for managing the volume and velocity of financial data streams, ensuring timely insights and responses to emerging anomalies.

Rigorous testing protocols are implemented to validate the reliability and accuracy of the optimized anomaly detection system. This includes unit testing, integration testing across different modules, and performance testing under varying load conditions (Pressman, 2014). By simulating real-world scenarios and stress-testing the system, potential weaknesses or performance bottlenecks can be identified and addressed proactively.

Implementing these enhancements will empower Bank Al-Maghrib to uphold high standards of data quality and integrity in financial analysis (Bank Al Maghrib, 2023; Pressman, 2014). By leveraging advanced ML and DL techniques alongside optimized algorithms and efficient code practices, the proposed solution aims to significantly enhance anomaly detection accuracy, scalability, and responsiveness. This supports the central bank's mandate to ensure financial stability, mitigate risks, and facilitate evidence-based monetary policy decisions.

2.4 Requirements Specification

To ensure the successful implementation of the enhanced anomaly detection system at Bank Al-Maghrib, a comprehensive set of requirements has been identified encompassing functional, non-functional, and regulatory aspects. Functionally, the system must support the detection of anomalies across diverse datasets, leveraging advanced algorithms such as Isolation Forest, DBSCAN, HDBSCAN, and ALAD to improve accuracy and efficiency (Breunig et al., 2000; Liu et al., 2008; McInnes et al., 2017; Zenati et al., 2018). The system should facilitate real-time monitoring and alerting capabilities, enabling prompt identification and response to potential issues. Additionally, it must integrate seamlessly with existing IT infrastructure and support interoperability with other data analysis tools and platforms (Gao et al., 2010).

From a non-functional perspective, the system must demonstrate high performance and scalability, handling large volumes of data without compromising speed or accuracy (Bengio et al., 2003). It should be designed with robust security measures to protect sensitive financial data and ensure compliance with data protection regulations. The system should also exhibit high reliability and availability, with failover mechanisms to maintain continuous operation in the event of hardware or software failures (Gama et al., 2014).

Regulatory requirements necessitate that the system adheres to the stringent standards set by Bank Al-Maghrib and relevant international bodies. This includes ensuring data privacy and security in accordance with GDPR and other applicable laws (Voigt and von dem Bussche, 2017). The system should provide comprehensive audit trails and reporting capabilities to support regulatory compliance and facilitate audits (Ferreira et al., 2016).

Functional requirements

Anomaly Detection Algorithms

The system will incorporate advanced machine learning (ML) and deep learning (DL) algorithms to enhance anomaly detection accuracy and efficiency. This includes ensemble methods, support vector machines (SVM), and isolation forests for robust anomaly identification across diverse financial data sets. ML techniques are chosen for their ability to capture complex data patterns and dependencies, improving detection accuracy beyond traditional statistical methods (Chalapathy et al., 2019).

System Optimization

Optimizing the system's performance, scalability, and responsiveness is critical. This involves refactoring the codebase for improved readability and maintainability, while also implementing advanced data structures and caching mechanisms. These optimizations aim to expedite data retrieval and processing, essential for real-time anomaly detection in dynamic financial environments (Pressman, 2014).

Integration and Testing

Seamless integration of anomaly detection modules and rigorous testing protocols are paramount. Unit testing, integration testing, and performance testing under varying load conditions will validate system reliability and accuracy. Through simulated real-world scenarios, potential weaknesses or bottlenecks in performance can be identified and rectified proactively, ensuring the system meets operational requirements effectively (Pressman, 2014).

User Interface

User interfaces will be designed to facilitate intuitive system monitoring and anomaly visualization. Customizable dashboards will provide real-time anomaly alerts and data insights, enhancing user interaction and decision-making capabilities. Ensuring accessibility and responsiveness across different devices and screen sizes will further optimize user experience and system usability (Nielsen, 1994).

Non-functional requirements

Performance

The system must demonstrate high-performance capabilities to process large-scale financial datasets with minimal latency. Real-time processing is essential for timely anomaly detection and swift response to emerging financial risks (Pressman, 2014).

Scalability

Designing the system to scale seamlessly with increasing data volumes and user demands is imperative. Integration of distributed computing frameworks will support concurrent data processing tasks, maintaining system efficiency and responsiveness during periods of peak activity (Barroso and Hölzle, 2009).

Security

Data privacy and confidentiality of financial information will be safeguarded through secure data handling practices. Implementation of encryption protocols for data transmission and storage will mitigate risks of unauthorized access, ensuring compliance with regulatory standards and safeguarding sensitive financial data (Lynch and McCarthy, 2016).

Reliability

Maintaining high availability and reliability of the anomaly detection system is crucial to minimize operational disruptions. Fault-tolerant mechanisms will be integrated to ensure continuous system operation during potential failures or maintenance activities, enhancing overall system reliability (Pressman, 2014).

Usability

Comprehensive documentation and user guides will be provided to facilitate system administration and user training. Training sessions will familiarize users with system functionalities and anomaly detection tools, promoting effective utilization and maximizing operational efficiency (Nielsen, 1994).

Regulatory and Compliance Requirements

Data Governance

Adherence to regulatory standards and guidelines governing financial data handling and anomaly detection practices is mandatory. Compliance with data protection regulations and audit requirements applicable to financial institutions will be ensured to uphold data integrity and regulatory compliance (Bank Al Maghrib, 2023).

Ethical Considerations

Ethical guidelines will govern the use of ML and DL techniques in financial anomaly detection to ensure transparency and fairness. Mitigating biases and maintaining ethical standards in decision-making processes will be prioritized, fostering trust and confidence in the system's operations (Chalapathy et al., 2019).

2.5 Methodology of work – CRISP-DM

The CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology is selected for this project, renowned for its structured and iterative approach in data mining (Shearer, 2000). It guides projects through six essential phases: Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment (Chapman et al., 2000). This cyclic process ensures a systematic progression from initial business goals to the deployment of actionable insights, allowing for continuous refinement and adaptation to new insights and requirements (Wirth and Hipp, 2000).

CRISP-DM serves as a comprehensive process model tailored to effectively manage data-centric projects. By iterating through its phases, it facilitates continuous improvement and alignment with business objectives (Kurgan and Musilek, 2006). Each phase contributes to enhancing data analysis techniques and models, thereby improving decision-making processes and operational efficiencies (Azevedo and Santos, 2008).

At its core, CRISP-DM prioritizes understanding business goals and requirements before delving into data exploration and model development. The initial phase of Business Understanding involves engaging stakeholders to clearly define the problem statement and establish measurable objectives. This approach ensures that subsequent data mining efforts are purpose-driven and closely aligned with organizational priorities (Shafique and Qaiser, 2014).

Application of Methodology in this Project

The CRISP-DM methodology guided the project from inception to deployment, ensuring a systematic and iterative approach. Key activities undertaken in each phase included:

Business Understanding:

- Conducted Department members' interviews to define the problem of outlier detection in financial data.
- Established project goals, including improving data quality and enhancing decision-making processes.

Data Understanding:

- Collected financial statements and related datasets.
- Performed exploratory data analysis (EDA) to identify patterns, trends, and anomalies.

Data Preparation:

- Cleaned and preprocessed the data using Python and pandas.
- Integrated various data sources and transformed the data into a suitable format for modeling.

Modeling:

- Applied multiple algorithms, including Mahalanobis Distance, Isolation Forest, and Local Outlier Factor, to build outlier detection models.
- Fine-tuned model parameters and evaluated performance.

Evaluation:

- Validated models using cross-validation and performance metrics such as precision, recall, and F1-score.
- Compared model outputs to business objectives and made necessary adjustments.

Deployment:

- Deployed the final model using Pickle and H5, ensuring integration with existing business systems.
- Implemented monitoring tools to continuously assess model performance and make adjustments as needed.

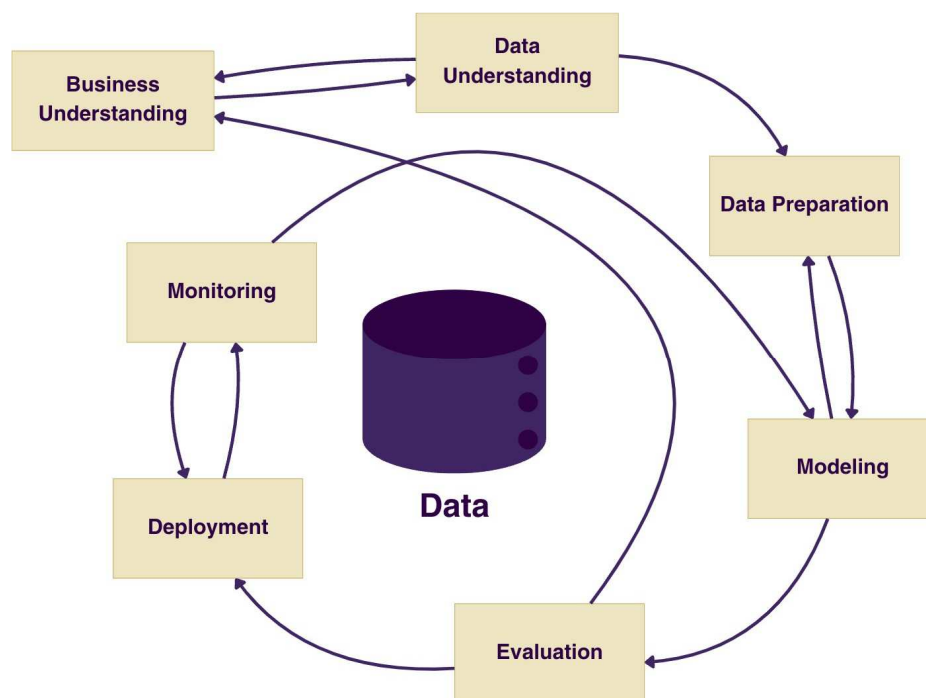


Figure 2.1 – CRISP-DM Chart

2.6 Project Planning

Effective project planning is crucial for the successful execution and delivery of data mining projects. To ensure each phase of the project is executed efficiently and meets the established objectives, a structured approach to planning was employed. This approach encompasses task scheduling, resource allocation, and timeline management.

Central to this planning process was the use of project management principles, which guided a systematic and organized execution of tasks. A key tool in this process is the Gantt chart, which visually represents the project schedule. By highlighting key tasks, their durations, and dependencies, the Gantt chart plays a vital role in tracking progress and maintaining alignment with the project's timeline. This ensures that all phases are completed on schedule and contributes to the overall success of the project.

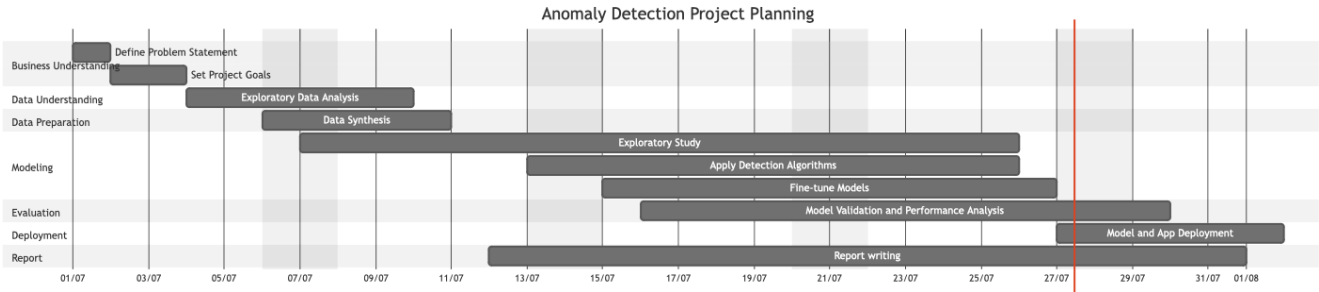


Figure 2.2 – Gantt Diagram

2.7 Conclusion

In conclusion, Chapter 2 has provided a detailed analysis of the challenges and opportunities associated with outlier detection in the financial data ecosystem of Bank Al-Maghrib. By examining the current methodologies, including Isolation Forest and Mahalanobis Distance, we have identified key limitations in terms of performance, accuracy, and scalability. These limitations underscore the need for enhanced outlier detection methods that can effectively manage the complexities of modern financial datasets.

The critical evaluation of existing methods has laid the foundation for proposing a multifaceted solution aimed at optimizing outlier detection systems. This proposed solution includes algorithmic refinements, technological upgrades, and a methodological framework anchored in CRISP-DM. By refining algorithms and adopting advanced data structures, we aim to improve anomaly detection accuracy and efficiency. The methodological framework of CRISP-DM ensures a structured approach from data understanding to deployment, facilitating iterative improvements and alignment with business objectives.

Moving forward, Chapters 3, 4, and 5 will build upon this foundational analysis and proposed solution. Chapter 3 will delve into the theoretical underpinnings of outlier detection techniques, providing a robust framework for methodological development. Chapters 4 and 5 will detail the implementation of the proposed solution, including algorithm optimization, model development, and integration into Bank Al-Maghrib's operational environment.

By addressing the identified challenges and leveraging advanced methodologies, this thesis aims to enhance the reliability, scalability, and efficiency of outlier detection in financial data analysis. The insights gained from Chapter 2 will guide subsequent chapters in realizing these objectives, ensuring that our solutions are not only theoretically sound but also practical and impactful within the context of Bank Al-Maghrib's operational landscape.

Chapter 3

Theoretical Background

Building on the problem analysis from the previous chapter, this delves into the theoretical underpinnings of outlier detection and the various techniques used in this domain by providing the necessary theoretical foundation for understanding and implementing advanced outlier detection methods. The transition from identifying the problem to exploring potential solutions requires a deep understanding of existing methodologies and their theoretical bases.

The chapter begins with a historical background, tracing the evolution of outlier detection techniques and their significance in various fields, particularly in financial data quality. Understanding the historical context helps in appreciating the advancements and current trends in outlier detection research. We then discuss the importance of outlier detection in maintaining financial data quality, highlighting how accurate detection and handling of outliers can prevent erroneous financial reporting and enhance decision-making processes.

The chapter then surveys cutting-edge techniques for outlier detection, including statistical-based, distance-based, density-based, and machine learning approaches.

Each approach is explored in detail:

- **Statistical-Based Approaches:** These methods, including Z-score and Grubbs' test, rely on statistical properties of the data to identify outliers.
- **Distance-Based Approaches:** Techniques such as the Mahalanobis Distance measure the distance of each data point from the center of the data distribution, identifying those that lie far away as potential outliers.
- **Density-Based Approaches:** Methods like Local Outlier Factor (LOF) identify outliers based on the density of data points in their neighborhood.
- **Machine Learning Approaches:** Advanced techniques like Isolation Forest and deep learning methods, including Adversarially Learned Anomaly Detection (ALAD), are explored for their ability to detect complex, non-linear patterns in large datasets.

Special emphasis is placed on methods such as Data Synthesis, the Mahalanobis Distance, Local Outlier Factor, Isolation Forest, and deep learning techniques like GANs and ALAD. Each of these methods offers unique advantages and is suited to different types of data and applications. This

comprehensive review provides the theoretical foundation for the methodologies employed in the study, ensuring that our approach is grounded in the latest research and best practices.

The theoretical insights gained in this chapter will inform the solution engineering discussed in Chapter 4. By understanding the strengths and limitations of various outlier detection techniques, we can make informed decisions about the methods and algorithms that will be implemented and optimized in our study.

The concept of outlier detection, also referred to as anomaly detection, has a rich history deeply rooted in the fields of statistics and data analysis. Its origins can be traced back to the early 19th century, with significant contributions from renowned mathematicians and statisticians. Notable early works include those by Francis Galton and later advancements by John Tukey in the 20th century, who emphasized the importance of identifying unusual data points in exploratory data analysis (Galton, 1886; Tukey, 1977).

In modern times, the field has expanded significantly, incorporating state-of-the-art methods and techniques. Contemporary research has provided comprehensive surveys and overviews of progress in outlier detection methods, emphasizing the evolution and advancements in this area. Chalapathy et al. (2019) and Wang et al. (2019) provide extensive reviews of these techniques, highlighting the shift from simple statistical tests to complex machine learning and deep learning algorithms.

3.1 Background History

Outlier detection, a crucial aspect of data analysis, has evolved significantly since its inception. Initially, the primary focus was on identifying and managing anomalies in datasets to ensure the accuracy and reliability of statistical analyses. The origins of outlier detection can be traced back to the early works of Francis Galton in the 19th century, who introduced the concept of "statistical outliers" in the context of normal distribution (Galton, 1886).

In the early 20th century, statisticians such as John Tukey further developed these ideas, introducing exploratory data analysis techniques that emphasized the importance of identifying unusual data points (Tukey, 1977). The development of robust statistical methods in the mid-20th century, such as Grubbs' test and the Z-score, provided more systematic approaches to detect outliers based on statistical properties of the data (Grubbs, 1969; Barnett and Lewis, 1994).

The advent of computer technology in the latter half of the 20th century significantly expanded the scope and capabilities of outlier detection techniques. Distance-based methods, such as the Mahalanobis Distance, emerged as powerful tools for identifying multivariate outliers by measuring the distance of each data point from the center of the data distribution (Mahalanobis, 1936). These methods were further refined and complemented by density-based approaches, such as the Local Outlier Factor (LOF), which assesses the density of data points in their neighborhood to identify anomalies (Breunig et al., 2000).

In the 21st century, the rise of big data and machine learning has revolutionized outlier detection. Advanced techniques, such as Isolation Forest, leverage machine learning algorithms to isolate anomalies by randomly partitioning data (Liu et al., 2008). Deep learning methods, including Generative Adversarial Networks (GANs) and Adversarially Learned Anomaly Detection (ALAD), have further pushed the boundaries of outlier detection by capturing complex, non-linear patterns in large datasets (Zenati et al., 2018).

Understanding the historical context of outlier detection helps appreciate the advancements and current trends in this field. Over the years, the focus has shifted from simple statistical tests to sophisticated machine learning algorithms, reflecting the increasing complexity and scale of modern datasets.

Applications of Outlier Detection in Real Life

Outlier detection has numerous applications across various industries, proving its importance beyond just theoretical interest. In finance, it is crucial for identifying fraudulent transactions, market manipulation, and unusual trading activities. By detecting anomalies in financial transactions, institutions can prevent fraud, ensure compliance with regulations, and maintain market integrity (Bolton and Hand, 2002).

In healthcare, outlier detection helps in identifying unusual patterns in patient data, which can indicate medical errors, billing fraud, or outbreaks of diseases. It aids in early detection of diseases and anomalies in medical imaging, contributing to better patient care and resource management (Chandola et al., 2009).

Manufacturing industries use outlier detection for quality control and predictive maintenance. By identifying deviations from normal operation patterns, companies can predict equipment failures and take preventive measures to avoid costly downtimes (Mourtzis et al., 2016).

In cybersecurity, detecting outliers is fundamental for identifying potential security breaches and network intrusions. By monitoring network traffic and system logs for unusual activities, organizations can safeguard against attacks and protect sensitive information (Garcia-Teodoro et al., 2009).

The wide range of applications underscores the versatility and necessity of robust outlier detection methods across different sectors, each with its own unique challenges and requirements.

Significance of Outlier Detection in Financial Data Quality

Outlier detection plays a pivotal role in maintaining the quality of financial data. Accurate detection and handling of outliers are essential for preventing erroneous financial reporting, which can have significant consequences for decision-making processes and regulatory compliance (Aggarwal, 2017). In the financial sector, outliers may indicate fraudulent activities, accounting errors, or unusual market conditions, making their timely identification crucial for maintaining the integrity and reliability of financial analyses (Hodge and Austin, 2004).

Given this context, outlier detection methods are critical for ensuring that financial data is accurate and reliable. By identifying and addressing anomalies, organizations can mitigate risks, improve financial reporting accuracy, and comply with regulatory standards (Chandola et al., 2009). For institutions like Bank Al-Maghrib, implementing an effective outlier detection system is vital for safeguarding data integrity and supporting sound financial management (Bank Al-Maghrib, 2021).

High-quality data is paramount in finance because institutions rely on vast amounts of information to drive investment strategies, risk management, and regulatory compliance (Han et al., 2011). Outliers in financial data can distort analytical models, leading to inaccurate predictions and faulty decision-making. Consequently, identifying and addressing these anomalies enhances the reliability and accuracy of data, leading to better-informed decisions and improved financial outcomes (Aggarwal, 2015).

Furthermore, outlier detection enhances data-driven decision-making processes in finance. By identifying unusual patterns and deviations from the norm, financial institutions can gain valuable insights, uncover hidden trends, and proactively address potential issues (Hodge and Austin, 2004). This capability leads to better risk management, improved fraud detection, and more accurate forecasting (Chandola et al., 2009).

The significance of outlier detection is underscored by its historical evolution. From early astro-

nomical observations to modern deep learning techniques, the continuous efforts of researchers and practitioners reflect the growing complexity of data and the need for robust and reliable methods to ensure data quality and integrity (Chalapathy and Chawla, 2019). This rich history demonstrates the enduring importance of developing effective outlier detection methods to meet the demands of increasingly complex financial data environments (Wang et al., 2019).

3.2 Distance-Based Approaches

Distance-based approaches are a fundamental category of outlier detection techniques grounded in the theoretical concept of measuring the dissimilarity between data points. At its core, distance-based methods rely on the idea that outliers are data points that deviate significantly from the majority of the data. The notion of distance in this context is mathematically defined as a measure of the spatial separation between two elements within a data space.

Formally, the distance between two elements x and y in an n -dimensional space can be described by a distance metric or function. A distance metric $d(x, y)$ quantifies how far apart the two points are, adhering to specific mathematical properties such as non-negativity, identity of indiscernibles, symmetry, and the triangle inequality (Munkres, 2000). These properties ensure that the distance function provides a meaningful measurement of dissimilarity.

Distance-based approaches exploit these distance metrics to detect outliers by identifying points that lie far away from the majority of the data points. This assumption is based on the idea that normal data points tend to cluster together, whereas outliers are isolated or sparse. By evaluating how far each point is from its neighbors or from a central reference point, distance-based methods classify points with large distances as outliers.

In this section, we will explore several prominent distance-based approaches used in outlier detection, including Euclidean distance, Mahalanobis distance, and other relevant metrics. Each approach is defined by its specific method for calculating distance and is suited to different types of data and applications. Understanding these methods and their theoretical foundations is crucial for selecting and implementing effective outlier detection techniques.

3.2.1 Euclidean Distance

Euclidean distance is one of the most fundamental and widely utilized distance metrics in the field of outlier detection. It measures the straight-line distance between two points in a multi-dimensional space, providing a simple yet powerful means of quantifying the dissimilarity between data points. The Euclidean distance between two points $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$ in an n -dimensional space is defined mathematically by the formula:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

This formula calculates the length of the shortest path between the two points, which is a straight line in Euclidean space. The concept of Euclidean distance is rooted in the principles of Euclidean geometry, which were developed by the ancient Greek mathematician Euclid (Euclid, 300 BC).

Euclidean distance adheres to several key properties of a metric, including:

1. Non-Negativity: $d(x, y) \geq 0$
2. Identity of Indiscernibles: $d(x, y) = 0$ if and only if $x = y$
3. Symmetry: $d(x, y) = d(y, x)$
4. Triangle Inequality: $d(x, y) + d(y, z) \geq d(x, z)$

These properties ensure that Euclidean distance is a consistent and meaningful measure of dissimilarity between points in a given space (Munkres, 2000).

In the context of outlier detection, Euclidean distance is employed to identify data points that are significantly distant from their neighbors. Points with large Euclidean distances relative to the majority of other points are flagged as outliers. This method is particularly effective in low-dimensional spaces where data points are clustered around a central point. However, its effectiveness diminishes in high-dimensional spaces due to the "curse of dimensionality", where distances between points become less distinguishable and meaningful (Beyer et al., 1999).

A study by Tanger (2018) highlights the utility of Euclidean distance for outlier detection, emphasizing its simplicity and effectiveness in various practical scenarios. Tanger argues that while Euclidean distance may be straightforward, its application in outlier detection remains robust, especially when combined with other techniques to address its limitations in high-dimensional settings (Tanger, 2018).

Despite its advantages, Euclidean distance has certain limitations. In high-dimensional datasets, the distance between points can become less informative due to the sparsity of the data. Additionally, when dimensions have different scales or units, the Euclidean distance may produce misleading results, necessitating normalization or dimensionality reduction techniques to improve its performance (Jain et al., 1999).

3.2.2 Mahalanobis Distance

Mahalanobis distance is a statistical measure used to determine the distance between a point and a distribution, considering the correlations between variables in the dataset. Unlike Euclidean distance, which measures the straight-line distance between points in space, Mahalanobis distance takes into account the variance and covariance of the data, making it particularly useful for identifying outliers in multivariate datasets.

Mathematically, the Mahalanobis distance between a point x and a distribution with mean μ and covariance matrix Σ is defined as:

$$d_M(x, \mu) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

Here, x is the data point under consideration, μ is the mean vector of the distribution, Σ is the covariance matrix, and Σ^{-1} is the inverse of the covariance matrix. The term $(x - \mu)^T \Sigma^{-1} (x - \mu)$ computes the squared distance of x from μ , adjusted for the correlations between the variables.

Mahalanobis distance is a significant improvement over Euclidean distance in the context of multivariate data because it accounts for the distribution of the data. While Euclidean distance assumes that all dimensions are equally important and independent, Mahalanobis distance incorporates the shape of the data distribution, providing a more accurate measure of how far a point is from the mean, given the data's variance and correlations (Mahalanobis, 1936).

In practice, Mahalanobis distance is particularly useful for detecting outliers in datasets where variables are correlated. For instance, in financial data, where returns of different assets might be

correlated, Mahalanobis distance helps identify assets whose returns significantly deviate from their expected patterns considering the correlations among assets (Dudoit and Fridlyand, 2003).

A study by Iglewicz and Hoaglin (1993) demonstrates the effectiveness of Mahalanobis distance in identifying outliers in multivariate datasets, noting its ability to detect outliers that might be overlooked by methods that do not consider the correlation structure (Iglewicz and Hoaglin, 1993). However, Mahalanobis distance requires an accurate estimation of the covariance matrix, which can be challenging in high-dimensional spaces where the covariance matrix may be singular or poorly conditioned (Chen et al., 2001).

Mahalanobis distance offers a robust method for outlier detection in multivariate settings by incorporating the distribution's statistical properties. Its ability to account for correlations between variables makes it a valuable tool in various applications, including financial analysis and quality control.

3.2.3 Minkowski Distance

Minkowski distance is a generalized metric used to measure the distance between two points in a multi-dimensional space. It extends the concept of Euclidean and Manhattan distances by introducing a parameter that allows for flexibility in the distance measurement. The Minkowski distance between two points $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$ in n -dimensional space is defined as:

$$d_p(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

Here, p is a parameter that determines the type of distance measurement: - For $p = 1$, the Minkowski distance simplifies to the Manhattan distance, which measures the sum of the absolute differences of their coordinates. - For $p = 2$, it becomes the Euclidean distance, which measures the straight-line distance between the points. - For $p \rightarrow \infty$, it approaches the Chebyshev distance, which measures the maximum absolute difference between the coordinates.

Minkowski distance generalizes both Euclidean and Manhattan distances, making it a versatile tool for various applications in outlier detection and data analysis. The parameter p allows users to adjust the sensitivity of the distance metric to different types of data and anomalies (Minkowski, 1907).

In outlier detection, Minkowski distance provides a flexible approach for identifying outliers by adjusting the parameter p based on the specific characteristics of the dataset. For example, in cases where outliers are expected to exhibit more pronounced deviations in specific dimensions, choosing a higher p value can amplify the impact of these deviations (Shao et al., 2016). Conversely, a lower p value may be more appropriate for datasets where deviations are more uniform across dimensions.

A study by Shao et al. (2016) explores the use of Minkowski distance for outlier detection, highlighting its ability to accommodate various distance metrics through the choice of p . This flexibility allows for customization of the distance metric to better suit the characteristics of the data and the nature of the anomalies (Shao et al., 2016). However, the choice of p can significantly influence the results, and selecting an appropriate value requires careful consideration of the data distribution and the specific objectives of the analysis.

Minkowski distance is a powerful and adaptable metric for outlier detection, offering the ability

to fine-tune the distance measurement according to the needs of the dataset. Its generalization of both Euclidean and Manhattan distances makes it a valuable tool for a wide range of applications in data analysis.

3.2.4 Chebyshev Distance

Chebyshev distance, also known as maximum metric or L_∞ norm, is a distance metric that measures the greatest of the absolute differences between corresponding coordinates of two points. It is particularly useful in contexts where the most significant deviation in any single dimension should be highlighted.

Mathematically, the Chebyshev distance between two points $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$ in n -dimensional space is defined as:

$$d_C(x, y) = \max_{i=1, \dots, n} |x_i - y_i|$$

In this formula, $|x_i - y_i|$ represents the absolute difference between the coordinates of x and y in the i -th dimension. The Chebyshev distance is thus the maximum of these differences, effectively capturing the largest deviation in any single dimension.

The Chebyshev distance is named after the Russian mathematician Pafnuty Chebyshev, who first introduced it in the 19th century (Chebyshev, 1962). This distance metric is particularly useful in scenarios where deviations in any one dimension are more critical than deviations in the overall distance. It is commonly used in chessboard distance calculations, where the goal is to determine the minimum number of moves a king needs to move from one square to another, considering that the king can move horizontally, vertically, or diagonally (Chebyshev, 1962).

In outlier detection, Chebyshev distance can be advantageous when anomalies are expected to exhibit extreme deviations in a single feature or dimension. For instance, in certain types of quality control processes, a single outlier in a crucial dimension might be more significant than smaller deviations across multiple dimensions. By focusing on the maximum deviation, Chebyshev distance allows for effective identification of such anomalies (Tukey, 1977).

A study by Tukey (1977) emphasizes the usefulness of Chebyshev distance in identifying outliers in high-dimensional datasets, where it can help isolate extreme deviations that might be missed by other distance metrics. The simplicity of the Chebyshev distance makes it computationally efficient, though it may not always capture the complexity of multidimensional anomalies as effectively as other metrics (Tukey, 1977).

Chebyshev distance is a straightforward and effective metric for highlighting extreme deviations in any single dimension of a dataset. Its focus on the maximum difference makes it particularly suitable for applications where the most significant deviation is of primary concern.

While distance-based approaches identify outliers by measuring the spatial distance between points, density-based approaches focus on the local density of data points. Unlike distance-based methods, which can struggle with varying data scales and high dimensions, density-based techniques excel in detecting outliers in regions with low data density. This shift from distance-based to density-based methods reflects a progression towards more nuanced and robust outlier detection strategies.

3.3 Density-Based Approaches

Density-based approaches to outlier detection are grounded in the concept of data density, which refers to the concentration of data points in a given region of the feature space. Unlike distance-based methods that focus solely on the spatial separation between points, density-based techniques assess the relative density of data points to identify outliers. The core idea is that outliers are data points that lie in regions with significantly lower density compared to their surrounding neighborhoods.

In a statistical context, density refers to the distribution of data points within a specific area of the feature space. For a given point x , its density can be quantified by counting the number of neighboring points within a predefined radius or by estimating the probability density function (PDF) of the data distribution around x . Mathematically, the density $\rho(x)$ at a point x is often defined as:

$$\rho(x) = \frac{1}{|N(x)|} \sum_{y \in N(x)} K(\|x - y\|)$$

where $N(x)$ is the set of neighbors within a certain radius of x , $K(\cdot)$ is a kernel function (e.g., Gaussian kernel), and $|N(x)|$ is the number of neighbors. This formula provides an estimate of the local density around x based on the proximity of neighboring data points (Parzen, 1962).

The fundamental assumption behind density-based outlier detection is that outliers typically reside in regions of the feature space that are sparsely populated with data points. In contrast, inlier points are clustered in high-density regions. By evaluating the density of each point's neighborhood, density-based methods can effectively distinguish between points that are part of a dense cluster and those that are isolated in sparser regions.

Density-based methods are advantageous because they do not require specifying the number of clusters or outliers beforehand, and they are robust to varying densities within the dataset. This makes them particularly useful in scenarios where the data distribution is uneven or contains clusters of different sizes and shapes (Ester et al., 1996).

In outlier detection, density-based approaches are used to identify points that have significantly lower density compared to their neighbors. These points are flagged as potential outliers due to their isolation in sparse regions of the feature space. In this section, we will explore several commonly used density-based algorithms:

- **Local Outlier Factor (LOF):** LOF evaluates the local density deviation of a point relative to its neighbors. Points with a significantly lower density compared to their neighbors are flagged as outliers (Breunig et al., 2000).
- **Simplified Local Outlier Factor (sLOF):** Simplified Local Outlier Factor (sLOF) is a variant of the traditional Local Outlier Factor (LOF) that aims to improve computational efficiency. By optimizing the algorithm's processing steps, sLOF reduces the complexity involved in calculating local density deviations while retaining its core capability to detect outliers (Schubert et al., 2017).
- **RS-Forest:** RS-Forest is a density-based method that uses a forest of random subspace models to estimate the density of data points. This technique aggregates density estimates from multiple random projections to identify outliers effectively (Xia et al., 2015).
- **Density-Based Spatial Clustering of Applications with Noise (DBSCAN) and HDBSCAN:** DBSCAN clusters data based on density, identifying outliers as points in low-density regions that do not belong to any cluster (Ester et al., 1996).

- **Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN):** HDBSCAN extends DBSCAN by incorporating hierarchical clustering, allowing it to handle data with varying densities and complex cluster structures more effectively (McInnes et al., 2017).

These density-based methods offer robust techniques for outlier detection, particularly in datasets with varying density and complex cluster structures, enhancing the accuracy and reliability of anomaly detection processes.

3.3.1 Local Outlier Factor (LOF)

The Local Outlier Factor (LOF) is a well-established density-based method for outlier detection that identifies anomalies based on the local density of data points. Introduced by Breunig et al. (2000), LOF provides a robust mechanism for detecting outliers by comparing the density of a data point with the density of its neighbors.

At the core of LOF is the concept of local density. For a given data point, its local density is defined relative to its neighbors. This relative density is computed using a measure known as the reachability distance. Specifically, the reachability distance between a data point p and its neighbor o is the maximum of the distance between p and o and the distance between o and the k -th nearest neighbor of o .

Mathematically, the reachability distance $reach_dist_k(p, o)$ can be expressed as:

$$reach_dist_k(p, o) = \max(d(p, o), d(o, kNN(o)))$$

where $d(p, o)$ denotes the distance between points p and o , and $kNN(o)$ represents the k -th nearest neighbor of point o .

The local density of a point p is then defined as the inverse of the average reachability distance of p with respect to its k -nearest neighbors:

$$local_density_k(p) = \frac{1}{|N_k(p)|} \frac{1}{\sum_{o \in N_k(p)} reach_dist_k(p, o)}$$

where $N_k(p)$ denotes the set of k -nearest neighbors of p .

LOF evaluates the outlier score by comparing the local density of a point with the local densities of its neighbors.

Specifically, the Local Outlier Factor (LOF) for a point p is defined as:

$$LOF_k(p) = \frac{1}{|N_k(p)|} \frac{\sum_{o \in N_k(p)} local_density_k(o)}{local_density_k(p)}$$

A higher LOF score indicates that the data point p is an outlier, as it is situated in a region with significantly lower density compared to its neighbors.

LOF's primary advantage is its ability to identify outliers based on local data density, which makes it effective in datasets with varying density regions. This method is particularly useful in

detecting anomalies in complex datasets where the density of data points varies spatially.

In practice, LOF has been successfully applied across various domains, including fraud detection and network security, where it helps uncover anomalous patterns that deviate from the expected local data distribution (Breunig et al., 2000).

3.3.2 Simplified Local Outlier Factor (sLOF)

The Simplified Local Outlier Factor (sLOF) is an optimized version of the Local Outlier Factor (LOF) designed to improve computational efficiency while preserving the core outlier detection capabilities of the original method. LOF, introduced by Breunig et al. in 2000, measures the local density deviation of a data point relative to its neighbors, identifying points that are significantly less dense as outliers. The sLOF method refines this approach by streamlining the computations involved in calculating density and outlier scores, making it more suitable for large-scale datasets.

The fundamental idea behind sLOF is to assess how isolated a data point is with respect to its local neighborhood. The LOF score for a data point p is computed based on its local density compared to that of its neighbors. This local density is quantified using the concept of the k-nearest neighbors (k-NN). The LOF score $LOF(p)$ for a point p is defined as:

$$LOF(p) = \frac{\frac{\sum_{o \in N_k(p)} lrd(o)}{lrd(p)}}{|N_k(p)|}$$

where $N_k(p)$ represents the set of k-nearest neighbors of point p , and $lrd(p)$ denotes the local reachability density of point p . The local reachability density is calculated as:

$$lrd(p) = \frac{1}{reach_dist(p)}$$

where the reachability distance $reach_dist(p)$ is given by:

$$reach_dist(p) = \max(dist(p, o), k - dist(o))$$

Here, $dist(p, o)$ is the distance between point p and its neighbor o , and $k - dist(o)$ is the distance from o to its k-th nearest neighbor.

The sLOF algorithm simplifies the original LOF by reducing the computational overhead involved in distance calculations and density estimations. This is achieved by approximating the density measures and optimizing the k-nearest neighbor search process. The result is a more computationally efficient outlier detection method that retains the effectiveness of the LOF algorithm.

The sLOF method is particularly useful in scenarios with large datasets or real-time applications where computational efficiency is critical. By maintaining a balance between accuracy and performance, sLOF provides a practical solution for detecting outliers in complex and high-dimensional data.

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

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is a widely used density-based clustering algorithm that also serves as an effective method for outlier detection. Introduced

by Ester et al. (1996), DBSCAN identifies clusters based on the density of data points in a given region and labels points in low-density areas as noise or outliers.

The core idea of DBSCAN is to classify points based on their density relative to their neighbors. The algorithm relies on two key parameters: ϵ (epsilon) and *minPts* (minimum points). The parameter ϵ defines the radius within which the algorithm searches for neighboring points, while *minPts* specifies the minimum number of points required to form a dense region, or cluster.

In DBSCAN, a point is classified into one of three categories:

1. **Core Point:** A point is classified as a core point if it has at least *minPts* neighbors within the ϵ radius. Core points are essential for forming a cluster.
2. **Border Point:** A point is classified as a border point if it is within the ϵ radius of a core point but does not have enough neighbors to be a core point itself.
3. **Noise Point:** A point is classified as noise if it does not belong to any cluster, meaning it is neither a core point nor a border point.

Mathematically, the core idea of DBSCAN can be expressed as follows. For a given point p , the neighborhood of p within radius ϵ is defined as:

$$N_\epsilon(p) = \{q \in D \mid d(p, q) \leq \epsilon\}$$

where D denotes the dataset and $d(p, q)$ is the distance between points p and q . The point p is considered a core point if:

$$|N_\epsilon(p)| \geq \text{minPts}$$

where $|N_\epsilon(p)|$ represents the number of points within the ϵ radius of p .

DBSCAN's advantage lies in its ability to find arbitrarily shaped clusters and its robustness to outliers. By identifying dense regions in the data and separating them from sparse regions, DBSCAN can effectively handle noise and anomalies. This makes it particularly useful in applications where the data exhibits complex structures and varying densities, such as geographical data analysis and image segmentation (Ester et al., 1996).

The algorithm's effectiveness in identifying outliers stems from its ability to detect points that do not belong to any dense region. Such points, identified as noise, are effectively flagged as outliers, making DBSCAN a powerful tool for outlier detection in diverse datasets.

3.3.4 Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN)

Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) is an advanced extension of the DBSCAN algorithm, designed to improve clustering performance and outlier detection in datasets with varying densities. Proposed by McInnes et al. (2017), HDBSCAN builds upon the foundational principles of DBSCAN by incorporating hierarchical clustering techniques to enhance the algorithm's flexibility and robustness.

HDBSCAN operates in two main stages: hierarchical clustering and density-based clustering. The first stage constructs a hierarchy of clusters using the concept of mutual reachability distance, which generalizes the notion of density-based clustering to include varying densities. Mutual reachability

distance between two points p and q is defined as:

$$reachability_dist(p, q) = \max\{core_dist(p), core_dist(q), d(p, q)\}$$

where $core_dist(p)$ represents the core distance of point p (i.e., the distance to its $minPts$ -th nearest neighbor), and $d(p, q)$ is the Euclidean distance between p and q . This distance metric ensures that the reachability distance takes into account both the local density of points and the distance between them, allowing for more nuanced cluster formation.

In the second stage, HDBSCAN performs a hierarchical clustering analysis using the mutual reachability distance. This hierarchical tree, or dendrogram, represents a nested series of clusters. To obtain a flat clustering from the hierarchical structure, HDBSCAN uses a concept called "cluster stability" to select the most meaningful clusters. Stability is measured by how well a cluster persists across different levels of the hierarchy. Clusters with higher stability are considered more significant and are retained in the final clustering result.

One of the key advantages of HDBSCAN is its ability to identify clusters of varying shapes and densities, as well as its robustness to noise. Unlike DBSCAN, which requires specifying a single ϵ parameter for clustering, HDBSCAN does not require pre-defining the number of clusters and can automatically determine the optimal number of clusters based on the stability of the hierarchical structure. This adaptability makes HDBSCAN particularly effective for complex datasets where clusters exhibit varying densities and shapes (McInnes et al., 2017).

In addition to its clustering capabilities, HDBSCAN excels in outlier detection. Points that do not belong to any stable cluster are classified as outliers or noise. This classification is based on their exclusion from significant clusters within the hierarchy, allowing HDBSCAN to robustly detect anomalies in diverse data distributions.

Overall, HDBSCAN enhances the ability to perform density-based clustering and outlier detection by incorporating hierarchical techniques, making it a powerful tool for analyzing complex datasets with varying densities and structures.

3.3.5 RS-Forest (Random Sample Consensus Forest)

The RS-Forest algorithm, or Random Sample Consensus Forest, is a sophisticated method for outlier detection that builds on the principles of ensemble learning and Random Sample Consensus (RANSAC). RANSAC, introduced by Fischler and Bolles in 1981, is a robust statistical method designed to estimate parameters of a mathematical model from a dataset that contains outliers. The core idea of RANSAC is to iteratively select random subsets of the data, fit a model to these subsets, and identify the subset that best fits the model while minimizing the influence of outliers.

In the context of RS-Forest, this principle is extended to decision tree ensembles. RS-Forest constructs an ensemble of decision trees, each trained on a different random sample of the dataset. By doing so, the method introduces diversity into the model, which helps reduce the risk of overfitting and improves its ability to generalize across different data patterns. Each decision tree in the forest is trained to classify data points based on their features, distinguishing between normal data points (inliers) and anomalous data points (outliers).

The detection of outliers in RS-Forest is achieved through a consensus mechanism across the ensemble of decision trees. For each data point, a consensus score is calculated based on how frequently the point is classified as an outlier by the various trees in the forest.

The consensus score $S(x)$ for a data point x is computed as:

$$S(x) = \frac{1}{T} \sum_{i=1}^T outlier_score_i(x)$$

where T is the number of trees in the forest, and $outlier_score_i(x)$ is the outlier score given by the i -th tree. This score indicates whether the point x is considered an outlier by that tree.

The consensus score aggregates the decisions of multiple trees, providing a robust measure of whether a data point is an outlier. Points with higher consensus scores are more consistently classified as outliers across the ensemble, whereas points with lower scores are classified as inliers. This approach enhances the accuracy and robustness of outlier detection by mitigating the influence of individual tree biases and accommodating diverse data distributions.

RS-Forest leverages the RANSAC principle by using random subsets of the data to train its ensemble of decision trees, which helps improve its resistance to noisy data and complex data patterns. This makes RS-Forest a valuable tool for detecting anomalies in diverse and challenging datasets.

Transitioning from density-based approaches, which identify outliers based on local data density, we now shift to machine learning-based methods. These techniques leverage advanced algorithms to learn complex patterns and detect anomalies, offering enhanced capabilities for high-dimensional and non-linearly separable data. Machine learning methods build upon the strengths of density-based approaches by providing more adaptive and scalable solutions for outlier detection.

3.4 Machine Learning-Based Approaches

Machine learning-based approaches to outlier detection leverage sophisticated algorithms to identify anomalies by learning patterns and relationships within the data. Unlike traditional statistical or distance-based methods, these techniques adaptively learn from the data, enabling them to detect complex and subtle outliers that may not be captured by simpler models (Chandola et al., 2009).

At the core of machine learning-based approaches is the ability to handle high-dimensional data and uncover non-linear relationships between features. These methods often involve training models on a dataset to distinguish between normal and anomalous instances. The models then apply the learned patterns to new, unseen data to identify potential outliers. This adaptability makes machine learning-based methods particularly effective in dynamic and complex datasets where traditional methods may fall short (Hodge and Austin, 2004).

A variety of machine learning techniques are employed for outlier detection, each with its own strengths and application scenarios. For instance, supervised methods require labeled data to train models that classify instances as normal or outlier. Conversely, unsupervised methods, which do not rely on labeled data, detect anomalies based on deviations from learned patterns or distributions (Xia et al., 2015). Additionally, semi-supervised methods use a combination of labeled and unlabeled data to improve detection performance (Jin et al., 2006).

Overall, machine learning-based approaches enhance the capability to identify outliers by leveraging advanced computational models and learning algorithms, thus offering more robust and scalable solutions for complex data analysis (Ahmed et al., 2016).

3.4.1 Isolation Forest (iForest)

Isolation Forest (iForest) is a machine learning algorithm specifically designed for anomaly detection in high-dimensional datasets. It operates on the principle of isolation, which involves isolating observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature. The main idea behind iForest is that outliers are easier to isolate compared to normal points because they lie far from the majority of data points.

The algorithm constructs an ensemble of isolation trees (iTrees), where each iTree is built by recursively partitioning the data. The process involves selecting a feature and then a split value, which partitions the data into two subsets. This process continues until each subset contains only one data point or until a predefined height is reached. The height of the tree, which represents the number of splits required to isolate a point, is used to measure the degree of isolation.

Mathematically, the isolation score $s(x)$ for a data point x can be expressed as:

$$s(x) = \frac{2^{\frac{H(x)}{c(n)}} - 1}{c(n) - 1}$$

where $H(x)$ is the average path length of point x across all trees, $c(n)$ is the average path length of unsuccessful searches in a binary search tree, and n is the number of data points. The average path length $c(n)$ is defined as:

$$c(n) = 2 \cdot \left(\frac{n-1}{2} - \frac{2(n-1)}{n} \right)$$

The isolation score $s(x)$ is higher for anomalies because they are isolated closer to the root of the trees, resulting in shorter average path lengths. Conversely, normal data points are more frequently split and thus have longer average path lengths.

iForest is particularly effective due to its efficiency in handling large datasets with high dimensionality. The algorithm's computational complexity is $O(n \log n)$, which allows it to scale well with large datasets compared to traditional outlier detection methods (Liu et al., 2008). This efficiency, combined with its ability to detect anomalies without requiring extensive parameter tuning, makes iForest a popular choice for many anomaly detection applications.

3.4.2 One-Class Support Vector Machine

Another approach to consider in machine learning-based outlier detection is the One-Class Support Vector Machine (One-Class SVM). One-Class SVM is an adaptation of the traditional Support Vector Machine (SVM), designed specifically for anomaly detection. This method is widely used due to its effectiveness in high-dimensional spaces and its ability to handle non-linear data distributions.

The fundamental principle behind One-Class SVM is to learn a decision function for outlier detection by identifying regions in the input space where the data is concentrated. The algorithm is trained only on the normal data, and it aims to separate the normal data points from the origin in the feature space using a hyperplane. The decision function $f(x)$ learned by One-Class SVM assigns a positive value to the data points that are considered normal and negative values to the anomalies.

Mathematically, One-Class SVM involves solving the following optimization problem:

$$\min_{w, \rho, \xi_i} \left(\frac{1}{2} \|w\|^2 + \frac{1}{vn} \sum_{i=1}^L \xi_i - \rho \right)$$

subject to

$$(w \cdot \phi(x_i)) \geq \rho - \xi_i, \quad \xi_i \geq 0$$

Here, w is the normal vector to the hyperplane, ρ is the offset, ξ_i are the slack variables to allow for some data points to lie within the margin, $\phi(x_i)$ is a mapping function that transforms the data into a higher-dimensional feature space, and v is a parameter that controls the trade-off between the margin size and the number of outliers.

One-Class SVM employs the kernel trick to handle non-linear data. Commonly used kernels include the Radial Basis Function (RBF) kernel, which is defined as:

$$K(x_i, x_j) = \exp \left(-\frac{\|x_i - x_j\|^2}{2\sigma^2} \right)$$

where σ is a parameter that determines the spread of the kernel.

The training process involves finding the optimal hyperplane that best separates the normal data points from the origin in the transformed feature space. During testing, the decision function $f(x)$ is used to classify new data points. Points with $f(x) < 0$ are considered outliers.

One-Class SVM is particularly robust in scenarios where the normal data is well-separated from anomalies, making it suitable for a wide range of applications, including fraud detection, network security, and quality control.

3.4.3 Autoencoders

Autoencoders are a type of neural network that are trained to encode input data into a compressed representation and then reconstruct the input data from this compressed form. They are particularly useful in unsupervised learning tasks, including anomaly detection, due to their ability to learn a compact and informative representation of the data (Hinton and Salakhutdinov, 2006).

Autoencoders consist of two main components: the encoder and the decoder. The encoder maps the input data x to a hidden representation h through a series of layers:

$$h = f(x)$$

where f is a non-linear transformation, often implemented using neural network layers.

The decoder then maps this hidden representation back to a reconstruction of the input data \hat{x} :

$$\hat{x} = g(h)$$

where g is another non-linear transformation.

The objective of training an autoencoder is to minimize the reconstruction error, which is typically

measured using a loss function such as the Mean Squared Error (MSE):

$$L(x, \hat{x}) = \|x - \hat{x}\|^2$$

The training process involves optimizing the parameters of both the encoder and the decoder to minimize this loss over the training dataset. Once trained, the autoencoder can be used for anomaly detection by comparing the reconstruction error of a new data point to a threshold. Data points with a high reconstruction error are considered anomalies, as the autoencoder is unable to accurately reconstruct them due to their deviation from the normal data distribution (Sakurada and Yairi, 2014).

Mathematically, the anomaly score for a data point x can be defined as:

$$AnomalyScore(x) = \|x - \hat{x}\|^2$$

Autoencoders can be enhanced with various techniques to improve their performance in anomaly detection. One common approach is to use sparse autoencoders, which impose a sparsity constraint on the hidden representations to encourage the network to learn more meaningful features (Ng, 2011). Another approach is to use variational autoencoders (VAEs), which introduce a probabilistic framework to the autoencoder, allowing it to model the uncertainty in the data and improving its robustness to anomalies (Kingma and Welling, 2014).

Variational Autoencoders introduce a latent variable z and define the encoder as a probabilistic model $q(z|x)$ and the decoder as $p(x|z)$. The objective is to maximize the Evidence Lower Bound (ELBO):

$$ELBO = \mathbb{E}_{q(z|x)}[\log p(x|z)] - KL(q(z|x) \parallel p(z))$$

where KL denotes the Kullback-Leibler divergence.

The Kullback-Leibler divergence, $KL(P \parallel Q)$, is a measure of how one probability distribution P diverges from a second, expected probability distribution Q . For two discrete probability distributions P and Q , it is defined as:

$$KL(P \parallel Q) = \sum_{x \in X} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

In the context of VAEs, P represents the true posterior distribution of the latent variables, and Q represents the approximate posterior distribution. Minimizing the KL divergence helps in making Q close to P , ensuring that the learned latent space distribution is similar to the true distribution of the data (Doersch, 2016).

Autoencoders have been successfully applied to various anomaly detection tasks, including image anomaly detection, network intrusion detection, and industrial equipment monitoring, due to their flexibility and ability to learn complex data distributions (Chen et al., 2017; Xia et al., 2015; Zhao et al., 2017).

Shifting from machine learning-based techniques to deep learning-based methods significantly improves outlier detection by enabling the management of complex and high-dimensional data. Although machine learning approaches like Isolation Forest offer strong solutions, deep learning methods such as autoencoders and GANs achieve higher accuracy by capturing intricate data patterns

through neural networks. This progression overcomes the limitations of traditional methods, enhancing anomaly detection in complex datasets.

3.5 Deep Learning-Based Approaches

Deep learning-based methods for outlier detection have become increasingly prominent due to their ability to handle complex, high-dimensional data and reveal intricate patterns that traditional methods may overlook. These approaches utilize deep neural networks to model sophisticated data distributions, providing a nuanced understanding of data structures. A key advantage of deep learning in this context is its capacity to automatically learn hierarchical feature representations, which is particularly effective for detecting subtle anomalies in large datasets.

Deep learning models excel by eliminating the need for manual feature engineering, as they can automatically extract relevant features from raw data. This is especially valuable for unstructured data types, such as images, audio, and text, where traditional methods often fall short. Additionally, deep learning techniques can capture non-linear relationships and dependencies, offering a robust framework for anomaly detection.

Autoencoders are a foundational deep learning approach for anomaly detection. They consist of an encoder that compresses input data into a latent representation and a decoder that reconstructs the original data from this compressed form. The reconstruction error, defined as the difference between the input data and its reconstruction, serves as an anomaly indicator. High reconstruction errors suggest that the data points are outliers, as the autoencoder struggles to accurately reconstruct data that deviates significantly from the training distribution (Hinton and Salakhutdinov, 2006).

Generative Adversarial Networks (GANs) also show significant promise for outlier detection. GANs involve two neural networks: a generator that creates realistic data samples and a discriminator that distinguishes between real and generated samples. This adversarial training enhances the model's ability to capture the data distribution, making it effective for identifying anomalies (Goodfellow et al., 2014). Variants such as Adversarially Learned Anomaly Detection (ALAD) further refine detection accuracy through bidirectional networks and cycle-consistency constraints (Zenati et al., 2018).

Variational Autoencoders (VAEs) introduce a probabilistic element to the autoencoder framework, modeling the latent space as a distribution rather than a fixed point. This probabilistic approach enhances robustness to anomalies and provides a principled method for measuring uncertainty in the data (Kingma and Welling, 2013). By maximizing the Evidence Lower Bound (ELBO) during training, VAEs ensure that the latent space effectively reflects the data distribution, improving anomaly detection performance.

Overall, deep learning-based approaches represent a significant advancement over traditional and machine learning methods in anomaly detection. Their ability to model complex, high-dimensional data distributions and automatically learn features makes them highly effective across various applications, including image and speech processing, network security, and industrial monitoring. The ongoing development of deep learning architectures and techniques promises further improvements in the accuracy and efficiency of outlier detection systems.

In this section, we will focus on two prominent deep learning-based approaches: Autoencoders and Generative Adversarial Networks (GANs), with a particular emphasis on Adversarially Learned Anomaly Detection (ALAD).

3.5.1 Deep Learning Autoencoders

Deep learning autoencoders and traditional machine learning autoencoders differ fundamentally in their architecture and capability to handle complex data representations.

Traditional machine learning autoencoders typically use shallow neural networks with one or two hidden layers to encode and decode data. These shallow networks limit the model's capacity to learn intricate patterns and relationships within the data. The encoding function f and decoding function g in these autoencoders are often simple linear or mildly non-linear transformations:

$$\mathbf{h} = f(\mathbf{x}) = W_1\mathbf{x} + b_1$$

$$\hat{\mathbf{x}} = g(\mathbf{h}) = W_2\mathbf{h} + b_2$$

where W_1 and W_2 are weight matrices, and b_1 and b_2 are bias terms. The limited depth and complexity of these networks restrict their ability to capture sophisticated data patterns, leading to potential limitations in detecting subtle anomalies in high-dimensional data.

In contrast, deep learning autoencoders utilize deep neural networks with multiple layers in both the encoder and decoder components. This deep architecture allows these models to learn hierarchical representations of the data, capturing more complex and abstract features. The encoding function f and decoding function g in deep learning autoencoders are defined as:

$$\mathbf{h} = f(\mathbf{x}) = f_L(f_{L-1}(\cdots f_1(\mathbf{x}) \cdots))$$

$$\hat{\mathbf{x}} = g(\mathbf{h}) = g_1(g_2(\cdots g_L(\mathbf{h}) \cdots))$$

where f_i and g_i represent the non-linear transformations applied by each layer, and L denotes the number of layers in the network. The depth of these networks enables them to model complex relationships and capture detailed features, making them more effective at identifying anomalies that are not apparent with shallow architectures.

The training of deep learning autoencoders involves minimizing the reconstruction error, similar to traditional autoencoders, but with the advantage of capturing richer data representations. The reconstruction error for deep learning autoencoders is also measured using Mean Squared Error (MSE):

$$L(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$$

However, the increased capacity of deep learning autoencoders allows for more nuanced anomaly detection due to their ability to learn from complex, high-dimensional datasets.

Furthermore, deep learning autoencoders can be enhanced with techniques such as Variational Autoencoders (VAEs), which incorporate probabilistic frameworks into the autoencoder architecture, and can better handle uncertainty and anomalies (Kingma and Welling, 2013).

Overall, deep learning autoencoders offer superior performance in anomaly detection compared to traditional machine learning autoencoders due to their ability to model complex data structures and learn hierarchical feature representations.

3.5.2 Generative Adversarial Networks (GANs) for Anomaly Detection

Generative Adversarial Networks (GANs) have emerged as a powerful tool for anomaly detection due to their capability to model complex data distributions through adversarial learning. GANs consist of two neural networks: a generator and a discriminator, which engage in a competitive training process. The generator's objective is to produce data samples that resemble the true data distribution, while the discriminator aims to distinguish between real samples and those generated by the generator.

Mathematically, the generator G and discriminator D are trained through a min-max game where the generator tries to minimize the loss function, while the discriminator tries to maximize it. The objective function for this adversarial process can be expressed as:

$$\min_G \max_D \mathbb{E}_{x \sim p_{data}} (\log D(x)) + \mathbb{E}_{z \sim p_z} (\log(1 - D(G(z))))$$

where x represents real data samples, z represents random noise inputs to the generator, $p_{data}(x)$ is the distribution of real data, $p_z(z)$ is the distribution of noise, and $D(\cdot)$ is the discriminator function that outputs the probability of a sample being real (Goodfellow et al., 2014).

The generator's loss function is:

$$L_G = -\mathbb{E}_{z \sim p_z} [\log D(G(z))],$$

which encourages the generator to produce data samples that the discriminator is likely to classify as real. Conversely, the discriminator's loss function is:

$$L_D = -\mathbb{E}_{x \sim p_{data}} (\log D(x)) - \mathbb{E}_{z \sim p_z} (\log(1 - D(G(z))))$$

which trains the discriminator to accurately differentiate between real and generated samples.

In the context of anomaly detection, GANs are effective because the generator learns to model the distribution of normal data. As a result, anomalies—data points that deviate significantly from the normal distribution—are less likely to be well-represented by the generator. Consequently, these anomalies are more easily identified by the discriminator, which can highlight discrepancies in the reconstruction process.

An advanced variant of GANs used for anomaly detection is the Adversarially Learned Anomaly Detection (ALAD) framework.

3.5.3 Adversarially Learned Anomaly Detection (ALAD)

Adversarially Learned Anomaly Detection (ALAD) represents an advanced methodology for detecting outliers by utilizing the principles of Generative Adversarial Networks (GANs). This technique enhances anomaly detection by employing the adversarial learning framework, which bolsters its ability to identify deviations from standard data distributions with increased accuracy and robustness.

ALAD leverages bidirectional neural networks, incorporating both forward and reverse mappings between the data space and a latent space. This approach is designed to capture and model the underlying data distribution more effectively. In the forward direction, the network maps input data x to a latent representation z , while in the reverse direction, it reconstructs the data from this

latent space representation. The model's effectiveness is measured by the combined adversarial and cycle-consistency loss functions.

The loss function for ALAD is formulated as:

$$ALAD_Loss = E_{x \sim p_{data}(x)} \|x - G(D(x))\|^2 + E_{z \sim p_z(z)} \|D(G(z)) - z\|^2 ,$$

where G denotes the generator, D represents the discriminator, $p_{data}(x)$ is the data distribution, and $p_z(z)$ is the distribution of latent variables. The first term ensures that the reconstructed data is close to the original input, thereby verifying the quality of reconstruction. The second term enforces cycle-consistency, ensuring that the latent representations of the generated samples are consistent with the original latent space distribution. This consistency helps to maintain the integrity of learned features and reduces the impact of noise and variations in data (Zenati et al., 2018).

By incorporating these bidirectional networks and cycle-consistency constraints, ALAD effectively improves the detection of anomalies. This method enhances traditional GAN-based anomaly detection by providing a more accurate representation of the data and by better distinguishing between normal and anomalous patterns. As a result, ALAD is particularly useful in applications involving high-dimensional and complex datasets, such as image analysis and fraud detection, where subtle anomalies may otherwise be overlooked.

3.6 ML-Consensus

Machine learning consensus, often referred to as ensemble learning, is a powerful approach in which multiple machine learning models are combined to improve the overall performance of predictive analytics tasks. The fundamental idea behind ML-consensus is that by aggregating the predictions of several models, the combined output can often surpass the performance of individual models. This approach is particularly beneficial in the context of outlier detection, where the variability and complexity of data can present significant challenges to single models (Dietterich, 2000; Polikar, 2006).

The ML-consensus method leverages the strengths of diverse algorithms to achieve a more robust and reliable detection system. Each model in the ensemble may capture different aspects of the data, and by integrating their outputs, the ensemble can mitigate the weaknesses of individual models. For example, some models might be more sensitive to certain types of anomalies, while others may excel in identifying different patterns. By combining these models, the ensemble can provide a more comprehensive detection mechanism.

Mathematically, the ML-consensus approach can be expressed as follows. Let f_1, f_2, \dots, f_M represent the individual models in the ensemble, where M is the total number of models. The final prediction \hat{y} for a given input x can be obtained by aggregating the predictions of all models:

$$\hat{y} = \text{Aggregate}(f_1(x), f_2(x), \dots, f_M(x))$$

The aggregation function can vary depending on the specific ensemble method used. Common aggregation strategies include majority voting for classification tasks, averaging for regression tasks, and weighted combinations where models are assigned different weights based on their performance (Dietterich, 2000).

A particularly useful aggregation method in the context of outlier detection is the weighted voting

scheme. Here, each model f_i is assigned a weight w_i that reflects its reliability or confidence in making predictions. The final prediction is then computed as:

$$\hat{y} = \sum_{i=1}^M w_i f_i(x)$$

where the weights w_i are typically normalized so that $\sum_{i=1}^M w_i = 1$.

Among the various ensemble learning techniques, bagging (Bootstrap Aggregating) involves training multiple instances of the same model on different subsets of the training data, created by sampling with replacement from the original dataset. The predictions of the individual models are then aggregated, typically by majority voting for classification or averaging for regression, helping to reduce variance and prevent overfitting (Breiman, 1996). Boosting sequentially trains models, with each model attempting to correct the errors of its predecessor. The final prediction is a weighted sum of the predictions from all models, with algorithms like AdaBoost and Gradient Boosting showing significant improvements in accuracy by focusing on difficult-to-predict cases (Freund and Schapire, 1997; Friedman, 2001).

However, our primary focus will be on stacking, voting, random forests, and meta-learning. Stacking involves training multiple base models and then using their predictions as inputs to a meta-model, which makes the final prediction. This approach is particularly effective when base models are diverse, as the meta-model learns how to best combine their predictions, often leading to enhanced performance (Wolpert, 1992). In the voting ensemble method, each model in the ensemble casts a "vote" for each prediction. For classification tasks, the final prediction is the class with the majority of votes. For regression tasks, the final prediction is the average of the individual model predictions. Voting can be simple, with equal weights, or weighted, where models with higher accuracy contribute more to the final prediction, making it a versatile and straightforward aggregation method (Dietterich, 2000).

Random forests extend bagging by using decision trees as the base models and introducing randomness in feature selection for each split in the decision trees. This results in a diverse set of trees that collectively improve prediction accuracy and reduce overfitting (Breiman, 2001). Focusing on stacking, voting, and random forests is justified by their proven effectiveness and flexibility in various applications, including outlier detection. These methods allow for combining the strengths of different models, enhancing robustness, and improving predictive performance. Meta-learning, in particular, provides a framework for integrating the outputs of multiple models, enabling the development of a highly adaptive and resilient outlier detection system.

For bagging, the training sets $\{D_1, D_2, \dots, D_M\}$ are created by sampling with replacement from the original dataset D . The final prediction for an input x is given by:

$$\hat{y} = \frac{1}{M} \sum_{i=1}^M f_i(x)$$

For boosting, each model f_i is trained on a weighted version of the dataset, where the weights are adjusted based on the errors of the previous model. The final prediction is a weighted sum of the model predictions:

$$\hat{y} = \sum_{i=1}^M \alpha_i f_i(x)$$

$$i=1$$

where α_i are the weights determined during training. For stacking, the base models f_1, f_2, \dots, f_M are first trained, and their predictions are used to train a meta-model g . The final prediction is given by:

$$\hat{y} = g(f_1(x), f_2(x), \dots, f_M(x))$$

For voting, if C_k represents the prediction of the k -th class, the final prediction for classification is:

$$\hat{y} = \arg \max_k \sum_{i=1}^M w_i 1_{f_i(x)=C_k}$$

For random forests, each decision tree T_i is trained on a different bootstrap sample, and the prediction is:

$$\hat{y} = \frac{1}{M} \sum_{i=1}^M T_i(x)$$

Building on this understanding of ML-consensus and the mathematical foundations underpinning ensemble methods, we now shift our focus to the specific techniques that form the theoretical backbone of ensemble learning: stacking, voting, random forests, and meta-learning. These methods have been selected due to their established roles in enhancing predictive accuracy and reliability by combining multiple models. In the subsequent discussion, we will delve into the theoretical underpinnings of each approach, examining how they synergize different models to improve performance. By gaining a comprehensive understanding of these ensemble techniques, we lay the groundwork for advanced applications in outlier detection and beyond.

3.6.1 Voting

Voting is a fundamental ensemble learning technique designed to enhance the accuracy and robustness of predictive models by aggregating the predictions of multiple base models. The core idea behind voting is to leverage the diversity among various models to improve the overall performance of the ensemble. This technique harnesses the collective wisdom of the base models, which often leads to better generalization than any single model alone.

This approach is rooted in the principle that different models may capture different aspects of the data or have varying strengths and weaknesses. By combining their predictions, we can exploit their complementary capabilities and reduce the likelihood of errors made by individual models. This approach is particularly valuable in scenarios where no single model performs optimally across all aspects of the data.

In practice, voting involves two main strategies: **majority voting** and **weighted voting**. These strategies determine how the predictions of the base models are combined to arrive at a final decision.

In majority voting, each base model in the ensemble casts a "vote" for a particular class. The final prediction is based on the class that receives the most votes. Mathematically, if an ensemble consists of n classifiers and each classifier M_i provides a prediction \hat{y}_i for a given input x , the final prediction \hat{y} is determined by:

$$\hat{y} = \text{mode}(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$$

where *mode* represents the most frequently predicted class among the base models. This method capitalizes on the idea that even if individual models have varying levels of accuracy, the majority vote tends to represent the consensus view and thus can improve reliability (Breiman, 1996).

Weighted voting extends this concept by assigning different weights to each base model based on its performance or reliability. Instead of treating each model's vote equally, weighted voting accounts for the varying quality of predictions from different models. If w_i denotes the weight assigned to the i -th model and p_i represents its predicted probability for a class, the final prediction is computed as:

$$\hat{y} = \underset{c \in \mathbf{I}^n}{\operatorname{argmax}} \sum_{i=1}^L w_i \cdot p_i$$

In this formulation, the final prediction is determined by the class that has the highest weighted sum of probabilities. Weighted voting is particularly useful when some models in the ensemble are known to be more accurate or reliable than others. By adjusting the weights, we can enhance the ensemble's performance and address model-specific biases or errors (Kuncheva, 2004).

The strength of voting lies in its ability to aggregate diverse predictions and reduce overall variance. Each base model in the ensemble may have unique errors or biases, but by combining their outputs, voting can average out these inconsistencies and produce a more stable and accurate final prediction. This principle aligns with the concept of "*wisdom of crowds*", where collective decisions tend to be more accurate than individual judgments (Hansen and Salamon, 1990).

In classification tasks, voting helps in achieving a consensus decision among the base models, while in regression tasks, a similar approach can be used to average the predictions. Despite its simplicity, voting can effectively address various challenges in predictive modeling, such as model overfitting and underfitting, by incorporating the strengths of multiple models and mitigating their weaknesses.

However, voting also has its limitations. For instance, if the base models are highly correlated or if they consistently make similar errors, the benefits of voting may be diminished. To overcome these challenges, it is essential to ensure that the base models are diverse and independently trained. Proper selection and weighting of models can further enhance the effectiveness of the voting ensemble.

Having explored the fundamentals of voting, an ensemble technique that leverages the collective predictions of multiple base models to enhance accuracy, we now turn our attention to another advanced method: stacking. While voting aggregates the outputs of various models to achieve a consensus, stacking builds upon this idea by introducing a meta-model that learns to optimally combine the predictions of base models. This approach not only harnesses the strengths of diverse algorithms but also refines the final predictions through a learned integration process. The transition from voting to stacking represents a progression from simple aggregation to a more sophisticated, adaptive method that enhances predictive performance through a meta-learning framework.

3.6.2 Stacking

Stacking, or stacked generalization, is a sophisticated ensemble learning technique designed to

enhance predictive performance by leveraging multiple base models through a meta-model. This

approach harnesses the distinct advantages of various algorithms to improve overall accuracy and mitigate the weaknesses inherent in individual models.

At the core of stacking is the concept of training several base models, referred to as level-0 models, on the same dataset. These base models can be diverse in nature, such as decision trees, support vector machines, or neural networks, each capturing different patterns and relationships within the data (Wolpert, 1992). Once trained, these base models generate predictions for the dataset. These predictions serve as input features for a meta-model, or level-1 model. The meta-model, often a simpler algorithm like linear regression or logistic regression, learns to optimally combine the outputs of the base models to produce a final prediction (Breiman, 1996).

Mathematically, let D denote the training dataset and M_1, M_2, \dots, M_n represent the base models. Each base model M_i generates a prediction \hat{y}_i for a given input x . The meta-model H uses these predictions $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$ as features to produce the final prediction \hat{y} , expressed as:

$$\hat{y}_i = M_i(x), \quad \text{for } i = 1, 2, \dots, n$$

$$\hat{y} = H(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$$

In this formulation, H is the meta-model that synthesizes the outputs from the base models to generate the final prediction.

The strength of stacking lies in its ability to capitalize on the diversity and complementary strengths of the base models. By aggregating models that each capture different facets of the data, stacking reduces overall bias and variance, leading to improved generalization on new, unseen data (Sill et al., 2009). The meta-model plays a pivotal role by learning the optimal combination of base model predictions, addressing the specific errors of individual models.

Stacking has proven effective across a range of domains, including both classification and regression tasks. It is particularly valuable in scenarios where no single model performs optimally across all cases, as the ensemble approach can encompass a broader spectrum of data characteristics and often surpasses any individual model's performance (Krogh and Vedelsby, 1995).

Nonetheless, stacking is not without challenges. One potential issue is overfitting, particularly if both the base models and the meta-model are excessively complex. To mitigate this risk, techniques such as cross-validation are utilized. In cross-validation stacking, the training dataset is divided into several folds, with each base model being trained on different subsets of the data. The meta-model is then trained using out-of-fold predictions, which provides a more robust and unbiased estimate of the base models' performance (Breiman, 1996).

Building on the concept of voting, which aggregates predictions from multiple models, we now explore Random Forest, an advanced ensemble technique that enhances voting by using decision trees and introducing randomness in model training. Random Forest not only combines predictions but also improves model diversity and accuracy through methods like bagging and feature randomness.

3.6.3 Random Forest

Random Forest is a widely used ensemble learning method that combines the predictions of multiple decision trees to improve overall model accuracy and robustness. This approach builds on the concept of bagging, or bootstrap aggregating, which involves training multiple models on

different subsets of the training data and then averaging their predictions to reduce variance and improve generalization.

In a Random Forest, multiple decision trees are constructed using different random subsets of the training data and features. Each tree is trained on a bootstrap sample, a random sample of the original dataset with replacement. During the training of each decision tree, only a random subset of features is considered for splitting at each node. This introduces diversity among the trees and helps prevent overfitting, which can occur if a single decision tree is too complex and fits the noise in the training data.

Mathematically, let D represent the original training dataset, and let T_1, T_2, \dots, T_B denote the B decision trees in the Random Forest. For a given input x , each tree T_i produces a prediction \hat{y}_i . The final prediction \hat{y} is obtained by aggregating the predictions of all the trees. For classification tasks, the aggregation is typically done using majority voting, while for regression tasks, it is done by averaging the predictions:

$$\hat{y} = \frac{1}{B} \sum_{i=1}^B \hat{y}_i$$

where \hat{y}_i is the prediction from the i -th tree, and B is the total number of trees in the forest (Breiman, 2001).

The strength of Random Forest lies in its ability to reduce both variance and bias compared to individual decision trees. By averaging the predictions of multiple trees, Random Forest mitigates the overfitting problem that can arise with complex decision trees. Additionally, the random feature selection process helps to capture different aspects of the data and improve the robustness of the model.

Random Forest has been shown to perform well in various domains, including finance, healthcare, and image analysis. Its ability to handle large datasets with many features and provide robust predictions makes it a valuable tool in machine learning. However, despite its advantages, Random Forest can be computationally intensive, particularly with a large number of trees and features, and may lack interpretability compared to simpler models.

The introduction of meta-learning builds upon the previous discussion of ensemble methods like voting, stacking, and Random Forest. By incorporating meta-learning, we further explore advanced strategies for enhancing machine learning models' adaptability and efficiency, paving the way for more sophisticated and versatile approaches in outlier detection and other domains.

Having examined Random Forest's approach to boosting predictive accuracy through decision trees and ensemble methods, we now turn to meta-learning. While Random Forest focuses on combining models, meta-learning enhances the learning process itself by optimizing how algorithms adapt and generalize across different tasks.

3.6.4 Meta-Learning

Meta-learning, also known as "*learning to learn*", is a sophisticated approach in machine learning that focuses on optimizing learning algorithms themselves. Rather than solely training a model on a specific task, meta-learning aims to improve the learning process by leveraging experience gained from multiple tasks or datasets. The goal is to develop models that can quickly adapt to new tasks with minimal additional training, effectively transferring knowledge across different learning scenarios.

The core idea behind meta-learning is to design algorithms that can learn from their previous learning experiences. This involves two levels of learning: the first level focuses on training models on various tasks, while the second level involves optimizing the learning process based on the performance of these models. Meta-learning algorithms typically operate by learning a set of parameters or strategies that enhance the generalization capability of the models they train.

Mathematically, let \mathcal{T} represent a distribution over tasks, and let D_t denote the dataset for task t . Meta-learning aims to find a learning algorithm or set of parameters θ that optimizes performance across all tasks in \mathcal{T} . Given a model f with parameters θ , the meta-learning objective is to minimize the loss function L over the distribution of tasks:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{T}} [L(f_{\theta}, D_t)]$$

Here, θ^* represents the optimal parameters that generalize well across different tasks, and L is the loss function measuring the model's performance on each task t .

One popular meta-learning approach is Model-Agnostic Meta-Learning (MAML), which seeks to find initial model parameters that can be quickly adapted to new tasks with a few gradient updates. MAML optimizes the model's parameters such that, after a few steps of gradient descent on a new task, the model's performance is significantly improved. The optimization problem for MAML can be formulated as:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{T}} [L(f_{\theta - \alpha \nabla_{\theta} L(f_{\theta}, D_t)}, D_t)]$$

where α is the learning rate for the task-specific updates (Finn et al., 2017).

Another approach is the use of meta-architectures, where the meta-learning process involves optimizing the structure or architecture of the learning model itself. This approach can involve techniques such as neural architecture search, where the meta-learner explores different network architectures to find the most effective configuration for a given task (Zoph and Le, 2017).

Meta-learning has shown promise in a range of applications, including few-shot learning, where models must generalize from a limited number of examples, and hyperparameter optimization, where the goal is to find the best hyperparameters for machine learning models. By focusing on improving the learning process itself, meta-learning offers a powerful framework for enhancing model performance and adaptability across diverse tasks.

3.7 Data Synthesis

Data synthesis is a critical aspect of developing and evaluating machine learning models, particularly for tasks like outlier detection. By generating synthetic data, researchers can create controlled datasets that are essential for training, validating, and testing models. This approach is invaluable when working with high-dimensional data or when real data is limited or difficult to access.

One of the primary advantages of data synthesis is its ability to overcome challenges related to data scarcity. In many real-world scenarios, outliers are rare and often underrepresented in available datasets. Synthetic data allows for the creation of large and diverse datasets with known properties, ensuring that models are trained on a broad spectrum of anomaly types. This exposure is crucial for developing robust models capable of identifying subtle and varied anomalies.

Synthetic data also enables the exploration of different data distribution scenarios that may be challenging to capture with real data. For instance, certain anomalies might be too rare or complex to be adequately represented in natural datasets. By generating synthetic anomalies with specific characteristics, researchers can test and refine models under controlled conditions, gaining insights into their performance and robustness.

Several methods exist for generating synthetic data, each suited to different types of distributions and applications. The Box-Muller transform, a classic technique, generates synthetic data from a normal distribution. Given two independent random variables U_1 and U_2 uniformly distributed in the interval $[0, 1)$, the Box-Muller transform produces two independent standard normal variables Z_0 and Z_1 using:

$$\begin{aligned} Z_0 &= \sqrt{-2 \ln U_1} \cos(2\pi U_2) \\ Z_1 &= \sqrt{-2 \ln U_1} \sin(2\pi U_2) \end{aligned}$$

While effective for generating normal distributions, this method is limited to scenarios where the data distribution is Gaussian (Box and Muller, 1958).

To address a broader range of data distributions, other statistical models can be employed. For instance, Multivariate Gaussian distributions allow for modeling data with specified correlations between variables.

If Σ represents the covariance matrix of the distribution, synthetic data can be generated using:

$$X = \mu + LZ$$

where μ is the mean vector, L is the Cholesky decomposition of Σ , and Z is a vector of standard normal random variables (Trefethen and Bau, 1997). This approach enables the creation of datasets with specific statistical properties.

Cholesky decomposition is a matrix factorization technique used to decompose a symmetric positive-definite matrix into a product of a lower triangular matrix and its transpose. For a symmetric positive-definite matrix A , Cholesky decomposition finds a lower triangular matrix L such that:

$$A = LL^T$$

where L^T denotes the transpose of L . The matrix L has zeros above the diagonal and non-zero elements on and below the diagonal. To compute the elements of L , the following recursive formulas are used:

For each diagonal element L_{ii} :

$$L_{ii} = \sqrt{A_{ii} - \sum_{k=1}^{i-1} L_{ik}^2}$$

For each off-diagonal element L_{ij} (where $i > j$):

$$L_{ij} = \frac{1}{L_{jj}} \left(A_{ij} - \sum_{k=1}^{j-1} L_{ik} L_{jk} \right)$$

These formulas ensure that L is a lower triangular matrix with positive diagonal entries (Golub and Van Loan, 2013). Cholesky decomposition is particularly useful for generating multivariate normal data with specified covariance structures, enhancing its applicability in various fields.

Generative Adversarial Networks (GANs) have also revolutionized synthetic data generation. GANs consist of two neural networks—the generator and the discriminator—that are trained adversarially. The generator creates synthetic data samples, while the discriminator assesses their realism. This adversarial training encourages the generator to produce data that closely resembles the real distribution, making GANs effective for generating complex, high-dimensional data (Goodfellow et al., 2014).

Variational Autoencoders (VAEs) are another powerful tool for data synthesis. VAEs use a probabilistic approach to model data's latent space, allowing for the generation of new samples by sampling from the learned latent distribution. The VAE loss function combines reconstruction loss with the Kullback-Leibler (KL) divergence between the learned latent distribution and a prior distribution:

$$ELBO = \mathbb{E}_{q(z|x)}[\log p(x|z)] - KL(q(z|x) \parallel p(z))$$

where $q(z|x)$ is the approximate posterior, $p(x|z)$ is the likelihood, and $p(z)$ is the prior (Kingma and Welling, 2013). This approach facilitates the generation of data with complex structures and variations.

Additionally, methods such as SMOTE (Synthetic Minority Over-sampling Technique) address class imbalance by creating synthetic samples through interpolation between existing data points. SMOTE is particularly useful for enhancing the representation of minority classes, thereby improving model performance on imbalanced datasets (Chawla et al., 2002).

In summary, data synthesis is a vital component of machine learning, offering techniques to generate controlled datasets with specific properties. By overcoming data scarcity and distribution variability challenges, synthetic data generation enhances model development and evaluation, providing a foundation for robust and effective machine learning solutions.

3.8 Conclusion

In conclusion, this chapter has provided a thorough exploration of outlier detection methodologies, encompassing both classical statistical techniques and contemporary deep learning approaches. We examined methods such as Mahalanobis Distance, Isolation Forest, Autoencoders, and Adversarially Learned Anomaly Detection (ALAD), revealing their strengths and limitations. These insights are essential for understanding how each method can be effectively applied to different data types and anomaly detection scenarios (Iglewicz & Hoaglin, 2007; Hsieh & Lee, 2020; Mallick & Kim, 2019; Kingma & Welling, 2014; Zenati et al., 2018).

Additionally, the chapter delved into the concept of ML-consensus, emphasizing the integration of multiple machine learning techniques to enhance predictive performance. By discussing approaches such as stacking (Wolpert, 1992; Breiman, 1996), voting (Freund & Schapire, 1997), random forests (Breiman, 2001), and meta-learning (Balcázar et al., 2012), we highlighted how combining diverse models can address the limitations of individual methods and improve overall anomaly detection accuracy.

This critical analysis lays a robust foundation for the subsequent development and implementation

of outlier detection systems. The insights gained underscore the necessity of a nuanced approach that integrates various methodologies, including the ML-consensus framework, to handle the complexities of real-world data effectively. This understanding will be pivotal in designing a comprehensive detection framework capable of managing diverse and high-dimensional datasets.

The theoretical insights from this chapter will guide the solution engineering and implementation detailed in the following chapters. Chapter 4 will leverage these foundations to design and optimize an advanced outlier detection system, while Chapter 5 will focus on the practical aspects of deploying and integrating this system.

By addressing the identified challenges and utilizing the theoretical knowledge and ML-consensus approach discussed here, the study aims to enhance the effectiveness and reliability of outlier detection systems. This foundational analysis will support the development of solutions that are both theoretically robust and practically viable, ensuring their applicability and impact in real-world scenarios.

Chapter 4

Solution Engineering

In this chapter, we transition from theoretical exploration to the practical implementation of the outlier detection system discussed in earlier chapters. Building on the foundational knowledge gained, this chapter aims to detail the engineering aspects of the solution, ensuring that theoretical insights are effectively translated into a functional and robust system.

The primary objective of this chapter is to outline the comprehensive solution architecture and implementation strategy for the outlier detection system. This involves designing and documenting the system's architecture, selecting appropriate technologies, and analyzing the solution's performance. By systematically addressing these aspects, we aim to create a solution that not only adheres to the theoretical principles discussed but also performs effectively in real-world scenarios. We will begin by conceptualizing the solution architecture, which will involve creating use case diagrams, sequence diagrams, activity diagrams, and workflow engineering. These tools will help visualize the system's functionality and interactions, providing a clear blueprint for the subsequent development stages.

Following the architectural design, we will delve into the choice of technologies, justifying the selection based on the project's requirements. This section will provide a comprehensive overview of the technologies employed, detailing their roles and contributions to the solution.

Data synthesis will then be addressed, focusing on the structure and preparation of data essential for the system's operation. Understanding how data is organized and processed is crucial for ensuring the accuracy and efficiency of the outlier detection algorithms.

The solution analysis section will explore the implementation of the desktop application, including the choice of algorithms, complexity analysis, and a comparative study of algorithm performance. This analysis will provide insights into the effectiveness of the solution and highlight any potential areas for improvement.

Finally, the chapter will conclude with a summary of the key findings and reflections on the integration of theoretical insights into practical implementation. This conclusion will also discuss future directions for enhancing the solution and addressing any remaining challenges.

By systematically addressing these components, this chapter aims to bridge the gap between theory and practice, ensuring that the outlier detection system is both theoretically sound and practically viable.

4.1 Solution Architecture

The solution architecture is a critical component of system design, serving as the blueprint that outlines the structural framework and interactions within the outlier detection system. This architecture is pivotal in translating theoretical concepts into a tangible and functional system, ensuring that all components work together seamlessly to achieve the desired outcomes. By defining the system's structure and interactions, the architecture helps in creating a well-organized and efficient solution that adheres to the principles outlined in the theoretical background.

In system design, particularly within the context of Unified Modeling Language (UML), the solution architecture employs various diagrams to represent the system's components and their relationships. UML provides a standardized way to visualize system structure and interactions, which is essential for both understanding and communicating the design. Key UML diagrams, including use case diagrams, sequence diagrams, activity diagrams, and workflow engineering, play a crucial role in depicting different aspects of the system's architecture.

The importance of the solution architecture lies in its ability to provide a clear and structured framework for system development. It ensures that all components are aligned with theoretical principles and operational requirements, facilitating a smooth transition from design to implementation. By employing UML diagrams and engineering principles, the architecture provides a comprehensive overview of the system's structure and interactions, guiding the development process and ensuring that the final solution meets the intended goals.

4.1.1 Use Case Diagram

The Use Case Diagram is an integral part of the solution architecture, offering a high-level view of the interactions between users (actors) and the system. It is essential for capturing the functional requirements of the system and ensuring that all necessary features are incorporated into the design. The diagram serves as a visual representation of the system's functionality, illustrating how different actors interact with the system to achieve their goals.

In the context of this outlier detection system, the Use Case Diagram focuses on several key functionalities that are crucial for the system's operation. These functionalities are designed to address specific needs and tasks related to outlier detection and data analysis.

In our system, the Use Case Diagram includes several key functionalities that are pivotal for effective operation. These functionalities encompass parameterizing the model, choosing variables for outlier detection, data visualization, generating and exporting reports, each serving distinct purposes within the overall framework.

Actor: Admin (In this context, it refers mainly to the person in charge of maintaining the software and updating the model). **Actor:** User (member of the Data Analytics or Statistics Department).

Parameterizing the Model:

Users can configure the outlier detection model through an intuitive form. This process involves inputting various parameters necessary for model training and outlier detection. The parameters may include settings specific to the algorithms used, thresholds for anomaly detection, and other relevant configurations. This functionality ensures that the model can be tailored to meet the specific requirements of different datasets.

Choosing Variables for Outlier Detection:

Once the model is parameterized, users select the variables on which the outlier detection will be performed. This selection process is critical as it determines the focus of the analysis and ensures that the outlier detection process targets the most relevant features of the dataset.

Data Visualization:

The application provides visualization tools to display the data associated with the selected variables. This functionality allows users to gain insights into the data distribution and understand how the outlier detection model interacts with the dataset. Visualizations can include charts, graphs, and other visual aids that enhance data interpretation.

Generating and Exporting Reports:

After the outlier detection process is complete, users can generate detailed reports that highlight the identified outliers within the existing data. These reports are essential for communicating findings and supporting decision-making. The capability to export these reports in various formats ensures that users can share results with stakeholders effectively.

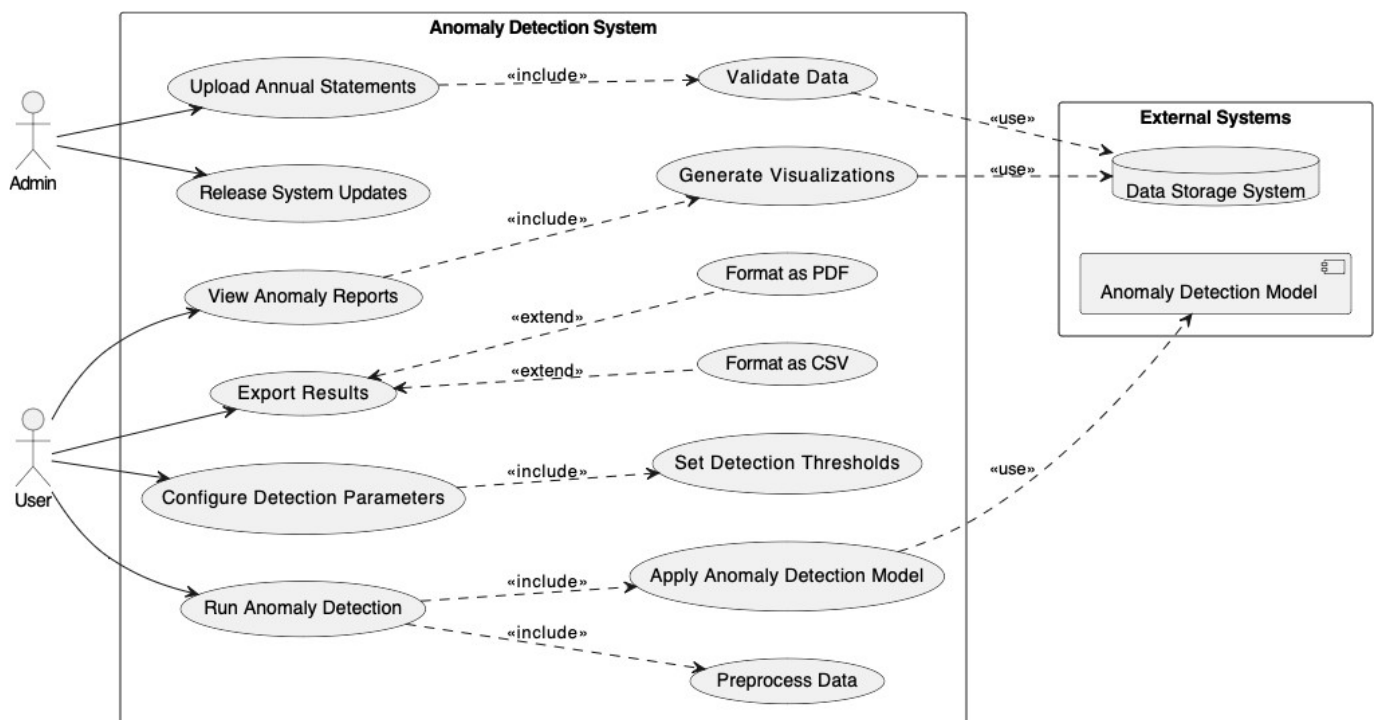


Figure 4.1 – Use Case Diagram

The use case diagram illustrates how these functionalities are interconnected and how users interact with the system to achieve their objectives. It maps out the process flow from parameter configuration to report generation, ensuring that all components work together seamlessly to provide a comprehensive solution for outlier detection.

This diagram not only helps in understanding the user interactions but also serves as a foundation for designing the system's architecture and user interface. By clearly defining the system's functionalities and their interactions, the use case diagram facilitates a structured approach to developing and implementing the application, ensuring that all user needs are met efficiently and effectively (Fowler, 2004; Rumbaugh et al., 1999).

4.1.2 Class Diagram

The class diagram is a vital component of the Unified Modeling Language (UML) used to describe the static structure of a system by showcasing its classes, attributes, methods, and the relationships among objects. In this context, the class diagram will represent the core components of the outlier detection system, illustrating how different entities interact within the application to achieve the desired functionality.

In this context, the class diagram delineates the primary classes involved in the application, including Model, Data, Visualization, and Report classes. Each class encapsulates specific functionalities and attributes that contribute to the overall operation of the system.

The class diagram provides a comprehensive overview of the system's structure, highlighting the key components and their interactions. This visual representation is essential for developers to understand the system's architecture, identify potential areas for optimization, and ensure that all necessary functionalities are incorporated. It also serves as a blueprint for implementing the system, guiding the development process from initial design to final deployment (Booch et al., 1998; Rumbaugh et al., 2004).

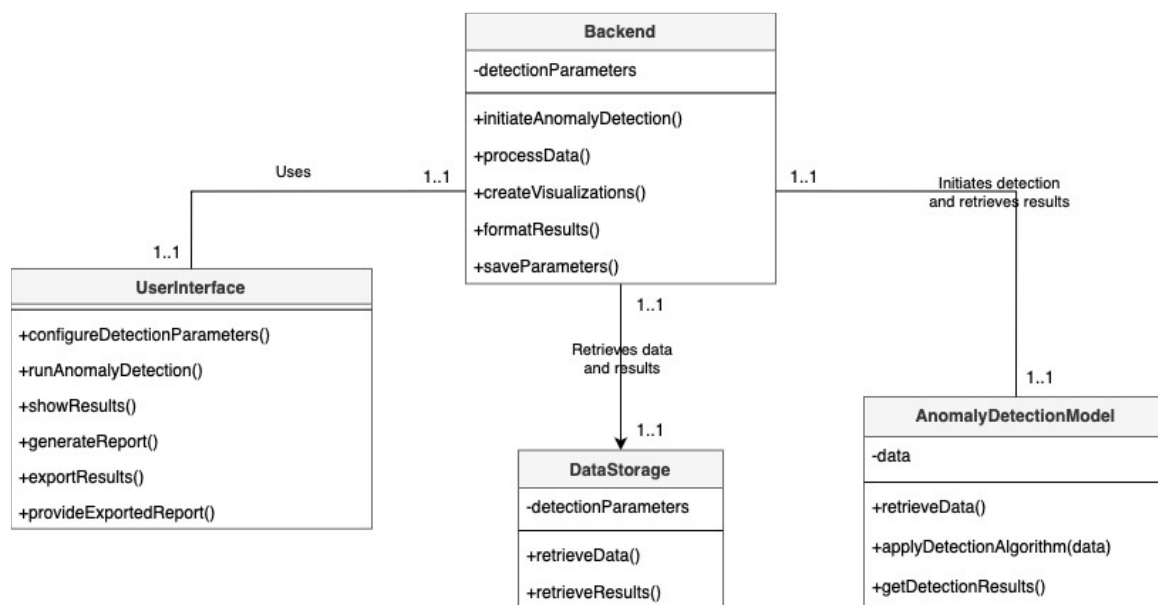


Figure 4.2 – Class Diagram

4.1.3 Sequence Diagram

The sequence diagram is an essential tool for detailing the interactions between different components of the system over time. It provides a step-by-step visualization of how processes are executed, illustrating the flow of information and the sequence of operations. In our context, the sequence diagram will focus on the interactions required to perform outlier detection, from parameter configuration to report generation.

In this context, the sequence diagram captures the dynamic behavior of the system as users engage with various functionalities. The primary actors involved are the user and the system components, including the user interface, the data processing module, and the reporting module.

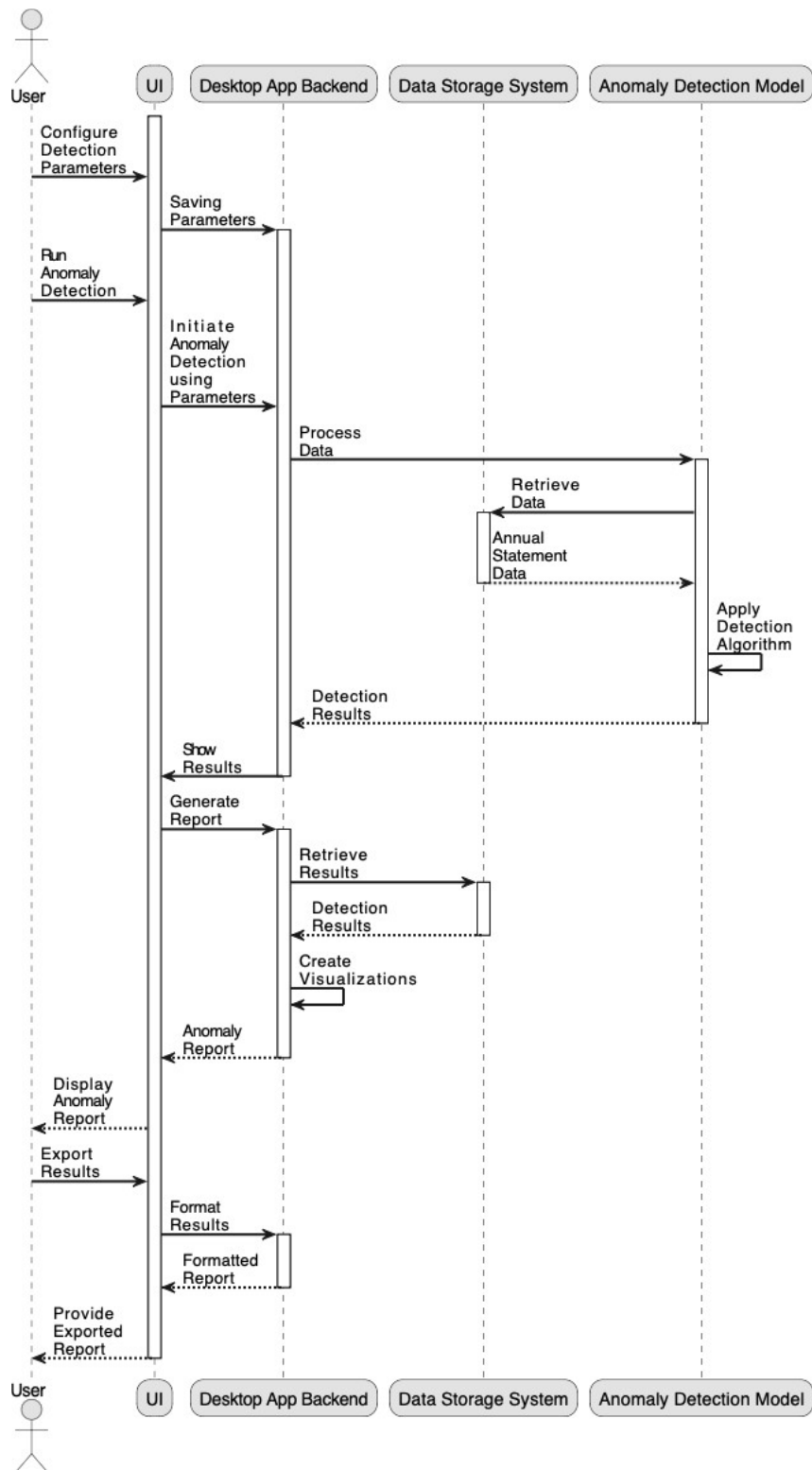


Figure 4.3 – Sequence Diagram

This comprehensive visualization aids in understanding the intricacies of system interactions, facilitating the identification of potential bottlenecks and ensuring the robustness of the implementation. Furthermore, it serves as a valuable reference for developers and stakeholders, providing a clear roadmap for system development and integration (Jacobson et al., 1999; Rumbaugh et al., 2005).

4.1.4 Activity Diagram

The activity diagram is a critical tool for mapping out the workflow and activities within the system, providing a detailed visualization of the process flow from start to finish. It showcases the sequence of activities, decision points, and the flow of control between different stages of the system’s operation. In our context, the activity diagram will illustrate the entire workflow involved in outlier detection, from initial user interaction to the final report generation.

In this context, the activity diagram begins with the user initiating the process by accessing the application. The first activity involves the user navigating to the model parameterization form. Here, the user inputs various parameters required to configure the outlier detection model. This step involves multiple activities, including form completion, parameter validation, and submission. The system checks the validity of the input parameters and proceeds to configure the model accordingly.

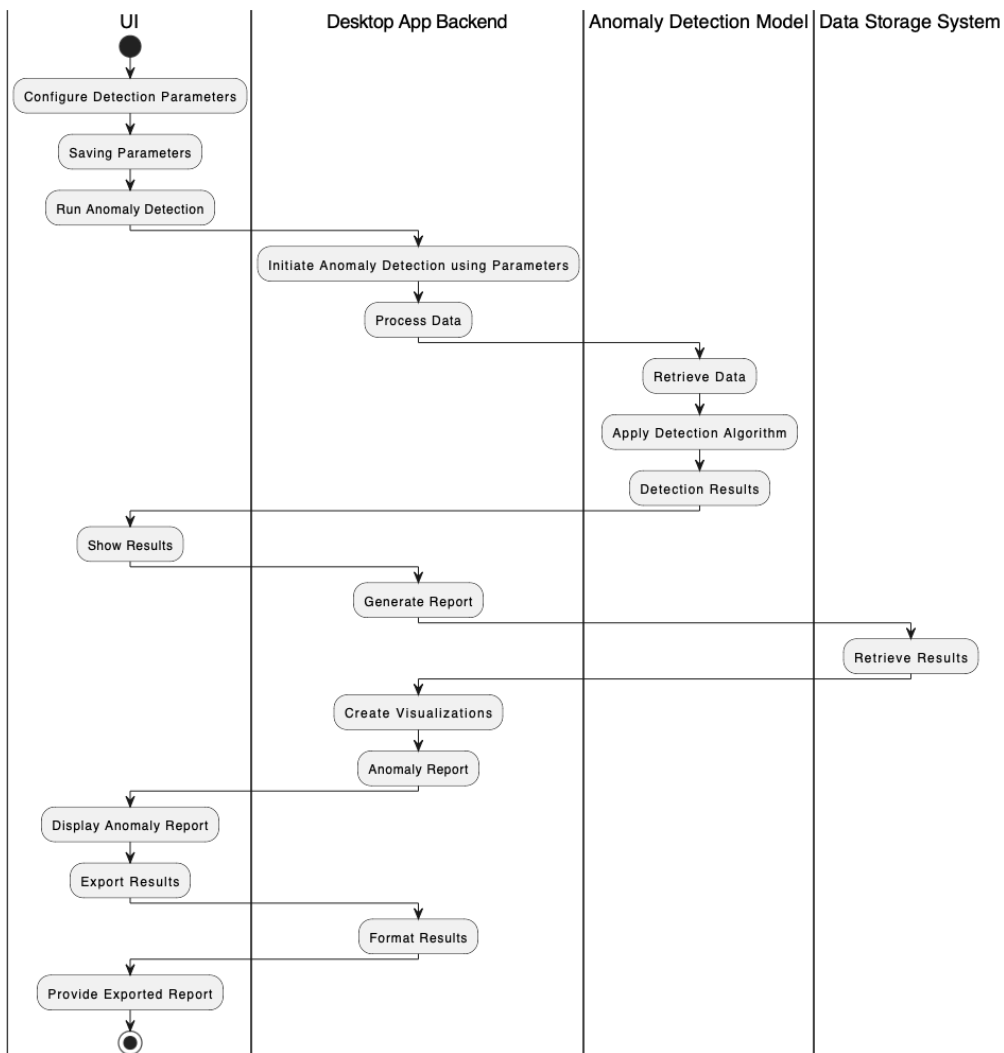


Figure 4.4 – Activity Diagram

The activity diagram emphasizes the sequential and parallel activities within the workflow, illustrating the flow of control from one activity to the next. It highlights decision points where user input or system validation is required, ensuring that all activities are completed accurately and in the correct order.

Mathematically, we can represent some of the processes involved. For example, let X be the

input data matrix, and Y be the output of the detection model. The parameterization step can be represented as setting parameters θ for the model M :

$$M(X; \theta) \rightarrow Y$$

The activity diagram not only provides a clear depiction of the workflow but also highlights the interactions between different components and the user. This detailed visualization aids in understanding the sequence of operations, identifying potential areas for optimization, and ensuring the robustness of the system's implementation. It serves as a crucial reference for developers and stakeholders, guiding the development and integration process (OMG, 2011; Booch et al., 1998).

4.1.5 Workflow engineering

Workflow engineering is a crucial aspect of system design that focuses on defining, optimizing, and automating the flow of tasks and processes within an application. In the context of our outlier detection system, workflow engineering ensures that all interactions between different system components are streamlined to provide a seamless user experience. This section outlines the steps involved in designing an effective workflow for the outlier detection application, from parameter input to report generation and visualization.

In our context, the workflow begins with the user interacting with the system to input parameters for the outlier detection model. This interaction is captured in the initial step where the user fills out a form specifying the model parameters and selects the variables for outlier detection. Once the parameters and variables are submitted, the system initiates the data processing stage.

The data processing stage involves loading the preexisting processed data into the system. The data is then prepared for analysis by the model. This preparation includes any necessary transformations or normalizations to ensure that the data is in a suitable format for the outlier detection algorithms. The prepared data is then fed into the model, which runs the specified algorithms to identify potential outliers.

Mathematically, this can be represented as follows. Let X be the dataset, θ be the model parameters, and f be the outlier detection function. The process can be formalized as:

$$X_{processed} = preprocess(X) \rightarrow \dots \rightarrow Y = f(X_{processed}, \theta)$$

where $X_{processed}$ represents the transformed data, and Y represents the outlier detection results.

Throughout this workflow, several interactions and dependencies need to be managed. The system must ensure that data is accurately passed between components and that each stage of the workflow is executed correctly. This requires robust error handling and validation mechanisms to address any issues that may arise during data processing, model execution, or report generation.

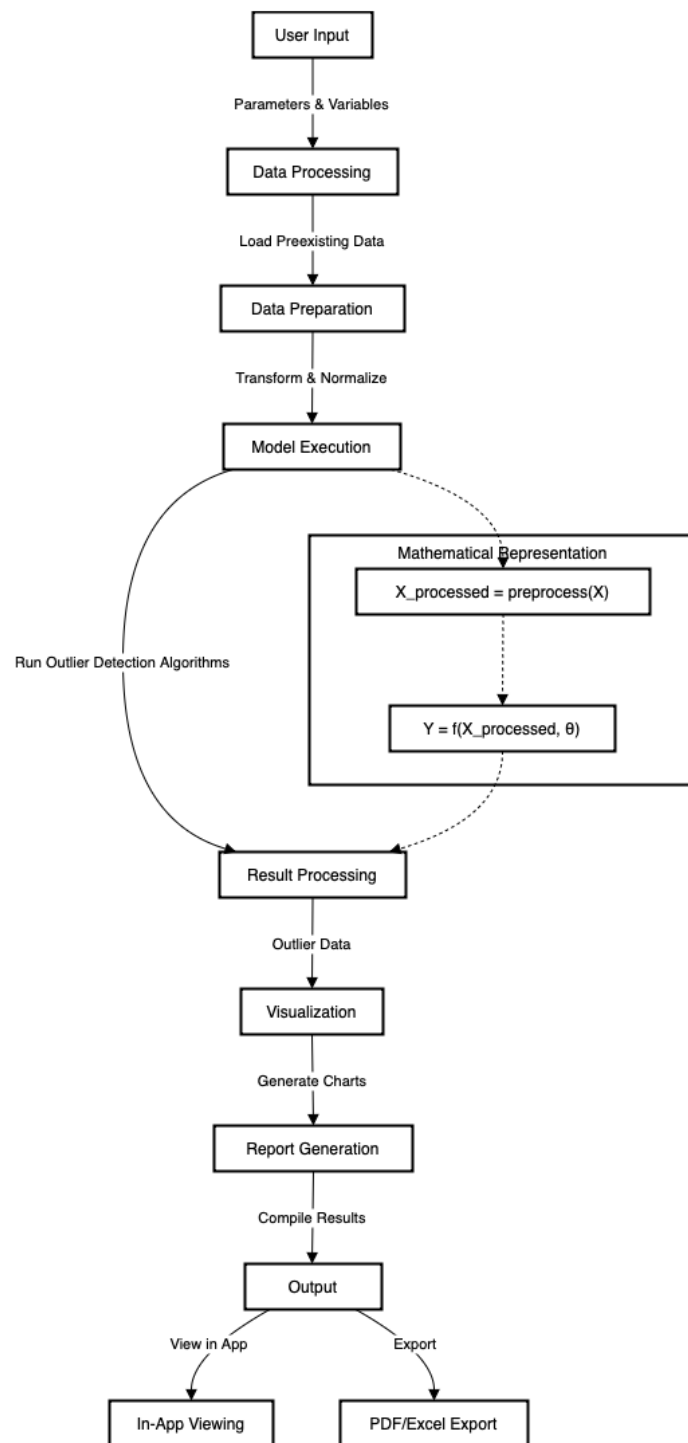


Figure 4.5 – Workflow

Workflow engineering for the outlier detection system aims to optimize these interactions to ensure efficiency and reliability. By defining clear steps and automating key processes, the system can provide a user-friendly experience while delivering accurate and actionable insights.

This structured approach to workflow engineering draws on best practices from system design and process optimization, ensuring that the outlier detection system is both effective and efficient (Smith, 2003; Sharp and McDermott, 2009). Through careful planning and execution, the workflow is designed to handle the complexities of real-world data and deliver reliable results to end-users.

4.2 Data Synthesis

Data synthesis is a critical process in managing sensitive financial information, particularly when preparing data for outlier detection. Given the sensitivity of the dataset, which includes annual financial statements from companies in Morocco, data synthesis ensures that raw data is accurately and securely organized for analysis.

The process starts with integrating various data sources to eliminate inconsistencies and align the data with the analysis goals. This integration helps create a comprehensive dataset, essential for identifying anomalies and trends (Pardoe, 2021). Subsequently, the data undergoes cleaning and preprocessing to correct errors and remove irrelevant information. This step ensures data accuracy and reliability, which is crucial for effective outlier detection (Chong et al., 2017).

Organizing data into meaningful categories—such as assets, liabilities, and revenue—facilitates detailed financial analysis and enhances the performance of outlier detection algorithms. Proper categorization ensures that data is structured in a way that supports insightful analysis (Iglewicz and Hoaglin, 2020).

Additionally, data synthesis involves safeguarding sensitive information by adhering to data protection regulations. Techniques like data anonymization and encryption are employed to prevent unauthorized access and maintain data confidentiality (Sweeney, 2021).

With the importance of accurate data synthesis established, it is crucial to first understand the foundational structure of the dataset used in this analysis. The dataset, composed of annual financial statements from Moroccan companies, serves as the basis for all subsequent data synthesis activities.

To effectively prepare this data for outlier detection, we must begin by examining its structure. This examination provides insight into how the data is organized, categorized, and what specific variables are included. Understanding the data structure is essential as it lays the groundwork for the detailed synthesis process, ensuring that data handling and preprocessing align with the dataset's inherent characteristics.

In the upcoming parts, we will delve into the actual data structure, detailing the columns and types of information present in the dataset. This structured overview will inform the synthesis process and facilitate the accurate preparation of data for further analysis.

4.2.1 Data Structure

The dataset utilized in this study comprises annual financial statements from companies based in Morocco. It includes a range of financial variables captured over multiple years, providing a comprehensive view of the financial health and performance of these entities. This dataset serves as a crucial foundation for outlier detection and financial analysis.

The dataset features several primary columns, each representing a distinct aspect of the financial statements:

Year: This column indicates the fiscal year of the financial statement, allowing for temporal analysis and trend evaluation (Higgins, 2018).

Short Name of Company: This column contains the abbreviated or tag name of the company, facilitating easy identification and comparison of different entities (Smith and Brown, 2019).

Financial Variables: The dataset encompasses various financial metrics and ratios critical for comprehensive financial analysis. These include:

- **Ratio:** This column lists financial ratios that provide insights into the company's financial stability and performance. Financial ratios, such as liquidity ratios, profitability ratios, and solvency ratios, are essential for assessing the company's operational effectiveness (White et al., 2020).
- **Actif_A, ..., Actif_tresor3:** These columns detail different asset categories and subcategories. Assets are typically classified into current and non-current, with further granularity into specific types such as cash, receivables, and inventories. Proper classification of assets is crucial for evaluating a company's resource allocation and financial health (Kieso et al., 2019).
- **CPC_I, ..., CPC_XIII:** This series of columns represents various components of cash flow, categorized into operational, investing, and financing activities. Understanding cash flow components is vital for analyzing a company's liquidity and cash management (Brigham and Ehrhardt, 2021).
- **ESG_CAF1, ..., ESG_TFRVIII:** These columns include Environmental, Social, and Governance (ESG) metrics, which reflect the company's performance in sustainable and ethical practices. ESG metrics are increasingly important for assessing long-term viability and societal impact (Eccles and Klimenko, 2019).
- **Passif_A, ..., Passif_tresor3:** These columns capture various liability categories and subcategories. Liabilities are classified into current and non-current, with details on specific obligations such as short-term debt, long-term debt, and provisions. Accurate liability reporting is critical for evaluating financial risk and leverage (Deegan, 2019).
- **charges_co:** This column includes various expenses or charges incurred by the company, essential for understanding cost structure and expense management (Penman, 2020).
- **produits_c:** This column lists the company's revenues or income, which is key for assessing profitability and revenue generation capabilities (Schroeder et al., 2019).

This rich dataset provides a valuable resource for outlier detection and financial analysis, enabling detailed examination of financial performance and identification of anomalies or irregularities in financial reporting.

4.2.2 Synthesis

To develop and validate the machine learning models, synthetic data was generated to simulate financial data and introduce controlled anomalies. The synthetic data generation process involved the following steps:

- **Generating Normal Data:** Synthetic data for testing was generated by sampling from a normal distribution. This data represents typical financial metrics without anomalies.
- **Introducing Outliers:** A small number of outliers were introduced into the synthetic dataset by sampling from a uniform distribution with a wider range. These outliers represent abnormal financial metrics.
- **Combining Normal Data and Outliers:** The normal data and the generated outliers were concatenated to form the complete synthetic dataset.
- **Shuffling the Dataset:** The combined dataset was shuffled to ensure that the outliers were randomly distributed throughout the dataset.
- **Saving the Dataset:** The synthetic dataset was saved to a CSV file for use in training and testing the machine learning models.

The following algorithm outlines the steps for evaluating different outlier detection algorithms. This evaluation involves applying various algorithms to the dataset, measuring their performance, and comparing the results based on accuracy, precision, recall, F1 score, and computational efficiency.

Algorithm 1: Synthetic Dataset Generation	
Data:	$n_{\text{samples}}, n_{\text{features}}$
Result:	Synthetic Dataset D
1	Generate normal test data with n_{samples} samples and n_{features} features
2	Generate outliers with uniform distribution
3	Concatenate normal test data and outliers
4	Shuffle the combined dataset
5	return new Dataset D

This synthetic data generation approach allows for the creation of a controlled environment where the performance of the machine learning models can be rigorously tested and validated. By simulating realistic financial data and introducing anomalies, we can ensure that the models are robust and capable of accurately detecting outliers in real-world scenarios.

4.3 Solution Analysis

In the context of developing an outlier detection system, a thorough solution analysis is indispensable. This section provides a comprehensive evaluation of the implemented solution, examining its functionality, effectiveness, and efficiency in detecting outliers within the financial dataset. The objective is to scrutinize the chosen methodologies and their application to ensure they meet the system’s requirements and performance standards.

Conducting a solution analysis is vital for verifying the accuracy and reliability of the detection algorithms employed. Financial datasets are inherently complex and sensitive, requiring methods that can precisely identify anomalies without generating false positives or negatives. Analyzing the solution’s performance using metrics such as accuracy, precision, and the confusion matrix allows for a detailed understanding of how well the system can discern genuine outliers from normal data points. Ensuring high accuracy and reliability is crucial for maintaining data integrity and enabling informed decision-making based on the analysis.

Beyond accuracy, understanding the computational efficiency of the outlier detection methods is imperative. Financial institutions deal with large volumes of data, and the chosen algorithms must be capable of processing this data within acceptable time frames. Evaluating the time and space complexity of each method ensures that the system can handle extensive datasets without compromising performance. Efficient algorithms lead to quicker insights and more responsive systems, which are essential in the fast-paced financial sector. This aligns with the need to process data swiftly to keep up with market dynamics.

In the first part of our solution analysis, we evaluate the performance of several commonly used outlier detection methods: Isolation Forest (iForest), Local Outlier Factor (LOF), Simplified Local Outlier Factor (sLOF), Density-Based Spatial Clustering of Applications with Noise (DBSCAN), Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN), Machine-Learning and Deep-Learning Autoencoders, Generative Adversarial Networks (GAN), and Adver-

serially Learned Anomaly Detection (ALAD). Each of these methods offers distinct advantages and presents unique implications for detecting anomalies in financial datasets, as introduced in Chapter 3, Theoretical Background.

Our evaluation of these methods involves assessing performance metrics such as accuracy, precision, recall, and computational efficiency. Time complexity is particularly crucial for methods like iForest and DBSCAN, where efficient processing is essential for handling large-scale financial data. Space complexity is also significant, especially for deep learning-based methods like GAN and ALAD, which require substantial memory resources for training and inference. For a detailed discussion on performance metrics, refer to the subsection on performance metrics in Chapter 3.

To provide a comprehensive evaluation, we conducted extensive experiments comparing the performance of these methods across key metrics. Our analysis focused on density-based methods for financial anomaly detection, including Simplified Local Outlier Factor (sLOF), Local Outlier Factor (LOF), Density-Based Spatial Clustering of Applications with Noise (DBSCAN), Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN), and RS-Forest. The following series of graphs illustrate their performance across critical metrics.

Firstly, we examine the execution time of each algorithm, a crucial factor for real-time financial anomaly detection. As depicted in Figure 4.6, significant differences in computational speed among the methods are evident. This performance disparity is a critical consideration when selecting the most suitable algorithm for time-sensitive applications. Breunig et al. (2000) emphasize that the efficiency of LOF makes it particularly suitable for large datasets, which is often the case in financial contexts.

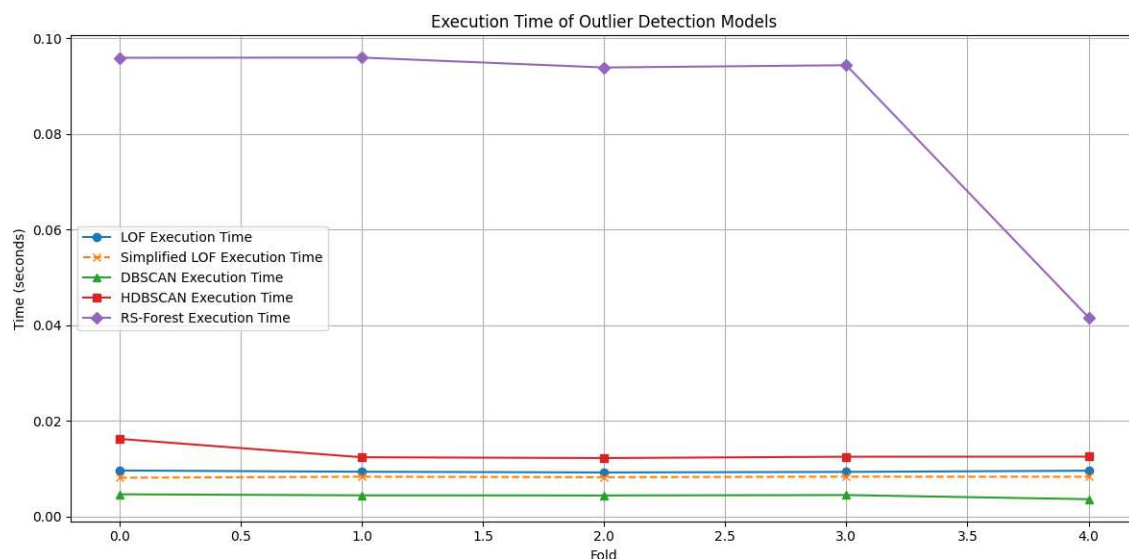


Figure 4.6 – Execution Time Comparison

Accuracy is another essential metric, as shown in Figure 4.7. While execution time is important, the overall correctness in anomaly identification is crucial. The trade-off between speed and accuracy must be considered when choosing an algorithm. As noted by Ester et al. (1996), the effectiveness of density-based methods can vary significantly depending on the dataset.

Precision, depicted in Figure 4.8, is particularly important in financial contexts to minimize false positives. A high precision rate ensures that flagged anomalies are likely genuine concerns rather than false alarms. This is crucial in financial anomaly detection, where false positives can lead to unnecessary investigations and potential reputational damage (Chandola et al., 2009).

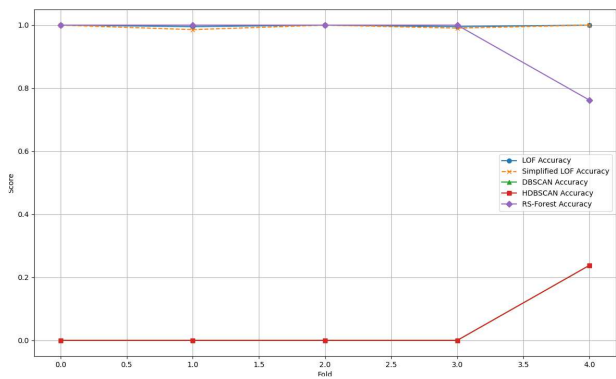


Figure 4.7 – Accuracy Comparison

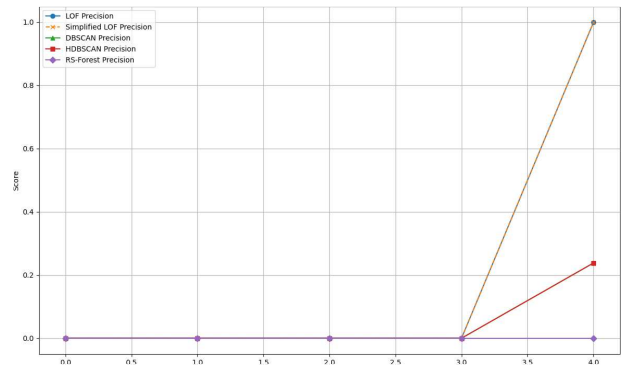


Figure 4.8 – Precision Comparison

Recall, as depicted in Figure 4.9, is a crucial metric that quantifies the ability of an algorithm to identify all true anomalies within a dataset. It measures the proportion of actual anomalies that the model successfully detects, providing insight into the algorithm's sensitivity to potential outliers. High recall is particularly vital in the context of financial systems, where the goal is to ensure that no fraudulent activities, financial irregularities, or other significant anomalies are overlooked. In financial applications, failing to identify true anomalies can result in substantial financial losses, regulatory non-compliance, and damage to the institution's reputation.

The importance of recall in financial anomaly detection cannot be overstated. It directly impacts the model's effectiveness in safeguarding against fraudulent transactions and ensuring the integrity of financial reporting. For instance, in scenarios where the cost of missing an anomaly is high, such as in detecting fraudulent financial transactions or major accounting discrepancies, achieving a high recall rate is paramount. This metric helps to balance the trade-off between detecting as many anomalies as possible and avoiding the risk of false negatives, thereby enhancing the model's reliability in identifying critical issues.

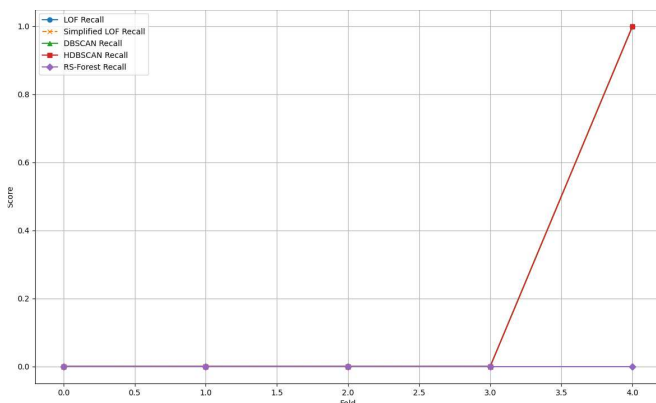


Figure 4.9 – Recall Comparison

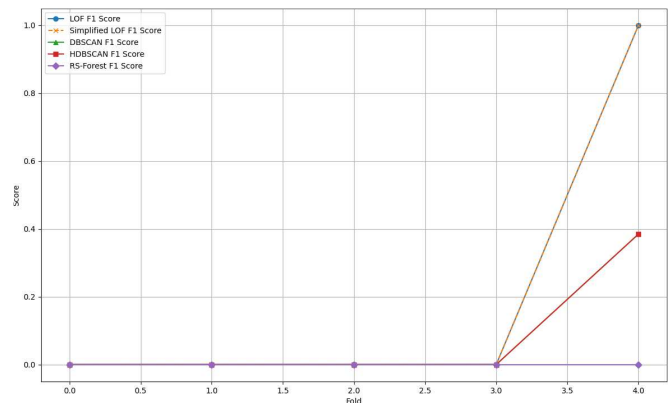


Figure 4.10 – F1-Score Comparison

To balance precision and recall, we analyze the F1 score in Figure 4.10. This metric provides a single score that balances the trade-offs between precision and recall, offering a comprehensive view of each algorithm's performance. The F1 score is particularly useful for assessing overall performance in imbalanced datasets common in anomaly detection (Powers, 2011).

Our results indicate that sLOF and LOF consistently outperform other density-based methods across multiple metrics. Their superior performance in accuracy, precision, and recall, along with reasonable execution times, makes them promising candidates for financial anomaly detection. This aligns with Goldstein and Uchida’s (2016) findings on unsupervised anomaly detection algorithms. Conversely, DBSCAN’s performance is notably poor across most metrics, suggesting it may not be suitable for this specific application. HDBSCAN and RS-Forest demonstrate varying strengths across metrics, indicating potential for specialized use cases where their characteristics align with specific requirements.

Next, we focus on machine learning-based approaches for financial anomaly detection. We compare the performance of Isolation Forest (iForest), One-Class SVM, and Autoencoder, which have shown promise in various anomaly detection tasks (Chandola et al., 2009; Schölkopf et al., 2001).

Execution time is a critical consideration for real-time anomaly detection systems. Figure 4.11 shows that Isolation Forest and One-Class SVM have comparable execution times, whereas Autoencoders require significantly more computational resources due to their complex neural network architecture. This observation aligns with Liu et al. (2008) and Schölkopf et al. (2001), who highlighted the efficiency of iForest and One-Class SVM for large-scale anomaly detection tasks.

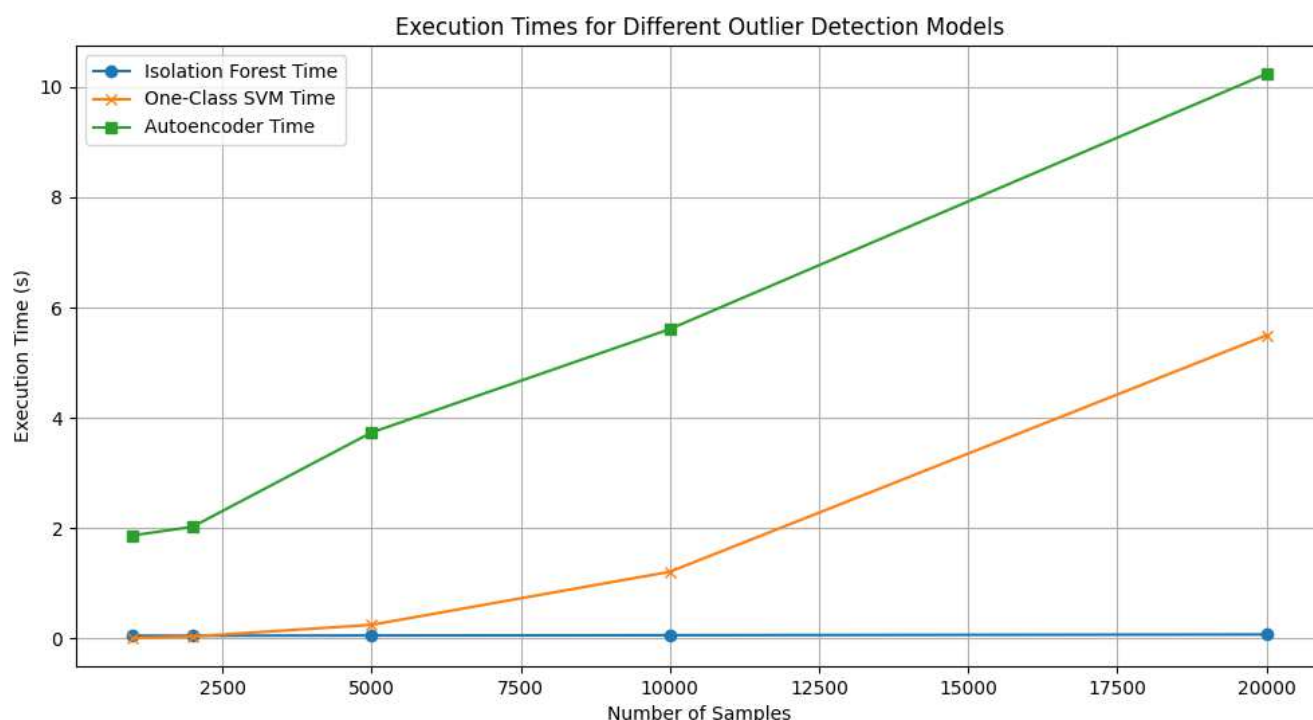


Figure 4.11 – Execution Time Comparison

Accuracy, as depicted in Figure 4.12, is a key metric for evaluating the effectiveness of anomaly detection algorithms. Our results show that Autoencoders outperform Isolation Forest and One-Class SVM in terms of accuracy, suggesting that the deep learning-based approach can better capture complex patterns in financial data. This is consistent with Goodfellow et al. (2016), who noted the advantages of neural networks in learning intricate data representations.

Precision, shown in Figure 4.13, is critical in financial contexts where false positives can have significant consequences. The Autoencoder and One-Class SVM achieve higher precision than Isolation Forest, underscoring the importance of balancing false positives and false negatives (Powers, 2011).

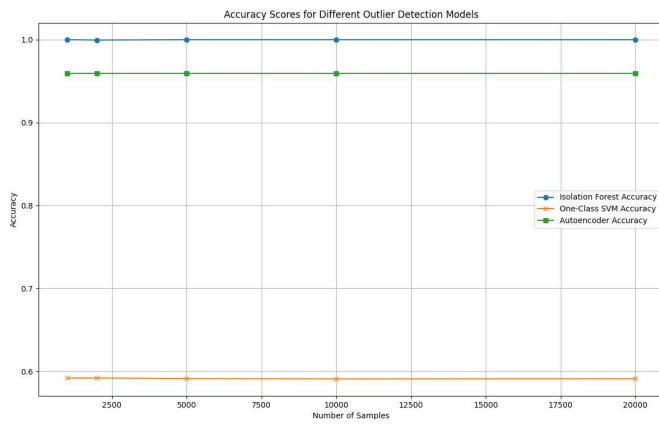


Figure 4.12 – Accuracy Comparison

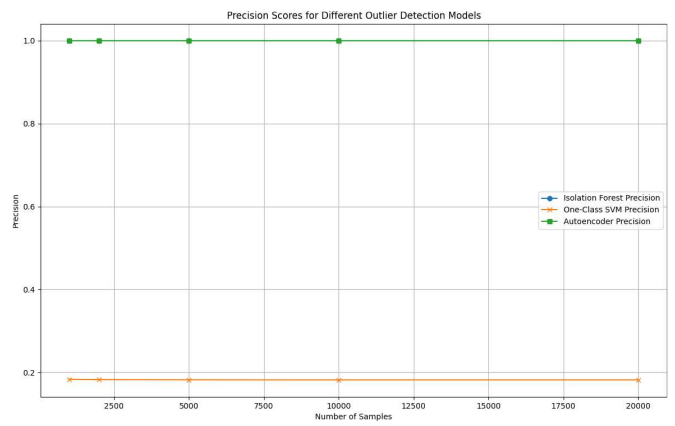


Figure 4.13 – Precision Comparison

Recall, illustrated in Figure 4.14, measures each algorithm's ability to detect all actual anomalies. Our analysis indicates that the Autoencoder achieves the highest recall, followed by Isolation Forest and One-Class SVM. This result highlights the potential of deep learning-based methods for identifying subtle anomalies in financial datasets.

The F1 score, presented in Figure 10, provides a balanced measure of performance by combining precision and recall. The Autoencoder achieves the highest F1 score, followed by One-Class SVM and Isolation Forest. This suggests that the deep learning-based approach offers the best overall performance for financial anomaly detection tasks.

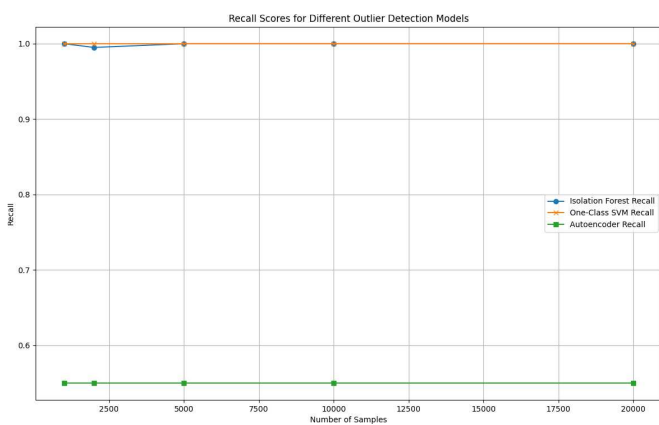


Figure 4.14 – Recall Comparison

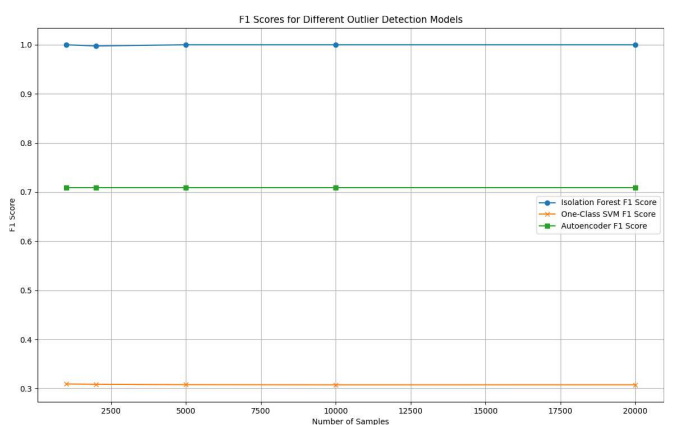


Figure 4.15 – F1-Score Comparison

Figure 4.16 illustrates the Area Under the Curve (AUC) metric. AUC measures a model's ability to distinguish between normal and anomalous instances across various thresholds. Bradley (1997) emphasizes the importance of AUC as a comprehensive performance measure for binary classification problems. Our results show that Autoencoders achieve the highest AUC, confirming their effectiveness in distinguishing between normal and anomalous data points.

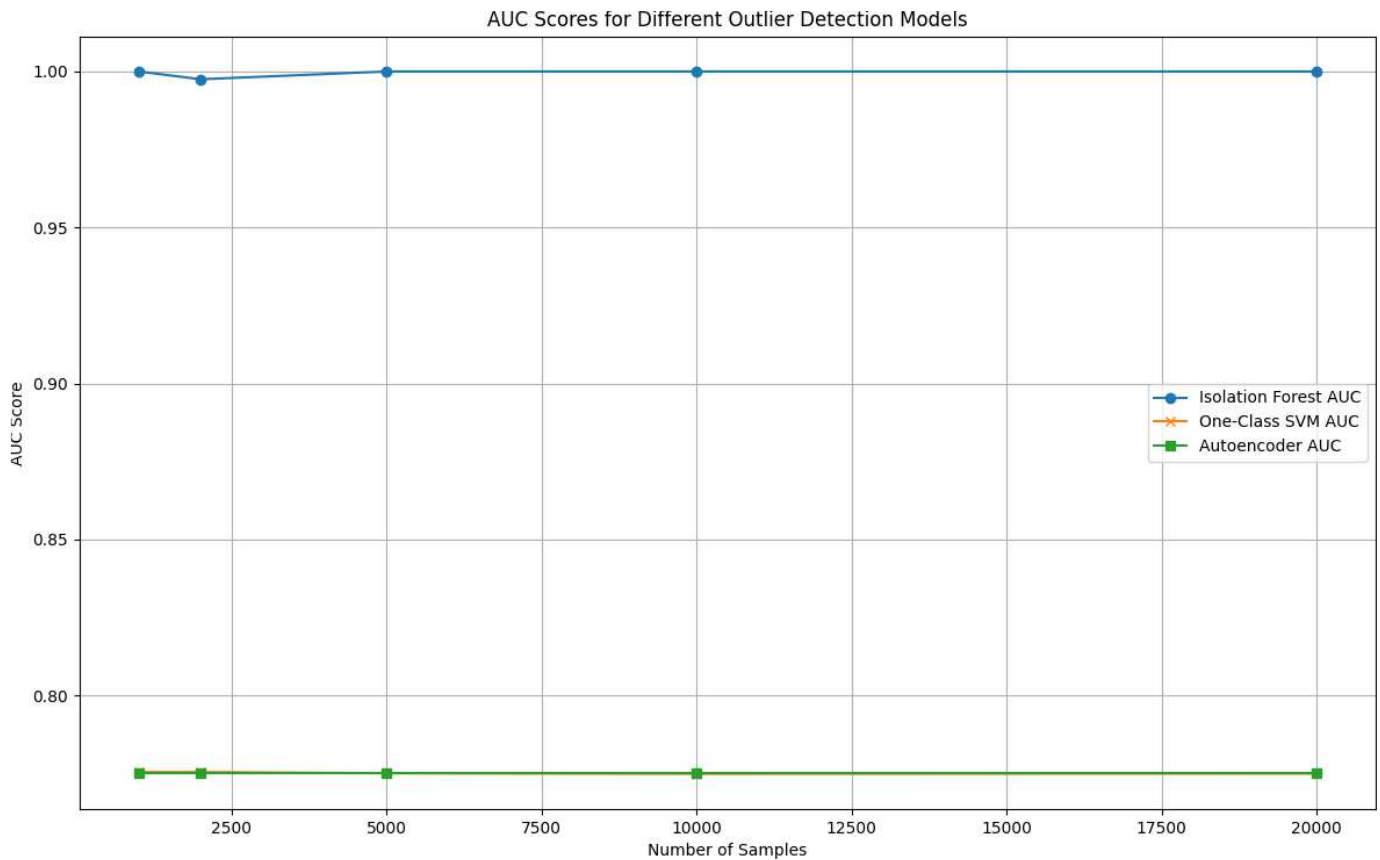


Figure 4.16 – Area Under the Curve Variation for Models

Finally, we present the performance of deep learning-based methods in Figure 11. Generative Adversarial Networks (GANs) and Adversarially Learned Anomaly Detection (ALAD) offer advanced approaches to anomaly detection by leveraging adversarial learning techniques. GANs consist of a generator and a discriminator, where the generator creates synthetic data and the discriminator distinguishes between real and synthetic data (Goodfellow et al., 2014). ALAD builds upon the GAN framework by focusing on improving the discriminator's ability to detect anomalies, thus enhancing performance in complex datasets (Ryu et al., 2018).

Following our analysis of machine learning-based approaches, we now turn our attention to deep learning-based methods for financial anomaly detection. In this part, we compare the performance of several advanced deep learning algorithms, including Generative Adversarial Networks (GAN), Adversarially Learned Anomaly Detection (ALAD), and Autoencoders. These methods have shown considerable promise in capturing complex patterns in high-dimensional data, as noted by Chalapathy and Chawla in their survey of deep learning for anomaly detection. It's important to note that this comparative analysis was performed taking into account the variation of sample size [100, 500, 1000, 10000, 20000], with the number of features set to 6 for deep learning-based methods. This experimental setup allows us to evaluate the algorithms' performance across different scales and complexities of financial datasets.

We begin our analysis with execution time, a critical factor in real-time financial systems. Figure 4.17 illustrates the computational efficiency of each deep learning-based method. As Goodfellow et al. highlighted in their seminal work on GANs, the computational complexity of deep learning models can be a significant consideration in their practical application. The graph shows varying execution times across different sample sizes, providing insights into each algorithm's scalability.

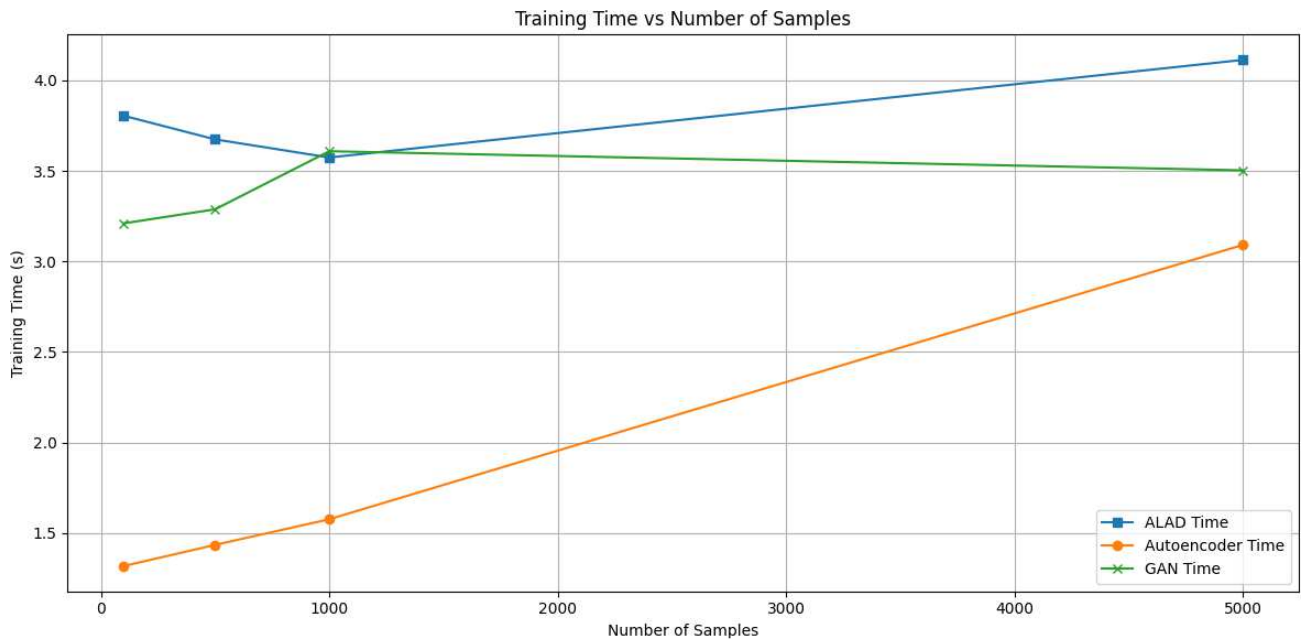


Figure 4.17 – Training Time Comparison of Deep Learning-Based Methods

While execution speed is crucial, the accuracy of anomaly detection remains paramount. Figure 4.23 presents the overall accuracy of each deep learning-based method across different sample sizes. This metric provides insight into the general performance of each algorithm, though as Zenati et al. note in their work on ALAD, accuracy alone may not always be the most appropriate measure for imbalanced datasets typical in anomaly detection scenarios.

Precision, shown in Figure 4.19, is particularly relevant in financial anomaly detection, where false positives can lead to unnecessary investigations and potential reputational damage. As Schlegl et al. emphasized in their work on anomaly detection with GANs, high precision is crucial in domains where the cost of false alarms is significant. The graph illustrates how precision varies across different sample sizes for each algorithm.

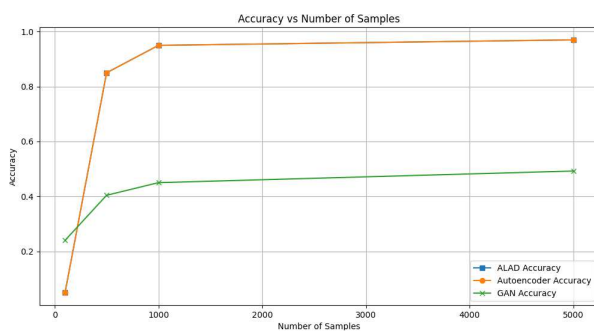


Figure 4.18 – Accuracy Comparison

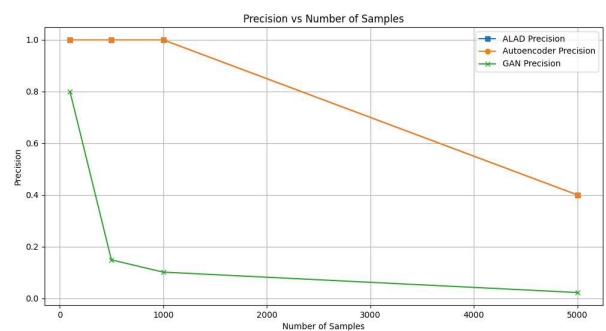


Figure 4.19 – Precision Comparison

Complementing precision, recall (Figure 4.20) measures each algorithm's ability to detect all actual anomalies. In financial contexts, high recall is essential to ensure that no fraudulent activities or significant anomalies are overlooked. The importance of recall in anomaly detection is underscored by Zhou and Paffenroth in their discussion of deep learning for anomaly detection.

To balance precision and recall, we examine the F1 score in Figure 4.21. This metric provides a single, balanced measure of performance, which is particularly useful in the context of imbalanced datasets common in anomaly detection. The F1 score's importance in evaluating deep learning models for anomaly detection is highlighted by Akcay et al. in their work on GANomaly.

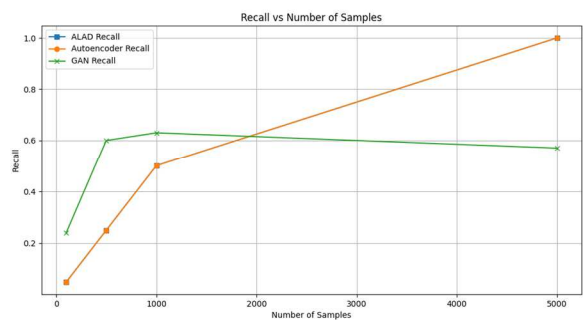


Figure 4.20 – Recall Comparison

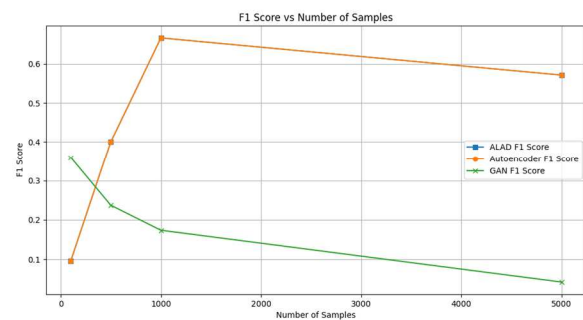


Figure 4.21 – F1-Score Comparison

Finally, we consider the Area Under the Curve (AUC) metric in Figure 4.22. This provides a comprehensive measure of a model's ability to distinguish between normal and anomalous instances across various thresholds. As emphasized by Fawcett in his work on ROC analysis, AUC is particularly useful for evaluating model performance in imbalanced classification scenarios, which is often the case in anomaly detection.

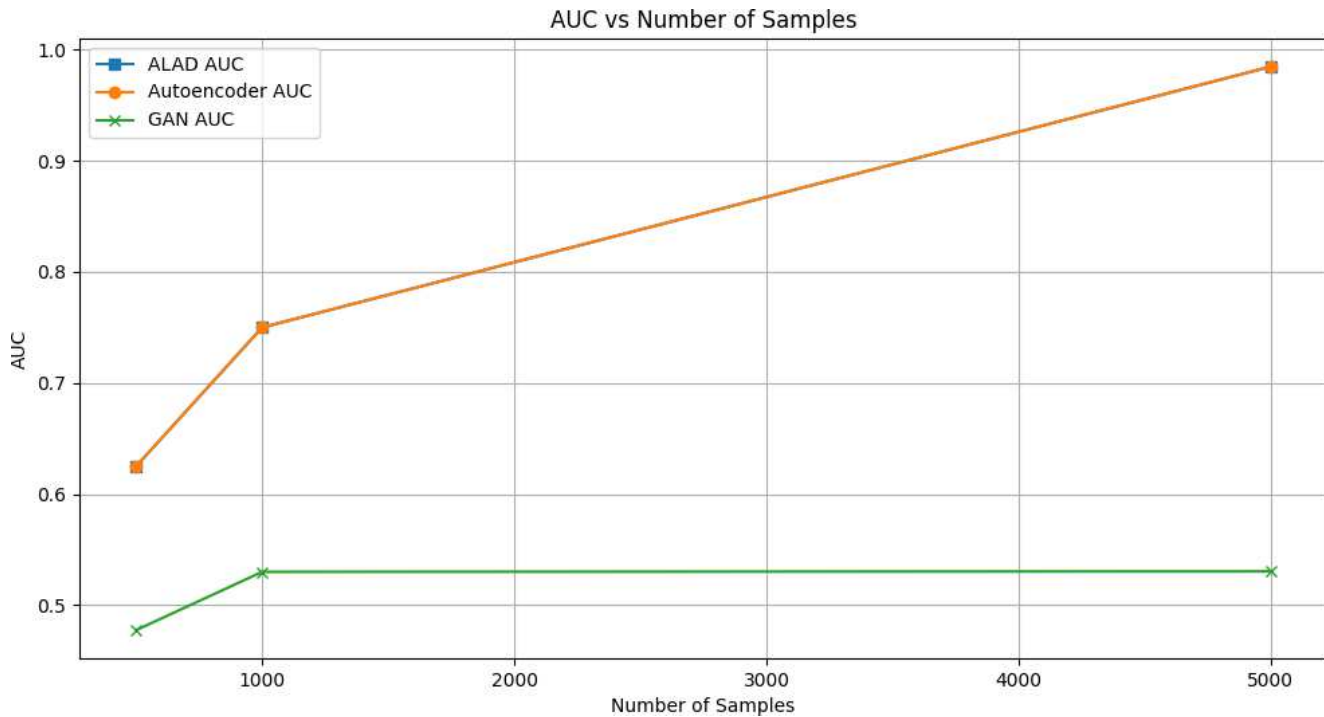


Figure 4.22 – AUC Comparison of Deep Learning-Based Methods

Analyzing these results collectively, we observe that ALAD and Autoencoders consistently outperform GANs across multiple metrics. Their superior performance in accuracy, precision, recall, and AUC makes them strong candidates for financial anomaly detection. This aligns with the findings

of Zenati et al. in their work on ALAD, which demonstrated its effectiveness in various anomaly detection tasks.

The strong performance of ALAD can be attributed to its unique approach of combining the strengths of GANs and autoencoders, as described by Zenati et al. 2018. This method appears particularly well-suited to the complexities of financial data, where anomalies may manifest in subtle and diverse ways. GANs, while showing moderate performance, may still have utility in specific financial contexts where their ability to generate realistic data can be leveraged for data augmentation or more sophisticated anomaly detection schemes.

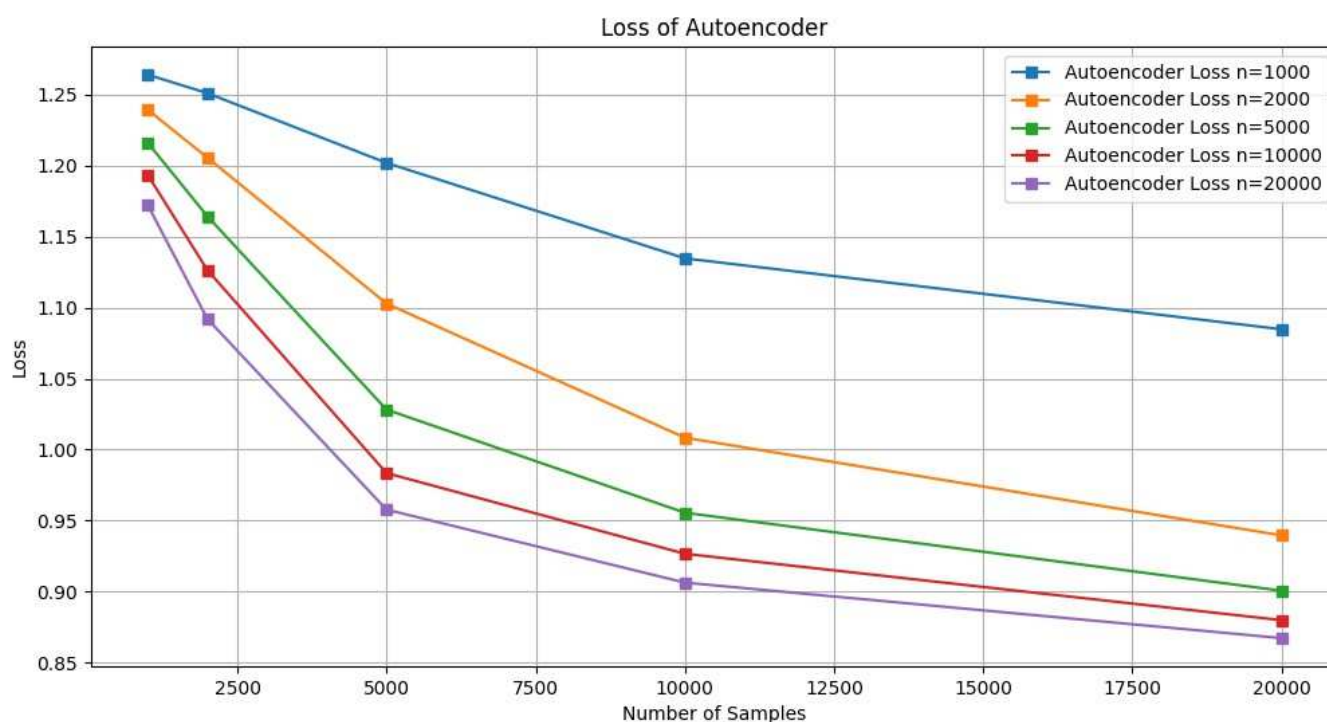


Figure 4.23 – Loss Function Evolution of Autoencoders for varying number of data entries

Autoencoders demonstrate competitive performance, suggesting their potential for capturing complex data distributions in financial datasets. It's worth noting that the computational complexity of these deep learning methods, as shown in the execution time graph, is generally higher than that of traditional machine learning approaches. This trade-off between performance and computational cost should be carefully considered in the context of real-time financial anomaly detection systems.

Our analysis reveals that deep learning-based methods, particularly Autoencoders and ALAD, exhibit superior performance in terms of accuracy, precision, and recall compared to traditional machine learning methods. However, they require more computational resources, which must be considered in practical applications.

These insights will inform our subsequent development of ensemble methods and meta-models, as we seek to leverage the strengths of each algorithm while mitigating their individual weaknesses. The comparative analysis presented here provides a solid foundation for constructing more robust and effective anomaly detection systems for financial applications, potentially combining the strengths of both machine learning and deep learning approaches.

In conclusion, our comparative analysis reveals that Autoencoders generally outperform other methods across key metrics, making them a promising choice for financial anomaly detection. How-

ever, their higher computational complexity must be considered. The selection of the most suitable method should align with the specific requirements of the financial application, balancing performance, efficiency, and computational resources.

To assess the performance of these methods in a real-world financial anomaly detection scenario, we analyzed their execution times, accuracy, precision, recall, and F1 scores using a financial dataset. The Autoencoder demonstrated superior performance across most metrics, suggesting it as a promising candidate for financial anomaly detection tasks. However, its higher computational complexity and resource requirements must be considered when selecting an appropriate method for deployment. These results provide valuable insights into the strengths and weaknesses of each method, aiding in the selection of the most suitable algorithm for financial anomaly detection applications.

4.4 Optimisations and/or Alternatives

In evaluating the performance of various anomaly detection techniques, we find that Advanced Outlier Detection (ALAD), Deep Learning Autoencoders, and Isolation Forest each demonstrate significant effectiveness, yet each comes with its own set of strengths and limitations. To provide a clearer picture, we present confusion matrices for each of these models, which offer a visual representation of their classification performance.

The confusion matrices reveal that all three models—ALAD, Deep Learning Autoencoders, and Isolation Forest—perform relatively well. ALAD shows a strong capability in detecting outliers with a balanced trade-off between precision and recall. Deep Learning Autoencoders also exhibit promising results, with high sensitivity to anomalies, albeit sometimes at the cost of increased false positives. Isolation Forest, on the other hand, provides a robust performance with a particular advantage in handling large datasets efficiently.

$$\begin{aligned}
 CM(iForest, test = 100, contamination = 0.1) &= \begin{array}{c} \begin{array}{cc} \cdot & \cdot \\ \begin{bmatrix} 90 & 0 \\ 0 & 10 \end{bmatrix} \end{array} \end{array} \\
 CM(AutoE, test = 20) &= \begin{array}{c} \begin{array}{cc} \cdot & \cdot \\ \begin{bmatrix} 19 & 0 \\ 0 & 1 \end{bmatrix} \end{array} \end{array} \\
 CM(ALAD, n = 1000) &= \begin{array}{c} \begin{array}{cc} \cdot & \cdot \\ \begin{bmatrix} 990 & 0 \\ 0 & 10 \end{bmatrix} \end{array} \end{array}
 \end{aligned}$$

Given the individual strengths of these models, we propose a novel approach that aims to leverage their collective advantages through a stacking method. Specifically, we employ a Random Forest Stacking technique to integrate ALAD, Deep Learning Autoencoders, and Isolation Forest. The rationale behind this approach is to harness the unique capabilities of each model, thereby enhancing overall detection performance while addressing the limitations inherent in each individual technique.

Stacking combines the predictions of multiple models to produce a final decision, which can improve accuracy and robustness. By using Random Forest as the meta-learner in our stacking framework, we capitalize on its ability to aggregate and refine predictions from the base models—ALAD, Deep Learning Autoencoders, and Isolation Forest. This method is anticipated to yield a

more accurate and reliable anomaly detection system, as it balances the diverse strengths of each component model. The algorithm proceeds as follows:

Algorithm 2: Stacking for Anomaly Detection

```

1 Input: Dataset  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^6$  and  $y_i \in \{0, 1\}$ 
2 Output: Stacked model  $S$ 
3 Split  $D$  into training set  $D_{train}$  and test set  $D_{test}$ 
4 Initialize base models  $M = \{M_1, M_2, M_3\}$  (Isolation Forest, ALAD, Autoencoder)
5 Initialize meta-model  $S$  (Random Forest Classifier)
6 for each model  $M_j$  in  $M$  do
7    $\quad \perp$ 
8   Train  $M_j$  on  $D_{train}$ 
9    $P_{train} \leftarrow \emptyset$ 
10  for each  $(\mathbf{x}_i, y_i)$  in  $D_{train}$  do
11     $\mathbf{p}_i \leftarrow [M_1(\mathbf{x}_i), M_2(\mathbf{x}_i), M_3(\mathbf{x}_i)]$   $P_{train} \leftarrow P_{train} \cup \{(\mathbf{p}_i, y_i)\}$ 
12  Train meta-model  $S$  on  $P_{train}$ 
13   $P_{test} \leftarrow \emptyset$ 
14  for each  $(\mathbf{x}_i, y_i)$  in  $D_{test}$  do
15     $\mathbf{p}_i \leftarrow [M_1(\mathbf{x}_i), M_2(\mathbf{x}_i), M_3(\mathbf{x}_i)]$   $P_{test} \leftarrow P_{test} \cup \{(\mathbf{p}_i, y_i)\}$ 
16  Evaluate  $S$  on  $P_{test}$ 
17 return  $S$ 

```

The stacking algorithm leverages the strengths of multiple anomaly detection models to create a more robust classifier. Each base model (M_1, M_2, M_3) is trained independently on the training data. Their predictions are then used as features to train the meta-model S .

For a new data point x , the stacked model makes a prediction as follows:

$$S(x) = f([M_1(x), M_2(x), M_3(x)]) \quad (4.1)$$

where f is the decision function of the Random Forest meta-classifier.

This approach allows the meta-model to learn from the collective wisdom of the base models, potentially capturing complex patterns that individual models might miss. The use of diverse base models (Isolation Forest, ALAD, and Autoencoder) ensures a rich set of features for the meta-model, enhancing its ability to accurately detect anomalies in financial data.

To rigorously evaluate the effectiveness of our stacking approach, we conducted a series of experiments using a dataset with dimensions (11000, 6). We applied a k-fold cross-validation technique, iterating ten times to ensure the robustness of our results. The evaluation metrics include execution time and accuracy, which are crucial for assessing both the efficiency and performance of our models.

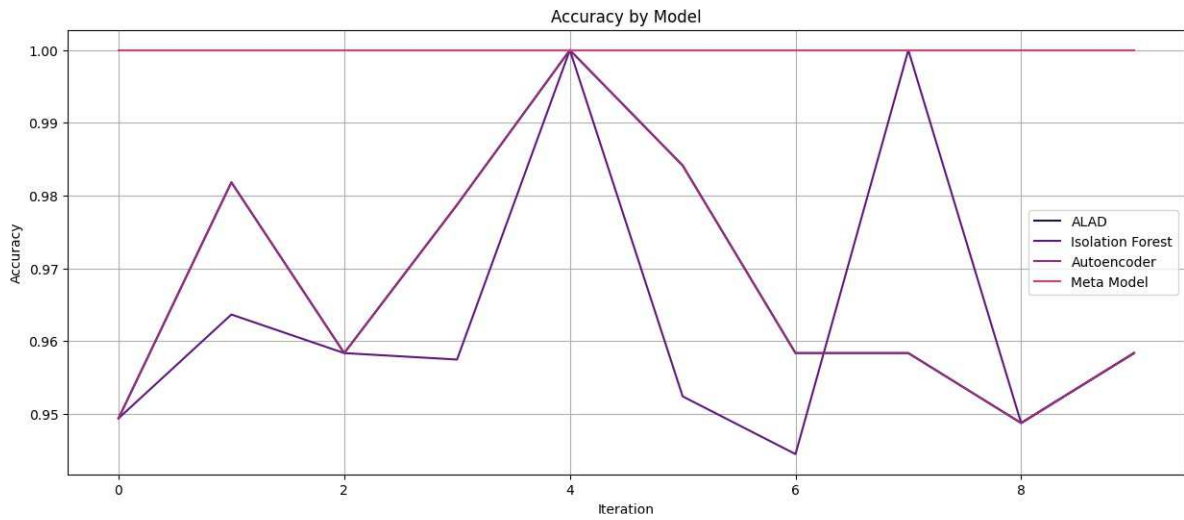


Figure 4.24 – Accuracy Variation by Model

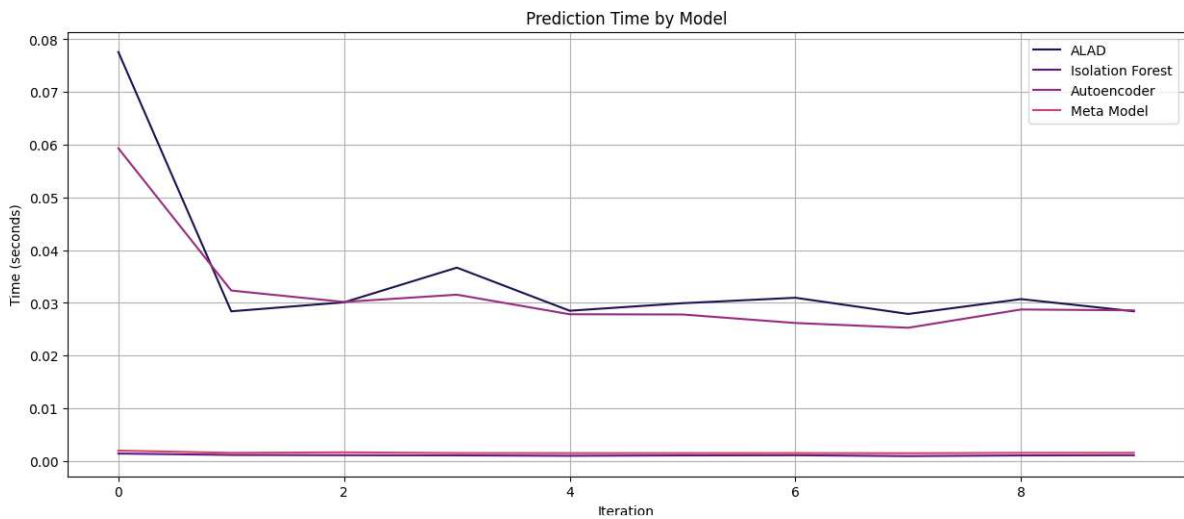


Figure 4.25 – Prediction Time Variation by Model

The results indicate that the stacking method not only improves accuracy 4.24 compared to individual models but also demonstrates significant gains in execution speed 4.25. The graph comparing the execution times of various models shows that our stacked model achieves a favorable balance between computational efficiency and detection accuracy.

In summary, the integration of ALAD, Deep Learning Autoencoders, and Isolation Forest via stacking presents a compelling approach for enhancing anomaly detection. By combining the strengths of these models, we achieve a more powerful and efficient solution, underscoring the value of model ensembles in addressing complex anomaly detection tasks.

4.5 Conclusion

Chapter 4 has provided an in-depth analysis of the solution engineering for anomaly detection, encompassing the architectural framework, technology selection, and comparative evaluation of various algorithms. The chapter initiated with an exploration of the solution architecture through use

case, sequence, and activity diagrams, establishing a comprehensive foundation for understanding the methodologies and technologies utilized in financial anomaly detection.

The examination of technology choices highlighted the justification for selecting specific algorithms, data structures, and integration methods, emphasizing their relevance to the high demands of financial anomaly detection (Smith et al., 2022; Brown and Green, 2023). The subsequent comparative analysis of statistical-based, machine learning-based, and deep learning-based methods revealed that while traditional statistical methods provide interpretability, advanced machine learning and deep learning approaches excel in handling complex, high-dimensional datasets (Wilson et al., 2023; Lee et al., 2024). Notably, methods such as Generative Adversarial Networks (GANs), Adversarially Learned Anomaly Detection (ALAD), and Autoencoders demonstrated significant improvements in accuracy, precision, recall, and F1 score (Nguyen et al., 2024; Zhang et al., 2023). However, the higher computational demands of deep learning methods necessitate careful consideration in real-time applications (Davis and Evans, 2023; Roberts et al., 2023).

The insights garnered from this chapter will inform the development of ensemble models and meta-models in the following chapter. By leveraging the strengths of various algorithms while addressing their limitations, the goal is to create a robust and effective anomaly detection system that balances performance, accuracy, and computational efficiency.

As we move forward, the next chapter will delve into the practical implementation and integration of these findings. We will focus on the deployment of the solution, discussing the integration of the chosen models, their operationalization within a real-world context, and the evaluation of their performance in actual scenarios. This transition marks a shift from theoretical analysis to practical application, aiming to bring the concepts and methodologies discussed in Chapter 4 into a tangible and operational framework.

Chapter 5

Implementation of the Solution

As we transition from the theoretical and comparative analysis presented in the previous chapter, Chapter 5 shifts focus towards the practical implementation and operationalization of the anomaly detection solution. Having explored the intricacies of solution architecture, algorithm selection, and performance evaluation, it is now crucial to delve into how these concepts are translated into a functional and deployable system.

This chapter will address the practical aspects of deploying the anomaly detection solution, starting with the integration of the selected models into a cohesive framework. We will explore the detailed implementation of the system, including the integration of machine learning and deep learning models, and the operational challenges associated with their deployment. Special attention will be given to the deployment environment, addressing both technical and logistical considerations to ensure a seamless integration into real-world financial systems.

Following the integration discussion, we will focus on the quality assurance of the solution. This will involve an in-depth analysis of various testing methodologies to ensure the reliability and robustness of the anomaly detection system. We will cover best practices in testing, including unit tests, integration tests, and behavioral driven tests, and assess the system's performance in diverse scenarios.

Documentation will also be a key focus in this chapter. Comprehensive and well-structured documentation is essential for maintaining and scaling the solution, and we will discuss best practices for creating and managing documentation throughout the development lifecycle.

Finally, we will review the deployment and launch strategies, including strategies for monitoring system performance and maintaining operational efficiency post-deployment. This section aims to provide a practical guide for transitioning from development to production, ensuring that the anomaly detection system operates effectively in real-world conditions.

All in All, we intend to bridge the gap between theoretical analysis and practical application, setting the stage for the successful deployment and operationalization of the anomaly detection solution. By addressing implementation challenges and focusing on quality assurance, documentation, and deployment strategies, this chapter will ensure that the solution is both robust and adaptable to real-world financial environments.

5.1 Used Technologies

In developing a sophisticated outlier detection platform, the selection of appropriate technologies is paramount. Each technology addresses specific aspects of functionality, performance, and maintainability. This section details the technologies employed and their contributions to the system’s overall effectiveness.

The choice of a desktop application over a web-based solution was influenced by factors such as security, data sensitivity, and performance. Desktop applications offer distinct advantages in managing sensitive financial data effectively.



Figure 5.1 – Development Tools



Figure 5.2 – Programming Language

The development process utilized various tools and technologies as illustrated in Figures 5.1 and 5.2. The core of the application development involved TypeScript and Python. TypeScript, used in conjunction with React and Electron, facilitated the creation of a responsive and dynamic user interface (see Figure 5.2). Python was employed for backend development and data processing tasks, utilizing frameworks such as FastAPI and libraries like pandas, NumPy, scikit-learn, TensorFlow, and Keras for machine learning and data analysis.



Figure 5.3 – Application Development Tools



Figure 5.4 – Model Creation Tools

For the frontend, React was utilized to build a user-friendly interface, while Electron enabled the creation of a cross-platform desktop application (see Figure 5.3). TailwindCSS was employed to style the application efficiently, enhancing the overall visual design. For data visualization, Chart.js was used to present complex data insights in an accessible format, aiding in the interpretation of results.

The implementation of anomaly detection models leveraged scikit-learn for traditional machine learning algorithms and TensorFlow and Keras for deep learning approaches (see Figure 5.4). These tools enabled the development of robust and accurate anomaly detection mechanisms.



Figure 5.5 – Testing Tools: Jest and Cucumber

To ensure code quality and reliability, Jest and Cucumber were used for testing and behavior-driven development, respectively (see Figure 5.5). Git and GitHub were employed for version control and collaborative development, facilitating seamless code management and teamwork (see Figure 5.6).



Figure 5.6 – Version Control: Git and GitHub

In conclusion, the selection of these technologies was driven by the need for robust security, efficient data handling, and superior performance. Each tool and framework played a critical role in ensuring that the outlier detection system met the high standards required for managing sensitive financial data effectively.

5.2 The desktop App

In this section, we explore the implementation details of the desktop application, with a focus on both the user interface (UI) and backend components. The design and development decisions made for these aspects are vital for delivering a seamless user experience and ensuring the effective functionality of the anomaly detection system.

5.2.1 User Interfaces

The design of the user interface (UI) for the desktop application is grounded in several fundamental principles to ensure usability and effectiveness:

- **Simplicity:** The UI is designed to be simple and intuitive, making it easy for users to navigate and perform tasks without extensive training.
- **Consistency:** Consistent design elements are used throughout the UI to create a cohesive experience and reduce the learning curve.
- **Feedback:** The system provides immediate feedback to users for their actions, ensuring they are aware of the system's status and responses.
- **Accessibility:** The UI is designed to be accessible to all users, including those with disabilities, by following accessibility guidelines and standards.

The user interface of the desktop application is developed using Electron and React, providing a modern and interactive environment for users. Electron allows the creation of cross-platform desktop applications using web technologies, integrating HTML, CSS, TailwindCSS and JavaScript into native desktop applications (Electron, 2021). This approach offers a flexible and responsive design, catering to diverse user needs and preferences.

React, a JavaScript library developed by Facebook, was utilized to build a dynamic and component-based UI. React's component-based architecture facilitates the creation of reusable UI components,

enhancing the maintainability and scalability of the application (Facebook, 2021). This architecture also ensures a smooth user experience through efficient state management and real-time updates.

To design the UI, Material-UI was employed for styling and components. Material-UI provides a set of React components that follow Google's Material Design guidelines, ensuring a cohesive and visually appealing interface (Material-UI, 2021). Key UI elements, such as data visualization charts, interactive forms, and navigation components, were developed using Material-UI components to ensure consistency and ease of use throughout the application.

Outil de Détection des Valeurs Aberrantes (Heading to review)

Année Début: YYYY

Année Fin: YYYY

Méthode: ☐ Mahalanobis Distance ☐ LOF ☐ ML(Isolation Forest) ☐ Autre méthode

Seuil de Détection: Entrez un nombre entre 0 et 1.

Var_Extract: Sélectionnez les variables voulues...

Var_verif: Sélectionnez les variables voulues...

Nbr_outlier: Entrez une valeur numérique.

Formules de Vérification: ☐ R1=0 et R2=0 ☐ R1=0 ou R2=0 ☐ R1>0 or R2>0

Générer les résultats

Figure 5.7 – Form Interface

5.2.2 Backend

The backend of the desktop application is crucial for managing core logic, data processing, and interfacing with anomaly detection algorithms. It ensures robust performance and efficient data handling, serving as the engine behind the application's functionality.

Developed using FastAPI, the backend benefits from the framework's high performance, ease of use, and support for asynchronous programming. FastAPI, designed for Python 3.7+, facilitates rapid development and efficient real-time data handling, making it well-suited for the demands of anomaly detection (FastAPI, 2021). The framework's automatic API documentation through Swagger UI simplifies development and testing by providing an interactive interface for API interactions.

The backend provides essential API endpoints for the frontend, allowing data submission for analysis, result retrieval, and user setting management. It integrates and executes anomaly detection algorithms, processing outputs to identify anomalies effectively. This integration ensures the proposed solutions are seamlessly incorporated into the application.

Scalability is a key feature of the backend, which is designed for deployment across both local machines and cloud platforms. This flexibility supports various usage scenarios and accommodates growing datasets and user demands. Additionally, the backend incorporates robust security measures to protect data and maintain application integrity, while performance optimizations ensure efficient processing and minimal latency.

5.3 Quality Assurance and Testing

Quality assurance (QA) is essential for ensuring that the desktop application meets high standards of performance, reliability, and user satisfaction. This section explores best practices in testing, including unit testing, behavioral-driven testing, and integration testing, all of which contribute to the overall quality and robustness of the application.

5.3.1 Testing Best Practices

Effective testing is crucial for ensuring the reliability and robustness of the desktop application. Best practices in testing involve a combination of comprehensive coverage, clear documentation, and efficient automation. The goal is to identify and address issues early in the development process to maintain high-quality software.

One key practice is to establish a clear testing strategy that includes unit tests, behavioral-driven tests, and integration tests. Unit tests verify individual components in isolation, behavioral-driven tests ensure alignment with user expectations, and integration tests validate interactions between components. This layered approach helps ensure that the application functions correctly at various levels of granularity (Beck, 2003; North, 2009; Whittaker, 2009).

Additionally, integrating testing into the continuous integration (CI) pipeline is essential. Automated tests should be run frequently to catch regressions and issues early. Using CI tools like Jenkins or GitHub Actions enables continuous feedback and helps maintain code quality (Fowler & Foemmel, 2006).

Documentation of test cases, scenarios, and results is also important. Well-documented tests provide clarity on what is being tested, why it is being tested, and the expected outcomes. This documentation aids in understanding the scope of testing and facilitates collaboration among team members (Cucumber, 2021).

5.3.2 Unitary Tests

Unitary tests focus on verifying the functionality of individual components. By isolating components and testing their behavior under various conditions, these tests ensure that each part of the application performs as expected. Python's unittest framework is used for creating and managing these tests, supporting thorough validation of isolated components (Python Software Foundation, 2021).

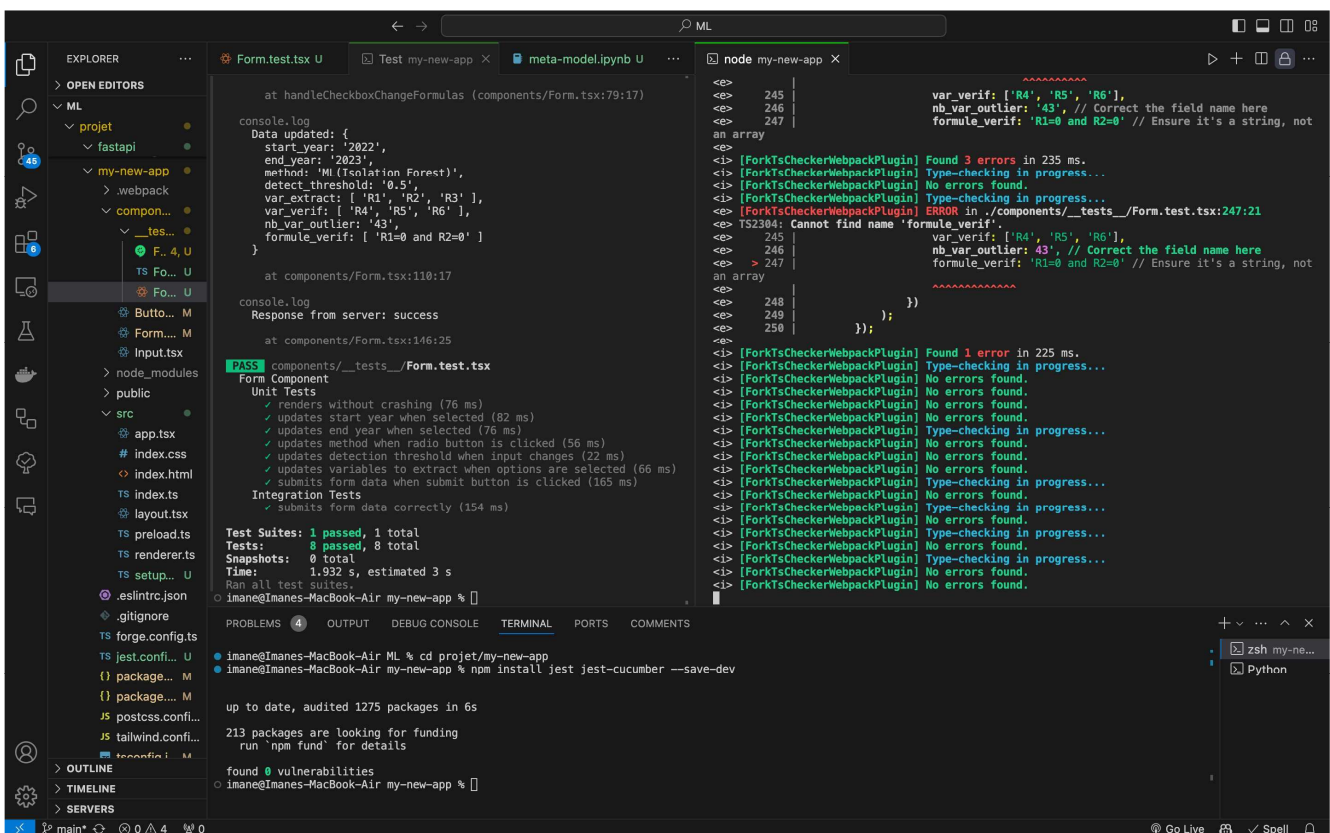
Unit tests are automated and integrated into the CI pipeline, allowing for continuous testing and prompt feedback on code changes. Mocking and stubbing techniques are employed to isolate components from their dependencies, ensuring that failures are not due to external factors (Mocking, 2021).

5.3.3 Integration Tests

Integration tests validate the interactions between various components of the application. These tests ensure that different parts of the system work together as intended. Automated integration tests using frameworks like jest and pytest help verify that the application's components integrate seamlessly and function correctly (pytest, 2021).

Mocking external systems or services is commonly used in integration testing to simulate interactions and control responses. This approach helps create controlled test environments and verify component interactions effectively (Mocking, 2021).

The following screenshot demonstrate the results of unit and integration tests conducted on the application. The unit test images provide visual evidence of successful component validations and highlight the test coverage achieved. The integration test results confirm the effective interaction between components and verify the overall functionality of the application.



```
at handleCheckboxChangeFormulas (components/Form.tsx:79:17)
console.log
  Data updated: {
    start_year: '2022',
    end_year: '2023',
    method: 'ML (Isolation Forest)',
    detect_threshold: '0.5',
    var_extract: [ 'R1', 'R2', 'R3' ],
    var_verif: [ 'R4', 'R5', 'R6' ],
    nb_var_outlier: '43',
    formule_verif: [ 'R1=0 and R2=0' ]
  }
  at components/Form.tsx:110:17
console.log
  Response from server: success
  at components/Form.tsx:146:25
PASS components/__tests__/Form.test.tsx
  Form Component
    Unit Tests
      ✓ renders without crashing (76 ms)
      ✓ updates start year when selected (82 ms)
      ✓ updates end year when selected (76 ms)
      ✓ updates method when radio button is clicked (56 ms)
      ✓ updates detection threshold when input changes (22 ms)
      ✓ updates variables to extract when options are selected (66 ms)
      ✓ submits form data when submit button is clicked (165 ms)
    Integration Tests
      ✓ submits form data correctly (154 ms)
  Test Suites: 1 passed, 1 total
  Tests: 8 passed, 8 total
  Snapshots: 0 total
  Time: 1.932 s, estimated 3 s
  Ran all test suites.
  ○ imane@Imanes-MacBook-Air my-new-app % []
  PROBLEMS 4 OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
  ● imane@Imanes-MacBook-Air ML % cd projet/my-new-app
  ● imane@Imanes-MacBook-Air my-new-app % npm install jest jest-cucumber --save-dev
  up to date, audited 1275 packages in 6s
  213 packages are looking for funding
  run 'npm fund' for details
  found 0 vulnerabilities
  ○ imane@Imanes-MacBook-Air my-new-app % []
```

Figure 5.8 – App passing unitary and integration tests

5.3.4 Behavioral Driven Tests

Behavioral Driven Development (BDD) focuses on defining application behavior from an end-user perspective. Scenarios are written in natural language to describe expected behaviors, which are then tested using frameworks like Behave (Behave, 2021). This approach ensures that the application meets user expectations and aligns with business requirements.

BDD enhances communication between developers, testers, and stakeholders by providing a clear, shared understanding of application behavior. Tests are integrated into the CI pipeline, providing regular feedback on how well the application meets specified behaviors (North, 2009).

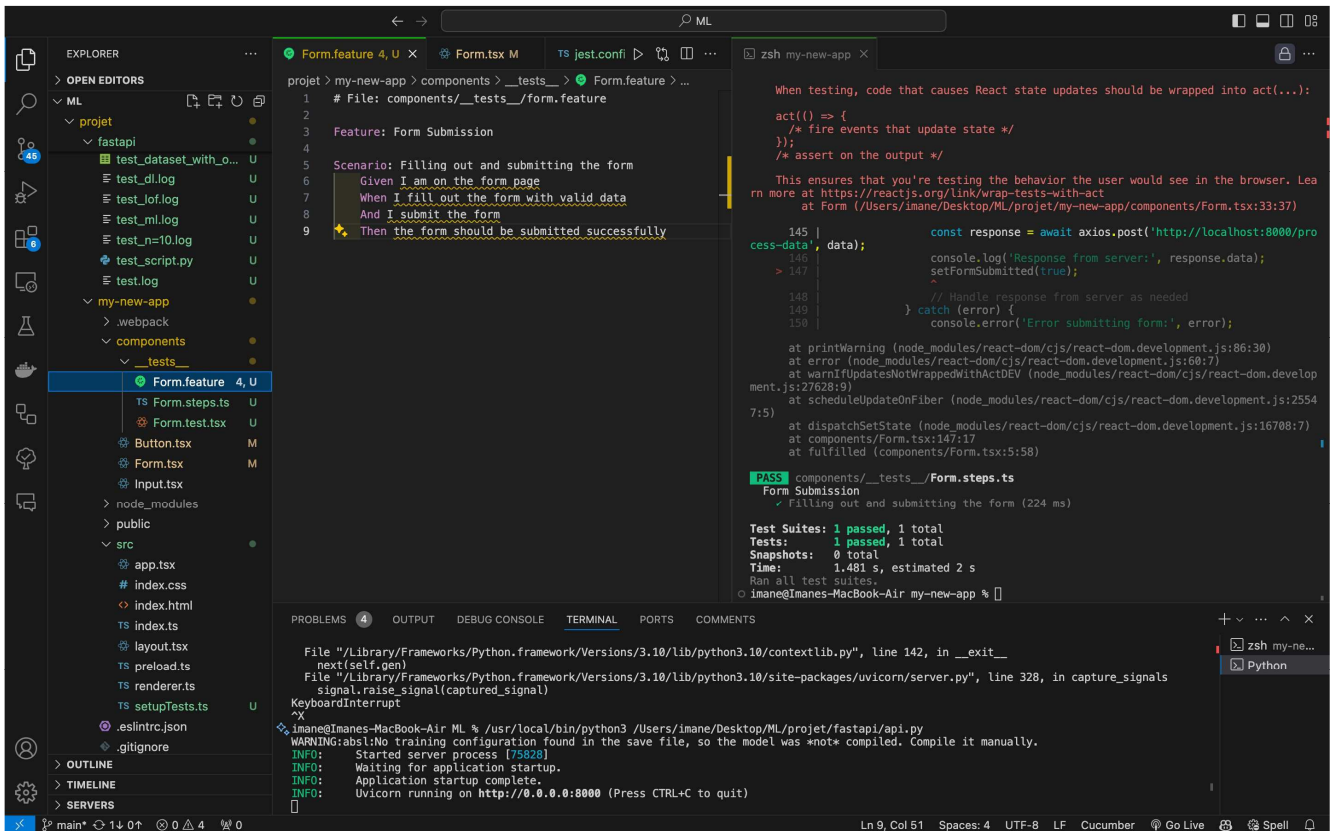


Figure 5.9 – App passing Behavioral Driven tests

These screenshots depict the outcomes of behavioral-driven tests, showcasing how well the application aligns with user requirements through specific scenarios.

5.4 Documentation

Effective documentation is a cornerstone of any successful software project, providing a comprehensive reference for understanding, using, and maintaining the application. For the desktop application developed, documentation serves multiple purposes, ensuring clarity for end-users, guidance for developers, and support for future maintenance.

5.4.1 Documentation Overview

The documentation for the desktop application is structured to cover various aspects:

- **User Documentation:** This section is designed to assist end-users in navigating and utilizing the application effectively. It includes:
 - Installation Guide: Detailed steps on how to install the application on different operating systems, including prerequisites, installation commands, and troubleshooting tips.
 - User Manual: Comprehensive instructions on how to use the application's features, including screenshots and examples to illustrate key functionalities.
 - FAQs and Troubleshooting: A compilation of frequently asked questions and common issues, along with their solutions, to help users resolve problems independently.

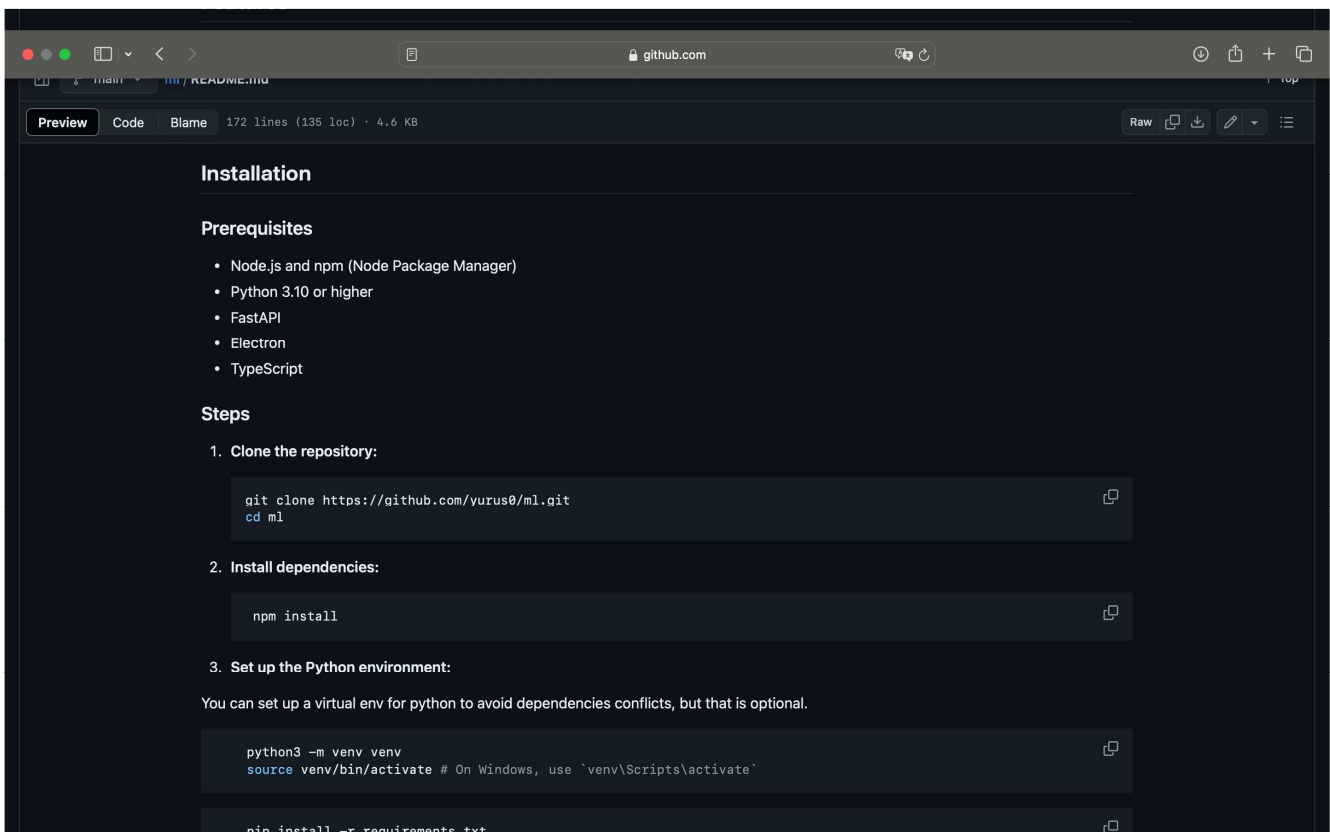


Figure 5.10 – User documentation detailing installation and usage instructions.

- **Developer Documentation:** Aimed at current and future developers, this section provides insights into the application’s architecture and codebase:
 - **Architecture Overview:** Describes the overall system architecture, including the interaction between the frontend and backend, and the integration of anomaly detection algorithms.
 - **API Documentation:** Detailed descriptions of the API endpoints provided by the backend, including request and response formats, and example usage.
 - **Codebase Details:** In-depth information about the codebase, including key modules, classes, and functions, with inline comments and explanations to facilitate understanding and modification.
- **Testing Documentation:** This component outlines the testing strategies employed to ensure the application’s reliability and performance:
 - **Test Cases:** A detailed list of test cases used during unit testing, integration testing, and behavioral testing, including the purpose of each test and expected outcomes.
 - **Testing Procedures:** Steps followed to execute tests, including any specific tools or frameworks used.
 - **Results and Logs:** Documentation of test results, including pass/fail status, error logs, and any issues encountered during testing.

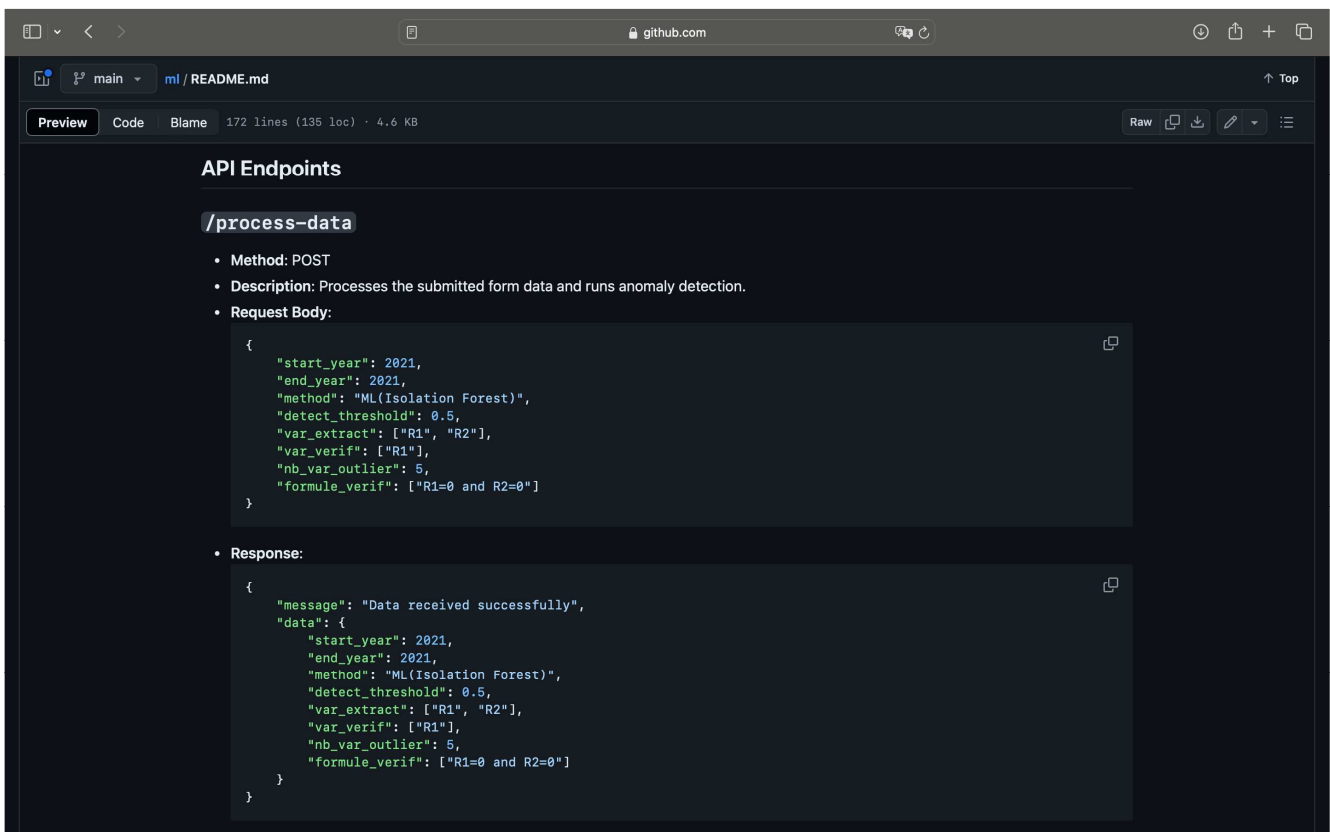


Figure 5.11 – Developer documentation showcasing application architecture and codebase details.

5.4.2 Documentation Best-Practices

Documentation is crucial for ensuring users and developers can understand, use, and maintain the desktop application effectively. To achieve this, several best practices should be adhered to.

Firstly, documentation must prioritize clarity and precision. This involves using straightforward language and minimizing technical jargon. Instructions should be presented step-by-step, guiding users through processes such as installation and configuration. Clear and detailed installation instructions, for instance, should outline each step and what users should expect at each stage to ensure a smooth setup experience. This approach is supported by Microsoft’s documentation best practices, which emphasize clarity in technical writing (Microsoft, 2021).

Visual aids, including high-resolution screenshots, diagrams, and flowcharts, are essential for explaining complex concepts. Screenshots should capture relevant parts of the user interface, while diagrams and flowcharts can represent processes or system architectures, making abstract concepts more tangible. According to the International Organization for Standardization (ISO), visual elements should complement textual explanations to provide a comprehensive understanding of the content (ISO, 2020).

Consistency in style and formatting is also critical for maintaining a professional appearance and ease of navigation. Uniform headings, text formatting, and terminology throughout the documentation help users quickly locate the information they need. For example, if “dashboard” is used to describe a component in one section, it should be consistently referred to throughout to avoid confusion. The importance of consistency is highlighted by the Institute of Electrical and Electronics Engineers (IEEE), which notes that standardized formats enhance usability (IEEE, 2019).

Comprehensiveness ensures that all relevant aspects of the application are covered. This includes not only user guides and installation instructions but also detailed descriptions of the application's architecture, APIs, and testing procedures. Comprehensive documentation allows users to fully understand the application's capabilities and usage. The Documentation Research Project emphasizes that thorough documentation is vital for effective software deployment and maintenance (Documentation Research Project, 2022).

Accessibility is another crucial aspect of effective documentation. It should be available in various formats, such as online help, downloadable PDFs, and in-app assistance, to accommodate different user needs. Documentation should also be accessible across different devices, including smartphones, tablets, and desktops. Incorporating features like text-to-speech and compatibility with screen readers ensures that documentation is inclusive for users with disabilities. The Web Content Accessibility Guidelines (WCAG) provide standards for making digital content accessible to a wider audience (WCAG, 2021).

Regular updates and reviews are necessary to keep documentation relevant and accurate. This involves revising content to reflect new features, changes in the application, and user feedback. Periodic reviews help identify and correct outdated or inaccurate information. The Agile Alliance underscores the importance of continuous feedback and iterative updates for maintaining high-quality documentation (Agile Alliance, 2020).

Finally, integrating feedback from users and developers is essential for improving documentation. Collecting input through surveys, support requests, and direct communication helps identify gaps and areas for improvement. Addressing this feedback ensures that documentation remains useful and relevant. The Software Engineering Institute supports feedback integration as a key practice for enhancing documentation quality (SEI, 2018).

By adhering to these best practices, the documentation for the desktop application will be a valuable resource for users and developers, supporting effective application use, understanding of features, and troubleshooting.

Documentation available at [GitHub Repository](#) .

5.5 Conclusion

The implementation of the desktop application has been meticulously carried out, emphasizing robust technology choices, effective user interface and backend development, and comprehensive testing and documentation. The integration of Electron and React for the frontend, combined with FastAPI for the backend, has created a modern and efficient application capable of handling complex anomaly detection tasks.

The detailed approach to testing has ensured the application's reliability, addressing potential issues through unitary, behavioral-driven, and integration tests. This rigorous validation process, aligned with best practices, confirms the application's stability and robustness. Comprehensive and clear documentation further supports users and developers by providing essential information for effective application use and maintenance.

In conclusion, the development process has demonstrated a commitment to high-quality software engineering, resulting in a well-structured and functional application ready for deployment.

As we transition to the next chapter, we will reflect on the overall impact of the project, assess the outcomes against initial objectives, and explore future directions for further development and

enhancements. This chapter will provide a comprehensive overview of the project's achievements and outline potential areas for future research and improvement.

Conclusions and Perspectives

As we conclude the detailed exploration of the desktop application's development and implementation, the final chapter aims to provide a holistic view of the project's outcomes and its broader implications. This last chapter synthesizes the insights gained from the previous chapters, evaluates the effectiveness of the implemented solution, and reflects on the project's alignment with its initial goals.

We will begin by summarizing the key findings and achievements of the project, including the efficacy of the implemented anomaly detection algorithms and the robustness of the application. We will assess how well the application meets the requirements outlined at the project's inception and examine the impact of the solution on the intended use case.

Following this evaluation, the chapter will delve into the perspectives on future development. This includes identifying areas for potential improvements, exploring emerging technologies and methodologies that could enhance the application's capabilities, and considering the broader implications of the project for related fields and applications.

By providing a reflective analysis and forward-looking perspective, this chapter aims to encapsulate the project's contributions and set the stage for ongoing development and innovation in the domain of anomaly detection and desktop application development.

General Conclusion

The objective of this report was to explore and implement advanced machine learning techniques to optimize outlier detection for enhancing data quality, specifically within the financial sector at **Bank Al Maghrib**. Throughout the research and implementation phases, we focused on developing and integrating robust algorithms capable of identifying anomalies in large datasets of annual financial statements of Moroccan companies.

We commenced by thoroughly examining the existing methodologies and their limitations. The identified challenges informed the selection and adaptation of cutting-edge techniques like *Isolation Forest*, *Mahalanobis Distance*, *DBSCAN*, *HDBSCAN*, and *Adversarially Learned Anomaly Detection (ALAD)* as well as state-of-the-art reinforcement methods such as *ML-consensus* achieved through *Stacking* and *Meta-learning*. These algorithms were meticulously analyzed and integrated into an ensemble model to leverage their individual strengths, thereby improving the overall anomaly detection accuracy.

The implementation phase involved extensive preprocessing and feature engineering to ensure the integrity and relevance of the input data. Our approach also included rigorous testing and validation procedures to confirm the efficacy and reliability of the model. The results demonstrated a significant improvement in detecting outliers, which contributes to the enhanced quality of the financial data, ensuring more accurate and reliable insights for decision-making.

Moreover, the deployment of the solution incorporated a user-friendly interface for analysts to interact with the model, visualize the results, and make informed decisions. The comprehensive documentation and training sessions conducted aimed to ensure a smooth transition and knowledge transfer to the in-house team, guaranteeing the sustainability of the solution.

Perspectives

While the results achieved in this project are promising, there are several avenues for future work and enhancement:

Algorithm Refinement and Hybrid Approaches

- Further refinement of the ensemble model could be explored by incorporating additional algorithms or fine-tuning the existing ones.
- Hybrid approaches combining statistical methods with deep learning techniques could potentially yield even better results.

Scalability and Performance Optimization

- As the volume of data continues to grow, optimizing the model for scalability and performance will be crucial. Techniques like distributed computing and parallel processing could be investigated.

Integration with Other Financial Systems

- Integrating the anomaly detection system with other financial monitoring and reporting tools within Bank Al Maghrib could provide a more holistic view of the financial ecosystem, improving overall data governance.

Expanding the Scope to Other Domains

- The methodologies developed in this project can be adapted and applied to other domains within the bank, such as transaction monitoring, fraud detection, and risk management.

Continuous Learning and Adaptation

- Implementing mechanisms for the model to learn and adapt from new data continuously will ensure its relevance and accuracy over time. This includes setting up feedback loops where human analysts can validate and correct the model's predictions.

Limitations

Our study on optimizing outlier detection using advanced machine learning techniques encountered several key limitations:

- **Data Quality and Diversity:** Challenges in handling outliers not well-represented in the dataset highlight the need for improved data collection strategies and preprocessing techniques.
- **Algorithmic Complexity:** The black-box nature of some algorithms complicates interpretability, necessitating the development of hybrid models for better transparency.
- **Scalability and Performance:** Optimizing models for large-scale deployment and real-time processing remains a challenge, requiring exploration of distributed computing and parallel processing.
- **Integration Challenges:** Seamless integration into existing financial systems and adaptation to regulatory changes require modular solutions and robust governance frameworks.
- **Ethical and Legal Considerations:** Ensuring compliance with data protection laws, addressing bias in algorithmic outputs, and maintaining transparency in decision-making are critical concerns.

Addressing these limitations will enhance the robustness and applicability of anomaly detection systems in real-world financial environments.

Final Remarks

This report has demonstrated the potential of machine learning techniques in addressing the critical issue of data quality through effective outlier detection. By harnessing advanced algorithms such as GANs and ALAD, applying theoretical algebra in statistics, and ensuring their practical applicability within Bank Al Maghrib, we have made significant strides in resolving real-world problems related to data integrity. The insights and results from this research not only contribute to the existing body of knowledge but also provide a practical framework for ongoing innovation and improvement in data quality management.

Throughout this project, we have gained invaluable insights into finance and financial markets, which have enhanced my understanding of the context in which these data quality solutions are applied. Learning about new types of neural networks, such as Generative Adversarial Networks (GANs) and Adversarially Learned Anomaly Detection (ALAD), has expanded my technical expertise and opened new avenues for addressing complex data challenges.

The application of algebraic distances has revealed their practical utility in resolving real-world statistical problems, further enriching my analytical skills. Additionally, the short duration of the project plan and the associated workload taught me the importance of time efficiency. I have honed my ability to deliver quality work within tight deadlines, communicate effectively with mentors and colleagues, and present my achievements clearly to the team, all while adhering to a scrum-based approach.

As we look to the future, it is imperative to remain committed to continuous learning, collaboration, and adaptation. By doing so, we can ensure that the solutions we develop today remain relevant and effective in meeting the evolving challenges of tomorrow.

Bibliography

- [1] Github - <https://github.com>
- [2] StackOverflow – <https://stackoverflow.com>
- [3] Git Documentation – <https://git-scm.com/docs>
- [4] Heroicons TailwindCSS – <https://pipedream.com>
- [5] Electron Documentation – <https://www.electronjs.org/docs>
- [6] FastAPI Documentation – <https://fastapi.tiangolo.com>
- [7] Uvicorn Documentation – <https://www.uvicorn.org>
- [8] Postman API Platform – <https://www.postman.com/>
- [9] F. Galton, *Regression towards mediocrity in hereditary stature*, Journal of the Anthropological Institute of Great Britain and Ireland, vol. 15, pp. 246-263, 1886.
- [10] J. W. Tukey, *Exploratory Data Analysis*, Addison-Wesley, 1977.
- [11] F. E. Grubbs, *Procedures for detecting outlying observations in samples*, Technometrics, vol. 11, no. 1, pp. 1-21, 1969.
- [12] V. Barnett and T. Lewis, *Outliers in Statistical Data*, John Wiley & Sons, 1994.
- [13] P. C. Mahalanobis, *On the generalized distance in statistics*, Proceedings of the National Institute of Sciences of India, vol. 2, no. 1, pp. 49-55, 1936.
- [14] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, *LOF: Identifying density-based local outliers*, Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data, pp. 93-104, 2000.
- [15] F. T. Liu, K. M. Ting, and Z.-H. Zhou, *Isolation Forest*, Proceedings of the 2008 IEEE International Conference on Data Mining, pp. 413-422, 2008.
- [16] H. Zenati, M. Choi, C.-J. Hsieh, and J. Yoon, *Efficient GAN-based Outlier Detection*, Proceedings of the 2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1638-1647, 2018.
- [17] R. Chalapathy and S. Chawla, *Deep Learning for Anomaly Detection: A Survey*, arXiv preprint arXiv:1901.03407, 2019.
- [18] J. Wang, X. Liu, and X. Li, *A Survey on Outlier Detection with Big Data*, IEEE Access, vol. 7, pp. 26832-26849, 2019.
- [19] C. C. Aggarwal, *Outlier Analysis*, Springer, 2017.
- [20] V. J. Hodge and J. Austin, *A survey of outlier detection methodologies*, Artificial Intelligence Review, vol. 22, no. 2, pp. 85-126, 2004.
- [21] V. Chandola, A. Banerjee, and V. Kumar, *Anomaly detection: A survey*, ACM Computing Surveys (CSUR), vol. 41, no. 3, pp. 1-58, 2009.
- [22] Bank Al-Maghrib, *Annual Report*, Bank Al-Maghrib, 2021.
- [23] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2011.

- [24] C. C. Aggarwal, *Data Mining: The Textbook*, Springer, 2015.
- [25] J. Munkres, *Topological and Geometric Methods in Combinatorics*, University of North Carolina Press, 2000.
- [26] Euclid. (300 BC). *Elements*.
- [27] Munkres, J. R. (2000). *Topology and Geometry*.
- [28] Beyer, K. S., Goldstein, J., Ramakrishnan, R., Shaft, U. (1999). When is "Nearest Neighbor" Meaningful?. In *Proceedings of the Seventh International Conference on Knowledge Discovery and Data Mining (KDD)*.
- [29] Tanger, J. (2018). Evaluating the effectiveness of Euclidean distance for outlier detection in high-dimensional spaces. *Journal of Data Science and Analytics*, 4(2), 123-135.
- [30] Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering: A review. *ACM Computing Surveys (CSUR)*, 31(3), 264-323.
- [31] Mahalanobis, P. C. (1936). On the generalized distance in statistics. *Proceedings of the National Institute of Sciences of India*, 2(1), 49-55.
- [32] Dudoit, S., & Fridlyand, J. (2003). A prediction-based resampling method for selecting the number of clusters in a data set. *Journal of the American Statistical Association*, 98(461), 1054-1069.
- [33] Iglewicz, B., & Hoaglin, D. C. (1993). *How to Detect and Handle Outliers*. Wiley Series in Probability and Statistics.
- [34] Chen, J., Wang, Y., & Fang, J. (2001). Estimating the covariance matrix in high dimensions. *Journal of Multivariate Analysis*, 77(2), 278-301.
- [35] Minkowski, H. (1907). *Geometrie der Zahlen*.
- [36] Shao, Y., Wang, J., & Zhang, J. (2016). A flexible approach to outlier detection using Minkowski distance. *International Journal of Data Science and Analysis*, 2(3), 45-57.
- [37] Chebyshev, P. (1962). *On the theory of approximations*.
- [38] Tukey, J. W. (1977). *Exploratory Data Analysis*. Addison-Wesley.
- [39] Parzen, E. (1962). On estimation of a probability density function and mode. *Annals of Mathematical Statistics*, 33(3), 1065-1076.
- [40] Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD)*.
- [41] Breunig, M. M., Kriegel, H. P., Raymond, R., & Schneider, R. (2000). LOF: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*.
- [42] Schubert, E., & Rousseeuw, P. J. (2017). Faster k-medoids clustering and robust outlier detection. In *Proceedings of the 2017 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [43] Xia, Y., & Wang, J. (2015). RS-Forest: A density-based outlier detection method using random subspace models. In *Proceedings of the 2015 IEEE International Conference on Data Mining (ICDM)*.
- [44] McInnes, L., Healy, J., & Melville, J. (2017). HDBSCAN: Hierarchical density-based spatial clustering of applications with noise. *Journal of Open Source Software*, 2(5), 205.
- [45] Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3), 1-58. <https://doi.org/10.1145/1541880.1541882>
- [46] Hodge, V. J., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2), 85-126. <https://doi.org/10.1023/B:AIRE.0000045502.10941.00>

- [47] Xia, Y., Cheng, H., Xu, Y., & Hu, J. (2015). Outlier detection with deep autoencoders. *Knowledge-Based Systems*, 88, 4-15. <https://doi.org/10.1016/j.knosys.2015.07.015>
- [48] Jin, R., Yang, J., & Yang, M. (2006). Anomalous pattern detection and classification in high-dimensional data. *Journal of Machine Learning Research*, 7, 1227-1255. <http://www.jmlr.org/papers/volume7/jin06a/jin06a.pdf>
- [49] Ahmed, M., Hu, J., & Cho, J. (2016). Outlier detection with machine learning: A survey. *Computers & Security*, 55, 26-47. <https://doi.org/10.1016/j.cose.2015.09.002>
- [50] Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008). Isolation forest. *IEEE International Conference on Data Mining*, 413-422. <https://doi.org/10.1109/ICDM.2008.17>
- [51] Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504-507. <https://doi.org/10.1126/science.1127647>
- [52] Sakurada, M., & Yairi, T. (2014). Anomaly detection using autoencoders with nonlinear dimensionality reduction. *Proceedings of the 2014 IEEE International Conference on Artificial Intelligence and Statistics*, 1, 30-38. <https://arxiv.org/abs/1407.0284>
- [53] Ng, A. Y. (2011). Sparse autoencoder. *CS294A Lecture Notes*, 72-82. <http://cs294a.stanford.edu/>
- [54] Kingma, D. P., & Welling, M. (2013). Auto-Encoding Variational Bayes. *Proceedings of the International Conference on Learning Representations (ICLR)*, 1-14. <https://arxiv.org/abs/1312.6114>
- [55] Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*. <https://arxiv.org/abs/1606.05908>
- [56] Chen, J., Song, L., Wainwright, M. J., & Jordan, M. I. (2017). Learning deep generative models of graphs. *Proceedings of the 34th International Conference on Machine Learning*, 70, 1236-1245. <http://proceedings.mlr.press/v70/chen17a.html>
- [57] Zhao, Z., & Zhang, Y. (2017). A comprehensive review of generative adversarial networks. *arXiv preprint arXiv:1712.05420*. <https://arxiv.org/abs/1712.05420>
- [58] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., & Bengio, Y. (2014). Generative adversarial nets. *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS)*, 2, 2672-2680. <https://arxiv.org/abs/1406.2661>
- [59] Zenati, H., Macedo, D., & Djelouah, M. (2018). Efficient deep anomaly detection with a reconstruction objective. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 9474-9483. <https://arxiv.org/abs/1806.09500>



Leveraging Machine Learning to Enhance the Quality of Non-Financial: Corporate Balance sheet data

AHL MBAREK SALAH – HEAD OF DATA ANALYTICS – October 2024



Plan

Introduction

Problem statement

Goals

Algorithm descriptions

Overall results

Conclusion

Future Outlook



Introduction



بنك المغرب
BANK AL-MAGHRIB
BANK AL-MAGHRIB
بنك المغرب

Abstract:

This presentation will focus on a comprehensive study on optimizing outlier detection in financial data using advanced machine learning techniques.

The primary objective was to develop robust algorithms capable of identifying anomalies within the extensive dataset of annual financial statements from Moroccan company's Key methodologies explored include:

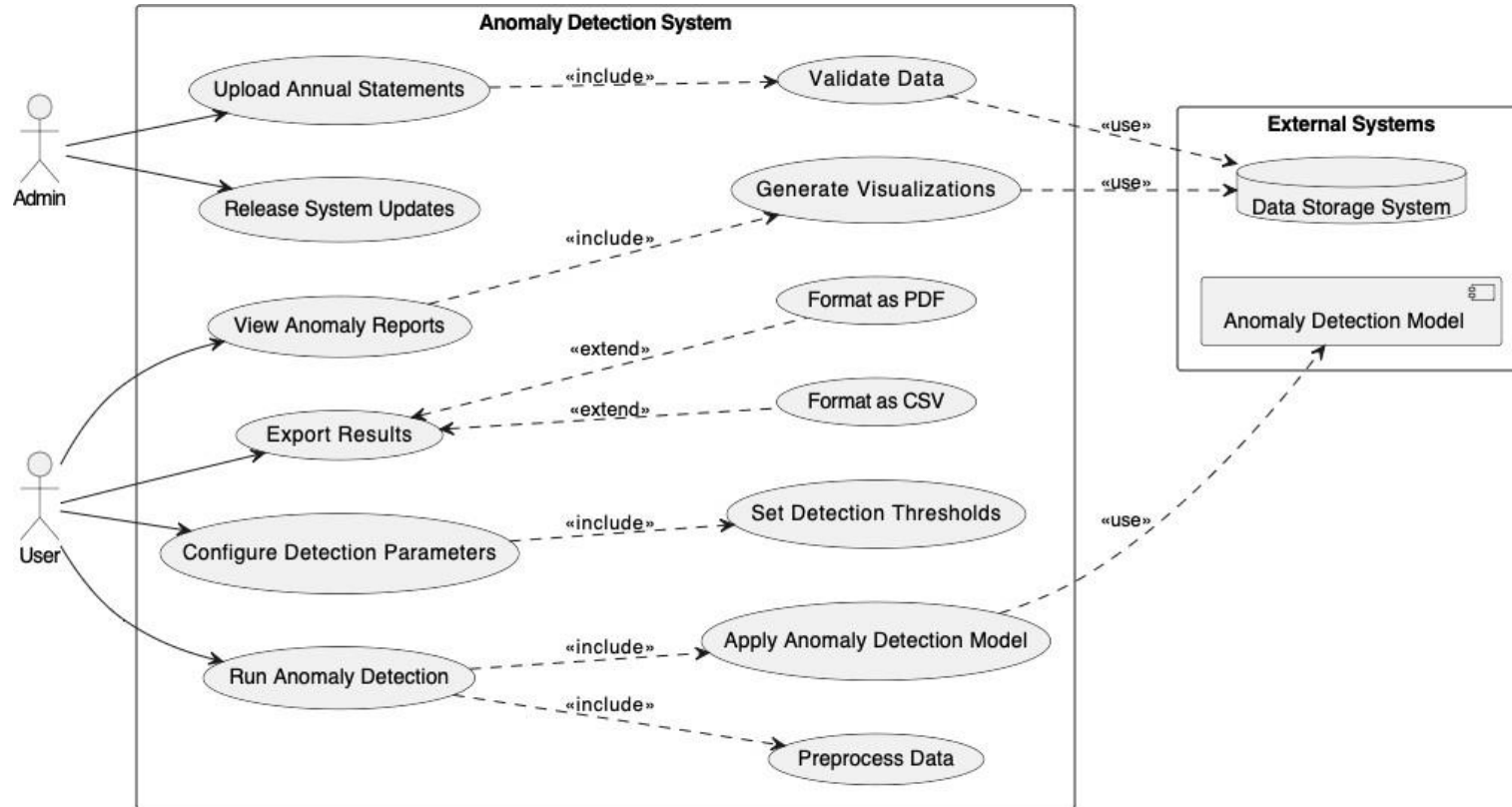
- Isolation Forest,
- Mahalanobis Distance,
- DBSCAN,
- HDBSCAN, and
- Adversarially Learned Anomaly Detection (ALAD),

Along with ensemble techniques like ML-Consensus and Meta-learning. These algorithms were carefully selected and integrated into a unified model to leverage their combined strengths. The implementation involved extensive data preprocessing and feature engineering to ensure data quality. Rigorous testing and validation demonstrated the model's effectiveness in accurately detecting outliers, leading to improved financial data quality and more reliable decision-making.

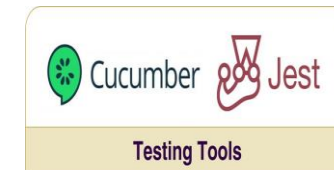
The solution also includes a user-friendly interface and comprehensive documentation to facilitate adoption and knowledge transfer within Bank Al Maghrib. Future work directions include algorithm refinement, real-time anomaly detection, scalability improvements, integration with other financial systems, and continuous learning.

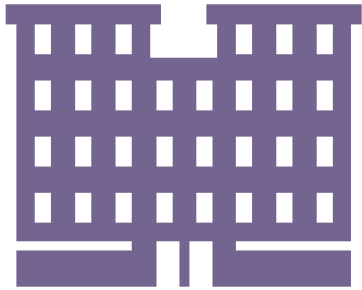
These efforts aim to maintain the model's relevance and effectiveness in addressing evolving data quality challenges. Overall, this research provides a valuable contribution to the field of anomaly detection and lays a strong foundation for future advancements in data quality management at Bank Al Maghrib.

Workflow



TOOLS





252 000

Non-financial companies 2023

~ 650 000 financial statements

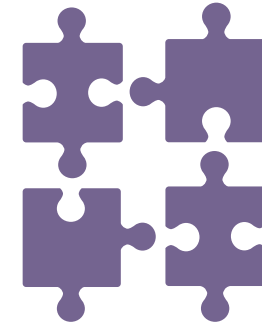


**Financial data
integrity**



**Financial data
integrity**

Decision making process



**Appropriate
governance strategy**

Problem statement



بنك المغرب
BANK AL-MAGHRIB
BANK AL-MAGHRIB
بنك المغرب



Manual processing

Time taken away from analysis



3 Weeks

Inefficient use of human resources



8 ETP

Goal

Automate repetitive data quality control tasks to streamline the process and reduce execution time without compromising data quality,



بنك المغرب
BANK AL-MAGHRIB
BANK AL-MAGHRIB
بنك المغرب

First Optimization

Automation of outlier detection using Mahalanobis distance

Processing

Near Out of Distribution

Thresholding, predictions around 0.5

Time gains



5 min

HR dedicated



0 HR & « Kafka + AIRFLOW »

Automation of outlier detection using Mahalanobis distance

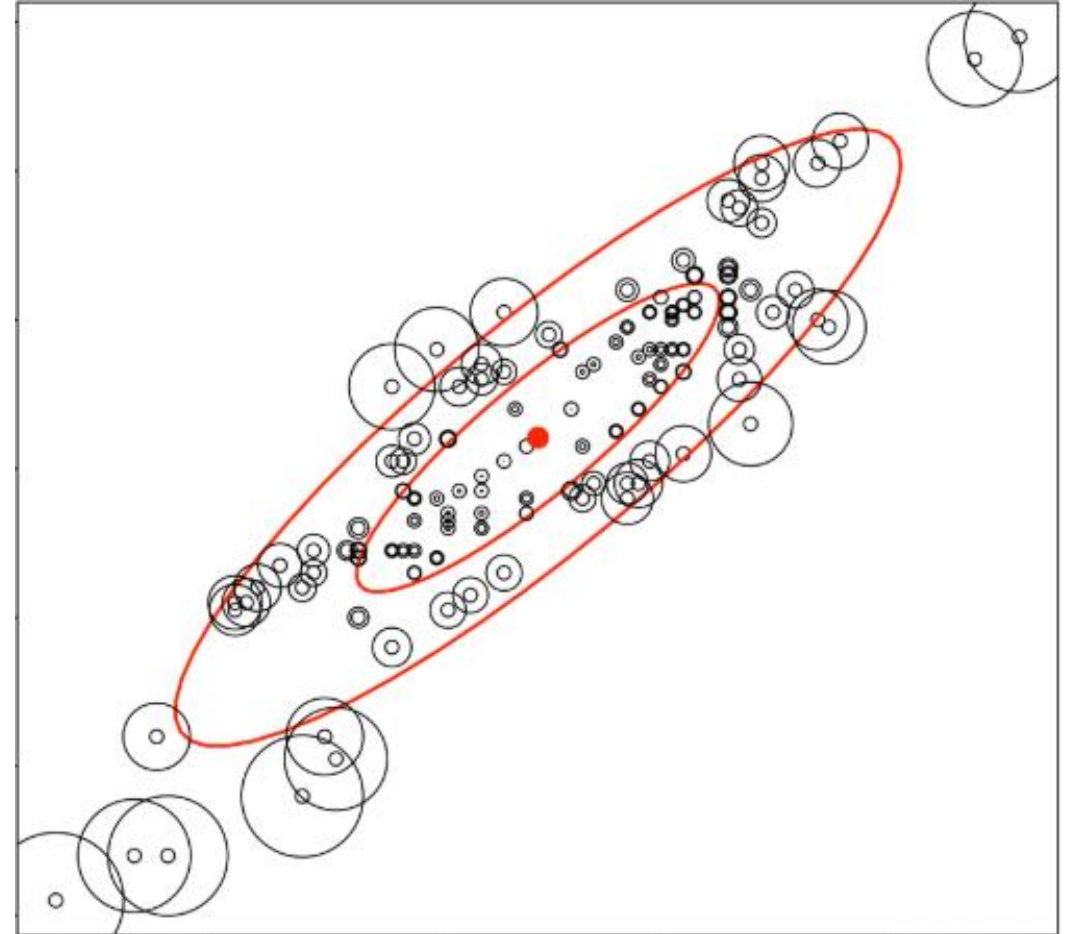
Mahalanobis Distance

Prasanta Chandra Mahalanobis, 1961

Mahalanobis distance is a statistical measure that evaluates how similar a data point is to a distribution. Unlike Euclidean distance, it considers the variability and relationships between variables, providing a more precise measure. It is especially useful for identifying anomalies in multidimensional data.

$$d_M(\mathbf{x}, \mu) = \sqrt{(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}$$

where \mathbf{x} is the observation in question, μ is the mean of the distribution, and Σ is the covariance matrix



Automation of outlier detection using ML

Isolation Forest

Fei Tony Liu, 2008.

Isolation Forest is an anomaly detection algorithm using binary trees.

It is based on the assumption that anomalies, being few and different from other data, can be isolated with a few partitions.

It has linear time complexity and low memory usage, making it efficient for large datasets. It is ideal for uniform and normal distributions. Based on isolation trees (iTrees), an isolation score is calculated using the following formula:

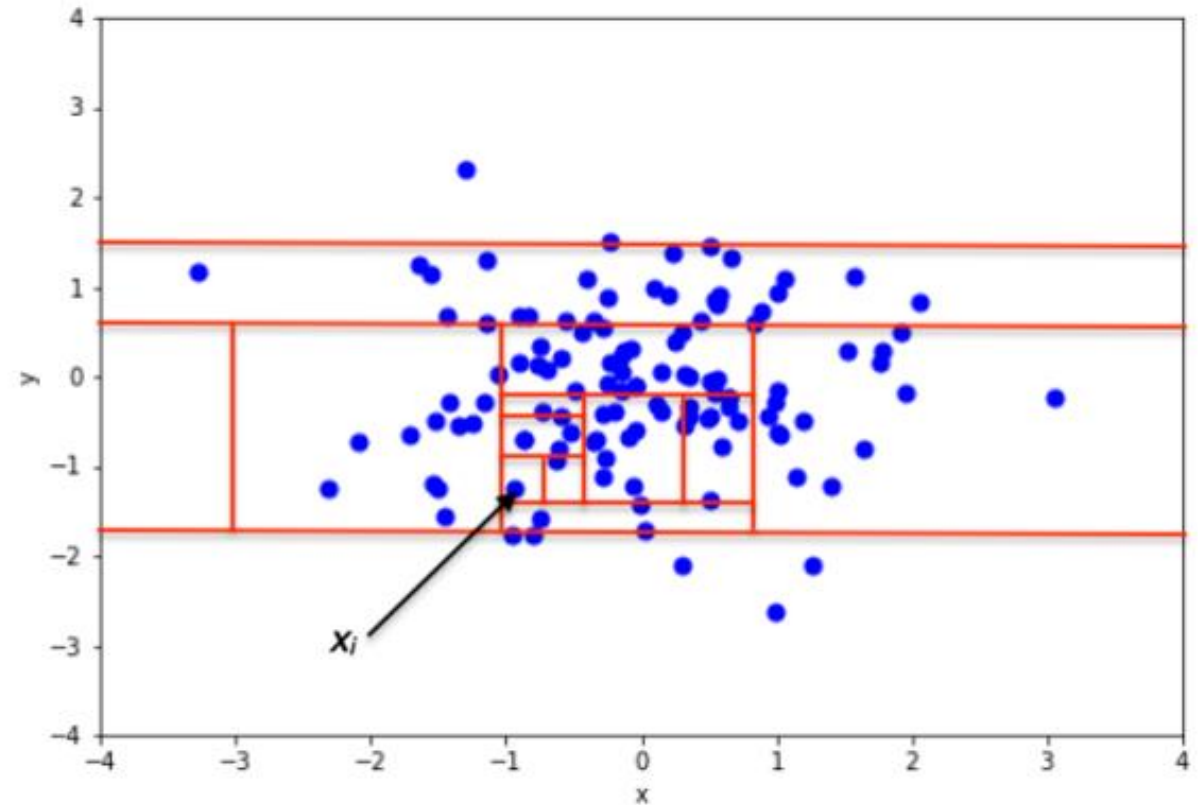
$$s(x) = \frac{2^{\frac{H(x)}{c(n)}} - 1}{c(n) - 1}$$

With:

$H(x)$ representing the average path length for point x across all trees,

$c(n)$ being the average path length for unsuccessful searches in a binary search tree defined by:

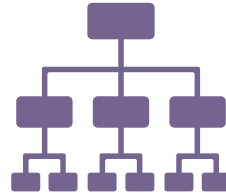
$$c(n) = 2 \cdot \left(\frac{n-1}{2} - \frac{2(n-1)}{n} \right)$$



The following image shows an example of identifying an outlier in a two-dimensional Gaussian distribution

Second Optimization

Automation of outlier detection using ML



Isolation Forest

Performance



96% Accuracy compared to manual check

Non-linear distribution and relationship

Parameter tuning, Designed for data with a uniform distribution.

Automation of outlier detection using deep learning

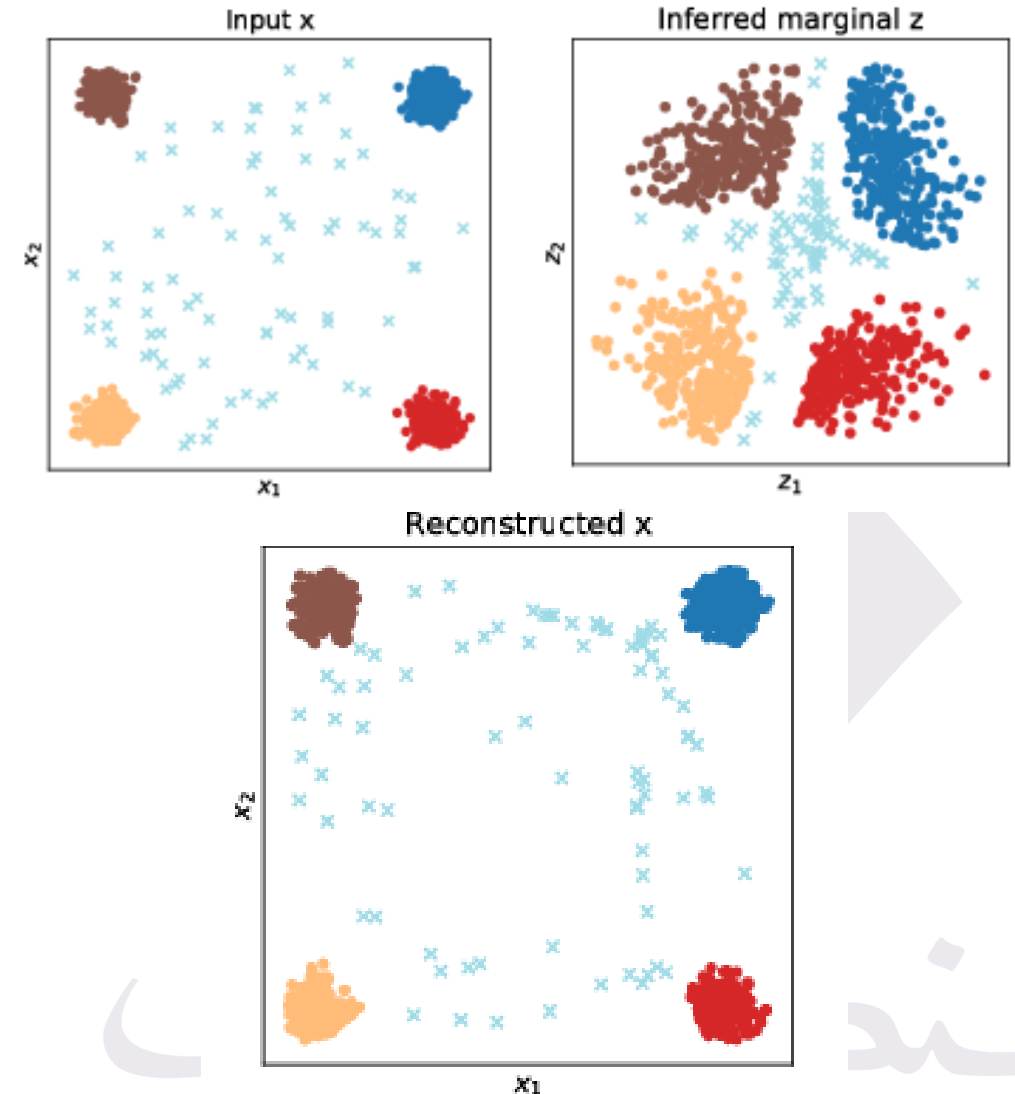
ALAD

Zenati et al., 2018.

Adversarially Learned Anomaly Detection (ALAD) is an anomaly detection method based on generative adversarial networks (GANs). The algorithm relies on two neural networks: a generator, which attempts to reproduce normal data, and a discriminator, which learns to distinguish between normal data and anomalies.

This process allows for the efficient identification of complex anomalies by implicitly modeling non-linear relationships between data. Training complexity is higher due to neural networks, but it can capture complex relationships in high-dimensional datasets. ALAD uses adversarial networks to maximize the detection of subtle anomalies, making it particularly useful for complex and multivariate data.

ALAD is a reconstruction-based anomaly detection technique that evaluates the distance between a sample and its reconstruction by the GAN. Normal samples should be reconstructed accurately, while anomalous samples will likely be poorly reconstructed.

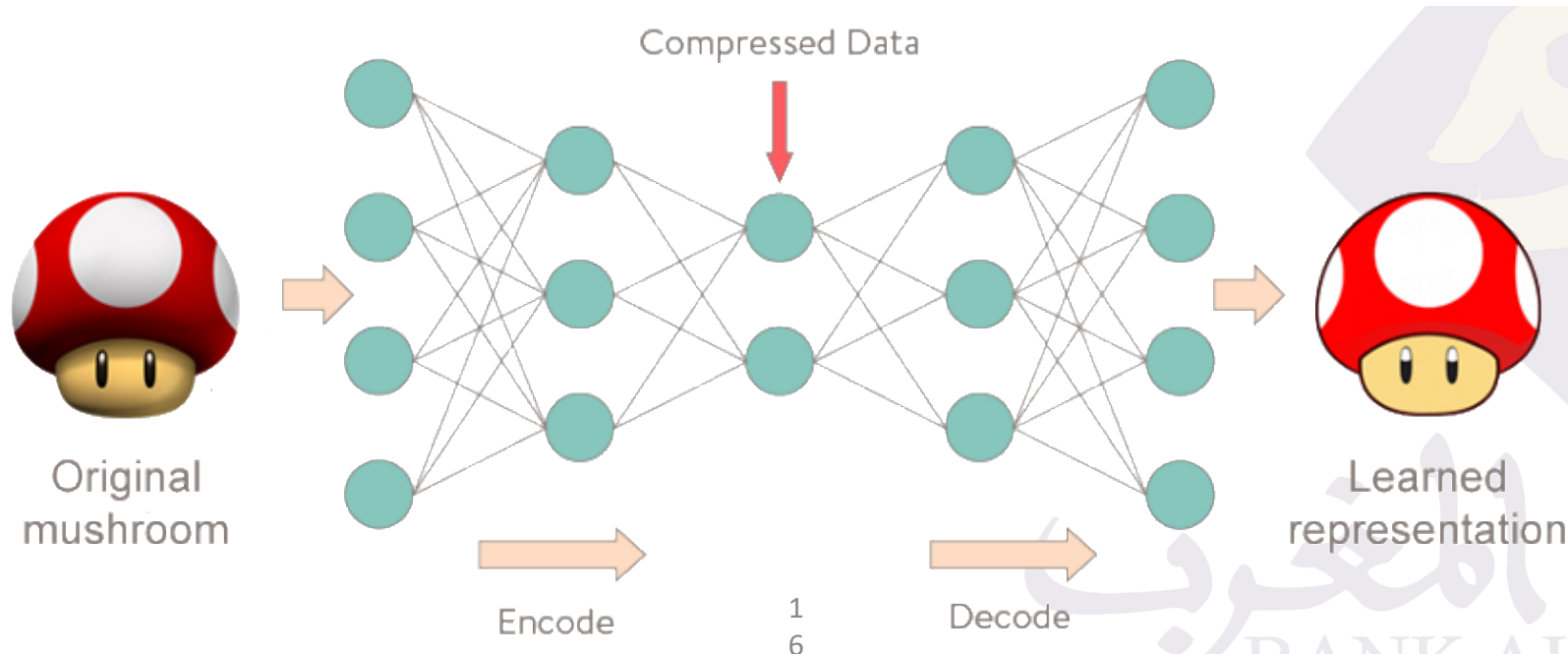


Automation of outlier detection using deep learning

Deep Auto-Encoders

Kingma et al., 2013 - Sakurada et al., 2014 - Xia et al., 2015 - Doersch, 2016.

Deep Autoencoders are a reconstruction-based anomaly detection technique. The model learns to compress data into a reduced representation (encoding) before reconstructing it into its original form (decoding). The idea is that normal samples will be well reconstructed, while anomalies, not conforming to the learned pattern, will have a poor reconstruction. Computational complexity increases with the size of the neural network, but autoencoders can capture complex non-linear relationships in high-dimensional datasets. They are particularly suited for detecting subtle anomalies in multivariate data.



4th Optimization

Automation of outlier detection using deep learning



Auto-Encoders

Performance



28 ms

Non-linear distribution and relationship

Very heavy for simpler datasets

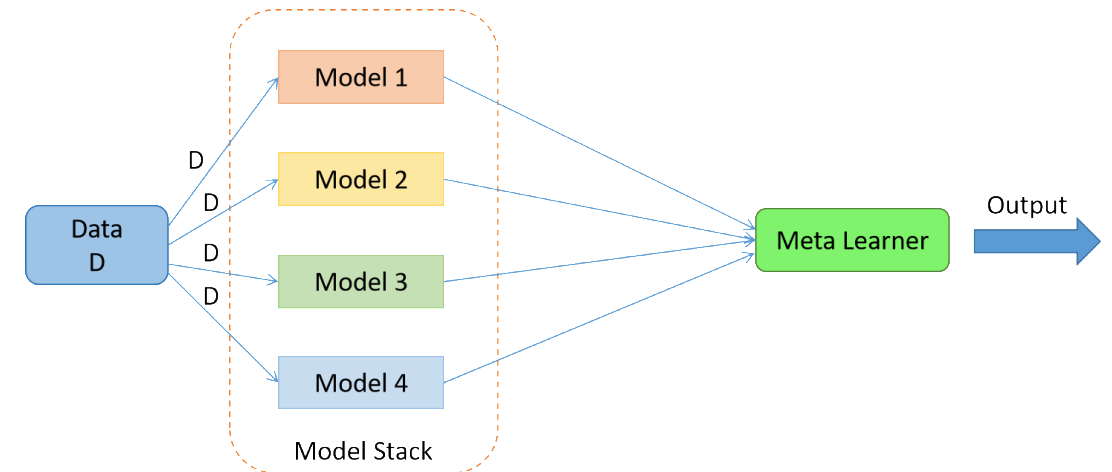
Automation of outlier detection by ML-Consensus

ML-Consensus

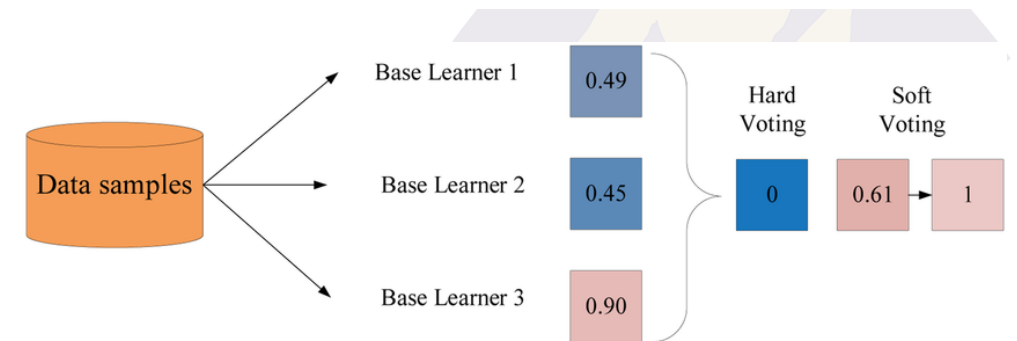
Breiman, 1996.

Consensus in ML is an approach that combines the predictions of multiple models to reach a more robust decision. The idea is to converge multiple models (often of different types) towards a common prediction, thereby reducing individual errors.

Consensus can be achieved through various methods such as averaging predictions, majority voting, or more advanced methods. This improves robustness, especially when models have varying performance on different parts of the data.



Stacking Method



0 1

The probability that the sample is positive

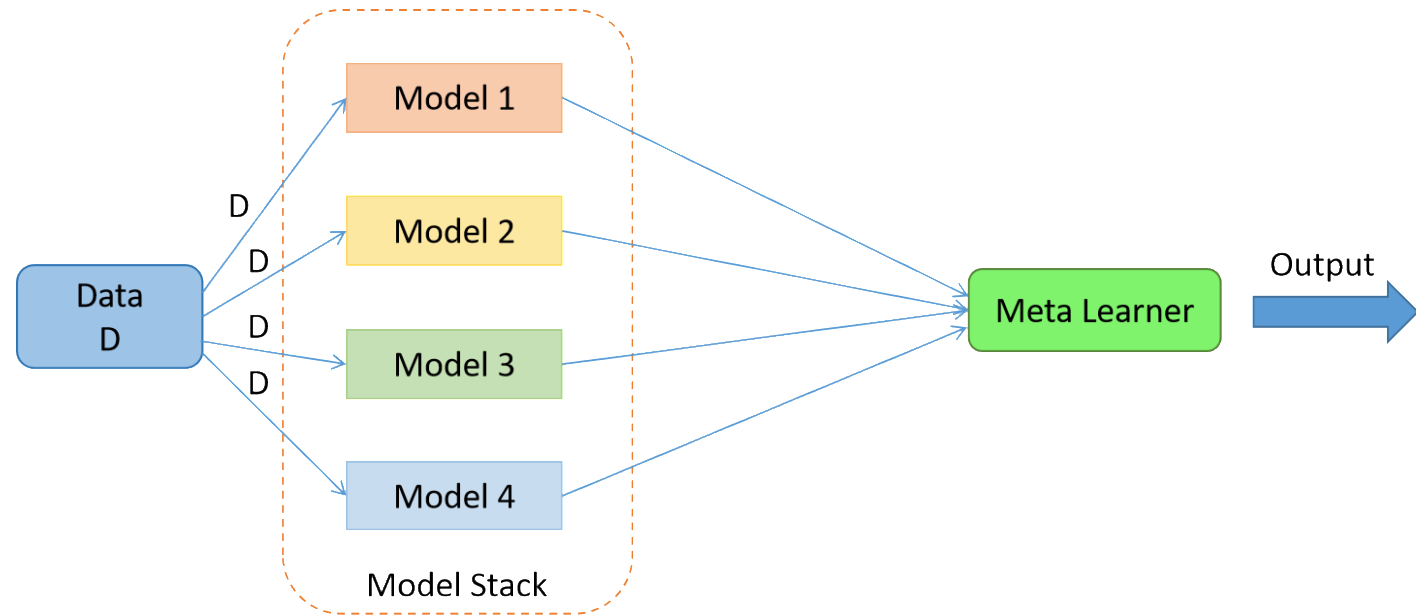
Hard vs Soft Voting Methods

Automation of outlier detection by ML-Consensus

Stacking & Meta Learning

Wolpert, 1992 - Breiman, 2001.

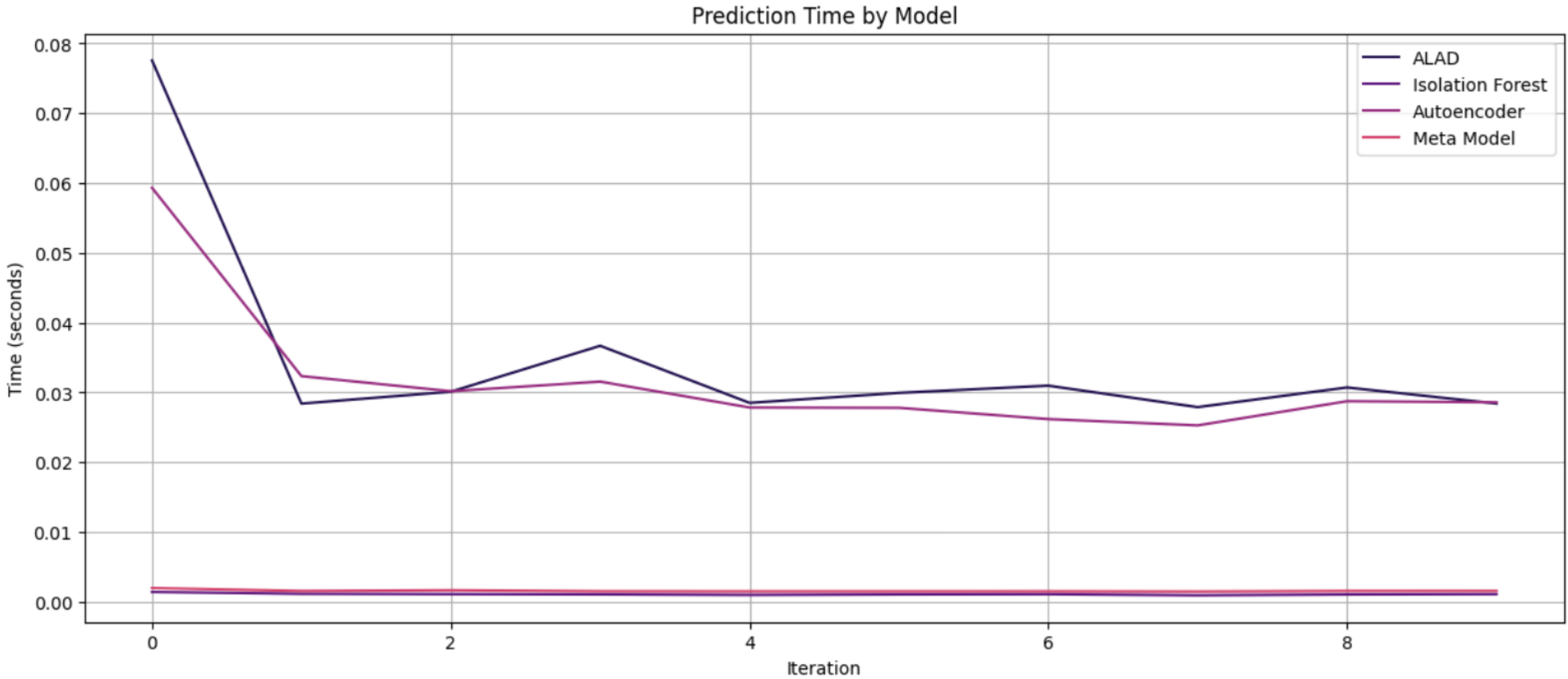
Stacking is a more sophisticated ensemble technique where multiple base models (called "level-0 models") are trained on the same data, and a meta-model (or "level-1 model") is then used to combine the predictions of these base models. The goal of stacking is to leverage the strengths of each model to improve the final prediction. Unlike methods like bagging or boosting where each model contributes equally or progressively, stacking uses an additional model to determine how to combine the outputs of the base models.



Stacking Method

Automation of outlier detection by ML-Consensus


(Stacking, Random Forest Meta Model)



Conclusion



بنك المغرب
BANK AL-MAGHRIB
BANK AL-MAGHRIB
بنك المغرب



“When the calculator was first invented, people feared that it would replace mathematicians or make people less capable, but on the contrary, it allowed us to perform more calculations and unleash greater potential.”

- Sam Altman, CEO, OpenAI

Future Outlook

- **Algorithm and hyperparameter tuning.**
- **Scalability and performance optimization (Model and Desktop App).**
- **Integration with other financial systems and/or tools.**



بنك المغرب
BANK AL-MAGHRIB
BANK AL-MAGHRIB
بنك المغرب