

**A Machine Learning Framework for Electoral Anomaly Detection: Case
Study Using Israeli Data for Kenyan Electoral Applications**

**By
Nicholas Mwadime
151147**

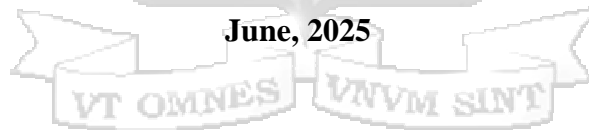
**Submitted in Partial Fulfilment of the Requirements for the Degree of Master of Science in
Data Science and Analytics at Strathmore University**

Strathmore Institute of Mathematical Sciences

Strathmore University

Nairobi, Kenya

June, 2025



This dissertation is available for Library use on the understanding that it is copyright material and that no quotation from the dissertation may be published without proper acknowledgement.

Declaration and Approval

Declaration

I declare that this work has not been previously submitted and approved for the award of a degree by this or any other University. To the best of my knowledge and belief, the dissertation contains no material previously published or written by another person except where due reference is made in the dissertation itself.

© No part of this dissertation may be reproduced without the permission of the author and Strathmore University

Student's name: Nicholas Mwadime

Sign: _____

Date: 22 - 05 - 2025

Approval

The dissertation of Nicholas Mwadime was reviewed and approved for examination by the following:

Dr. John Olukuru

Head of Data Science and Analytics.

Strathmore University.

Dr. Godfrey Madigu

Dean, Institute of Mathematical Sciences,

Strathmore University

Prof. Bernard Shibwabo,

Director of Graduate Studies,

Strathmore University

Abstract

This dissertation investigates how machine learning algorithms can be harnessed to strengthen electoral integrity by proactively detecting anomalies in voting data. Motivated by recurrent concerns about fairness and transparency in Kenya's electoral processes, the study applies a data-driven approach to uncover patterns indicative of irregularities such as mismatches between registered voters and votes cast, suspicious voting turnouts, and inconsistencies in valid and invalid ballots. To model and test the detection framework, historical electoral data from Israel (1996–2015) was used as a proxy due to its completeness and availability. The research followed the CRISP-DM methodology, encompassing phases of data understanding, preprocessing, algorithm training, and system deployment. The Isolation Forest algorithm, known for its unsupervised anomaly detection capabilities, was selected and adapted for the electoral context. The model successfully flagged 9,856 data points as anomalous across various election cycles, validating its applicability. To enhance usability, the algorithm was integrated into a Streamlit web application designed for interactive analysis, visualization, and stakeholder engagement. Through this deployment, electoral practitioners can upload datasets, visualize irregularities, and download reports in real time. The study contributes to the growing body of research on AI in public governance by presenting a practical, replicable model for anomaly detection in elections. It also proposes governance policy recommendations for adopting such tools in the Kenyan context, with the ultimate goal of fostering fairer, data-informed electoral oversight.

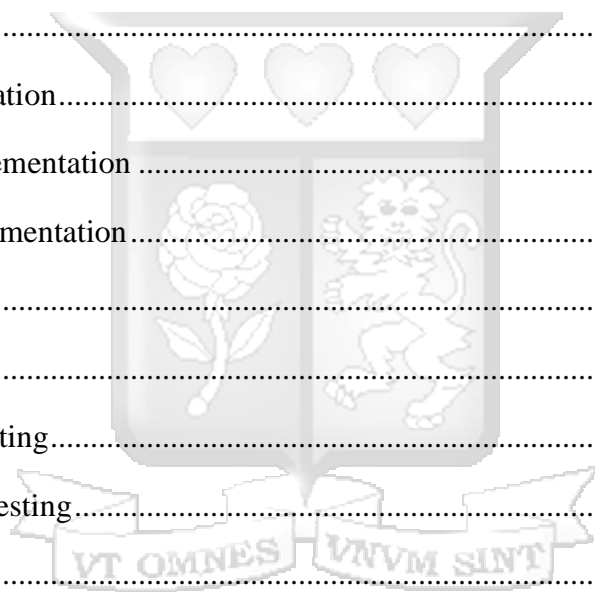
Table of Contents

Declaration and Approval	ii
Abstract	iii
Table of Contents	iv
List of Figures	ix
List of Tables	x
List of Abbreviations	xi
Acknowledgements	xii
Dedication	xiii
Chapter 1: Introduction	1
1.1 Background of the study	1
1.1.1 Technological Transformation in Democratic Discourse	1
1.1.2 Impact of AI on Political Dynamics in Western Nations	2
1.1.3 Impact of AI on Political Dynamics in Africa	2
1.1.4 Anomaly Detection in Electoral Processes	3
1.1.5 The Aim of the Study	4
1.2 Problem statement	4
1.3 Objectives	5
1.3.1 Main Objective	5
1.3.2 Specific Objectives	5
1.4 Research Questions	5
1.5 Scope of the study	6
1.6 Justification of the Study	7
Chapter 2: Literature Review	8

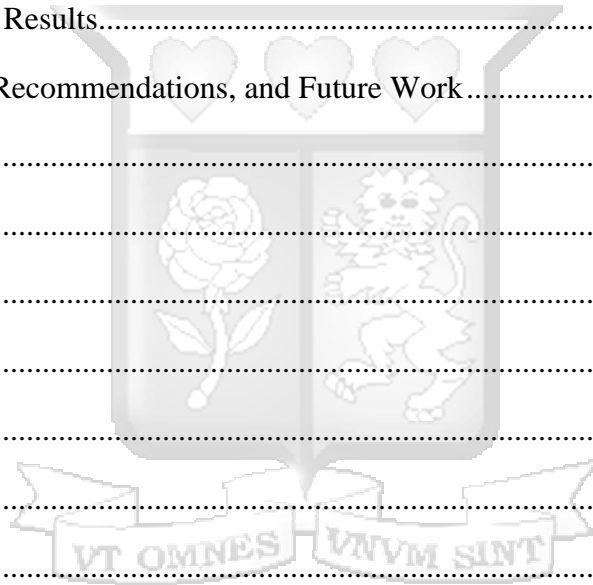
2.1 Introduction	8
2.2 Theoretical literature	8
2.2.1 Anomaly detection.....	9
2.2.2 Clustering.....	12
2.2.4 Neural networks.....	12
2.2.5 Natural Language Processing (NLP)	13
2.2.6 Ensemble Methods	14
2.3 Empirical literature.....	16
2.3.1 Types of Election Data	16
2.3.2 Challenges and Solutions in Anomaly Detection for Electoral Integrity	17
2.3.3 Existing ML Algorithms Used in Elections, and Their Strengths and Weaknesses.....	17
2.3.4 Assessing the Performance of Existing Machine Learning Algorithms.....	18
2.3.5 Proposal of Governance Policies for Anomaly Mitigation	18
2.4 Research Gap.....	19
2.5 Conceptual Framework	19
Chapter 3: Methodology	22
3.1 Introduction	22
3.2 Research design.....	22
3.2.1 Business Understanding	22
3.2.2 Data Understanding	22
3.2.3 Data Preparation:	22
3.2.4 Exploration Target Identification	23
3.2.5 Modeling.....	23
3.2.6 Evaluation.....	23
3.2.7 Deployment	24

3.3 Data loading and Inspection	25
3.4 Data Preprocessing.....	25
3.5 Feature Engineering and Scaling.....	26
3.5.1 Exploration Target Identification	26
3.6 Data Visualization.....	27
3.6.1 Line Plots	27
3.6.2 Scatter Plots	27
3.6.3 Correlation Heatmap.....	27
3.7 Comparative Model selection and development	27
3.8 Evaluation of model performance	28
3.8.1. Percentage of Anomalies	28
3.8.2. Silhouette Score for Cluster Evaluation	29
3.8.3 Visualization of Anomalies	29
3.9 Model Deployment via Streamlit Web Application.....	29
3.10 Policy Development Methodology.....	30
Chapter 4: System Design and Architecture	31
4.1 Introduction	31
4.2 System Requirements	31
4.2.1 Model API Requirements	31
4.3 Overview of System Architecture	32
4.3.1 Data Layer	32
4.3.2 Model Layer	32
4.3.3 API Layer	32
4.3.4 Fronted Layer	32
4.3.5 Database.....	32

4.4 Frontend Development.....	33
4.4.1 User Interface Design.....	33
4.5 Backend Development.....	33
4.5.1 API Integration for the Machine Learning Model.....	33
4.6 Physical Architecture.....	34
4.7 Security and Data Integrity.....	34
4.8 Deployment Strategy for Anomaly Detection Model.....	34
Chapter 5: System Implementation and Testing.....	37
5.1 Introduction.....	37
5.2 System Implementation.....	37
5.2.1 Frontend Implementation.....	37
5.2.2 Backend Implementation.....	39
5.3 Testing Procedures.....	41
5.3.1 Unit Testing.....	41
5.3.2 Integration Testing.....	42
5.3.3 Performance Testing.....	42
5.4 Results of Testing.....	43
5.4.1 Frontend Test Results.....	43
5.4.2 Backend Test Results.....	43
5.4.3 System Performance.....	44
Chapter 6: Discussion and Results.....	45
6.1 Introduction.....	45
6.2 Anomaly Detection Results.....	45
6.2.1 Overview of Detected Anomalies.....	45
6.2.2 Feature Engineering and Scaling.....	46



6.2.3 Data Visualizations	48
6.3 Model Evaluation	53
6.4 Interpretation of Anomalies	55
6.4.1 Electoral Integrity Implications	55
6.4.2 Broader Impact on Electoral Research	55
6.5 System Design and Deployment	56
6.5.1 System Design	56
6.5.2 Deployment	57
6.6 Governance Policy Results.....	58
Chapter 7: Conclusions, Recommendations, and Future Work.....	60
7.1 Introduction.....	60
7.2 Conclusions.....	60
7.2.1 Key Findings.....	60
7.3 Limitations	61
7.4 Recommendations	62
7.5 Future Work	63
References.....	64
Appendices.....	72
Appendix I: Similarity Report.....	72
Appendix II: Ethical Clearance Release Letter	74
Appendix III: Election Streamlit Application	75



List of Figures

Figure 2. 1: Conceptual Framework	21
Figure 4: 1 Schematic diagram illustrating the system architecture for your electoral anomaly detection system.....	36
Figure 5. 1: Upload mechanisms using csv format.....	38
Figure 5. 2: t-SNE visualizations showing well-separated clusters of normal data and anomalies	39
Figure 5. 3: Backend training of the model before deployment	40
Figure 5. 4: Model training using isolation forest.....	41
Figure 5. 5: Model accuracy using the silhouette score	44
Figure 6 1: The number of anomalies detected using isolation forest model	46
Figure 6 2: Feature selection.....	47
Figure 6 3: The implementation of feature scaling	48
Figure 6 4: Trends of registered voters, votes, invalid votes and valid votes over the years	49
Figure 6 5: the relationship between registered voters and votes, invalid votes and valid votes .	50
Figure 6 6: Distribution of voter anomalies.....	51
Figure 6 7: Relationships between numerical features such as votes, registered voters, valid votes, and invalid votes.....	52
Figure 6 8: Percentage of anomalies.....	54
Figure 6 9: Silhouette Score.....	55

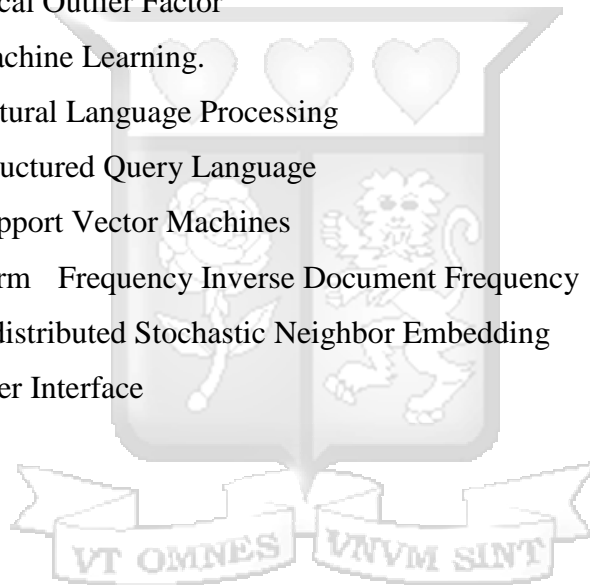
List of Tables

Table 6. 1: Comparative Performance of Anomaly Detection Algorithms on Israeli Electoral Dataset (1996–2015).....	53
---	----



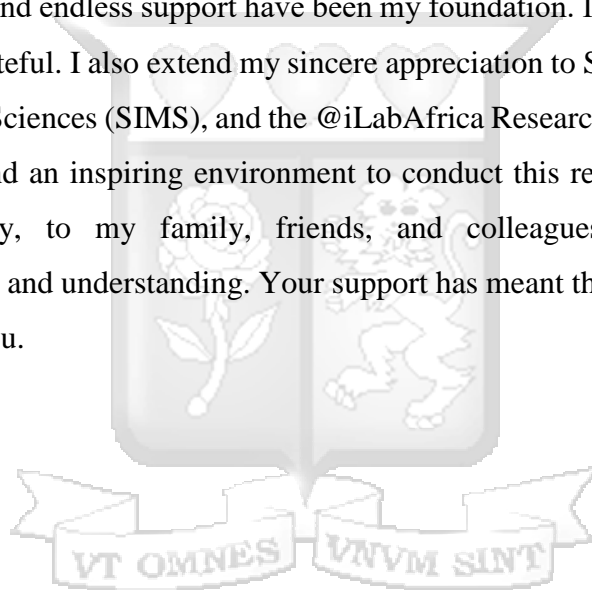
List of Abbreviations

2D	2- Dimensional
AI	Artificial Intelligence.
API	Application Programming Interface
CRISP-DM	Cross Industry Standard Process for Data mining
CSV	Comma Separated Values
EDA	Exploratory Data Analysis
HTTPS	Hyper-Text Transfer Protocol Secure
LOF	Local Outlier Factor
ML	Machine Learning.
NLP	Natural Language Processing
SQL	Structured Query Language
SVM	Support Vector Machines
TF-IDF	Term Frequency Inverse Document Frequency
t- SNE	t- distributed Stochastic Neighbor Embedding
UI	User Interface



Acknowledgements

First and foremost, I thank God for giving me the strength, wisdom, and perseverance to see this research through. I am deeply grateful to my supervisor, Dr. John Olukuru, for his invaluable guidance, patience, and encouragement. Your insights and support have been instrumental in shaping this work, and I truly appreciate your dedication. To my beloved wife, Stella Njue, your unwavering love, patience, and constant belief in me have been my greatest source of strength. Thank you for standing by me through this journey. To my loving mother, Proscovia Mwadime, your sacrifices, prayers, and endless support have been my foundation. I wouldn't be here without you, and I am forever grateful. I also extend my sincere appreciation to Strathmore University, the School of Mathematical Sciences (SIMS), and the @iLabAfrica Research Centre for providing the resources, knowledge, and an inspiring environment to conduct this research. Your support has been invaluable. Finally, to my family, friends, and colleagues—thank you for your encouragement, patience, and understanding. Your support has meant the world to me, and I truly appreciate each one of you.



Dedication

I dedicate this work to my beloved wife, Stella Njue—your love, patience, and unwavering support have been my greatest source of strength. To my loving mother, Proscovia Mwadime—your sacrifices, prayers, and endless encouragement have shaped my journey, and I am forever grateful. I also dedicate this research to everyone working to uphold electoral integrity and ensure a fair and just democracy for future generations.



Chapter 1: Introduction

1.1 Background of the study

1.1.1 Technological Transformation in Democratic Discourse

Elections are fundamental to democratic governance, serving as a direct expression of the public's will. However, the integrity of electoral processes has increasingly come under scrutiny due to the potential for irregularities, fraud, and manipulation. In response, researchers and policymakers are exploring advanced technological solutions to safeguard election credibility. Among these, Machine Learning (ML) and Artificial Intelligence (AI) have emerged as powerful tools in identifying and mitigating electoral anomalies (Chandola et al., 2009; Ahmed et al., 2016). Anomaly detection, a branch of ML, focuses on identifying rare patterns or outliers that deviate significantly from the majority of data. This approach has been successfully employed in various domains such as fraud detection (Bolton & Hand, 2002), cybersecurity (Pacha & Park, 2007), and healthcare diagnostics (Pimentel et al., 2014). Notably, models like the Isolation Forest (Liu et al., 2008) have demonstrated exceptional performance in detecting outliers within large and complex datasets. However, their application in electoral systems, where data complexity and subtle irregularities pose significant challenges, remains underexplored. In the electoral context, irregular voting patterns can signal potential fraud or administrative errors. Traditional methods for detecting such anomalies often rely on manual auditing and statistical sampling, which are labor-intensive and prone to human error (Levine & Lopez, 2018). The introduction of ML models offers a scalable, data-driven alternative capable of automatically identifying suspicious patterns in vast datasets, enhancing the transparency and reliability of elections (Binns, 2018).

Despite the promising capabilities of machine learning (ML) in anomaly detection, a research gap persists in its tailored application to electoral systems, particularly in politically sensitive contexts like Israel. While ML models have gained traction in sectors such as finance, cybersecurity, and healthcare, their adaptation to election anomaly detection remains underexplored (Chandola et al., 2009; Chalapathy & Chawla, 2019). Recent works by Xing et al. (2020) and Mushkin & Peleg (2019) have demonstrated the potential of applying ML in elections, highlighting the need to further refine these techniques for robust detection of fraudulent activities and irregularities.

This study contributes to filling this gap by examining Israeli election data and deriving lessons applicable to electoral systems in emerging democracies like Kenya.

1.1.2 Impact of AI on Political Dynamics in Western Nations

In the realm of decision-making, artificial intelligence (AI), particularly machine learning, has emerged as a formidable force, significantly impacting electoral processes worldwide. Its applications have redefined how citizens interact with political information, reshaping the dissemination, consumption, and sharing of content (Tucker, Theocharis, Roberts, & Barberá, 2017). The intricate interplay between machine learning algorithms and the vast reservoirs of data available on social media and online platforms has become a defining characteristic of contemporary political landscapes (Bennett & Livingston, 2018). Leveraging AI for content organization has enabled platforms to provide users with dynamic, personalized experiences, fundamentally altering the dynamics of modern information dissemination (Ferrara, 2020).

Globally, the integration of ML algorithms for anomaly detection in elections has gained significant momentum. In the United States, ML models have been used to analyze voter registration data to detect duplicate registrations and inconsistencies, thereby improving voter roll accuracy (Chawla & Davis, 2010). Additionally, social media platforms have employed ML techniques to monitor and limit the spread of misinformation during election periods, mitigating the risk of voter manipulation. European countries have similarly incorporated ML algorithms into their election monitoring frameworks. For example, advanced anomaly detection systems have been used to identify irregular voter turnout patterns in specific regions by combining demographic data with voting records (Aggarwal, 2015). These proactive measures have bolstered the transparency and accountability of democratic processes

1.1.3 Impact of AI on Political Dynamics in Africa

The impact of AI and machine learning technologies is not confined to Western nations; it has made inroads into the African continent, shaping political events and movements. The role of AI in the "Arab Spring" in Tunisia, where algorithms were employed to coordinate and catalyze political change, serves as a stark example. The influence of these technologies extends beyond geopolitical boundaries, challenging the traditional paradigms of political engagement and activism (Malevolent Soft Power, AI, and the Threat to Democracy, 2018).

Within the African context, the rising internet penetration has provided a platform for the confluence of AI and political processes. The 2017 elections in Kenya witnessed the use of AI and machine learning to disseminate information, both authentic and misleading, across various geographic areas and political fronts. The increasing prevalence of social media participation during the election period amplified existing conversations, with voters receiving tailored messages based on sophisticated big data and machine learning forecasts. This personalized approach, varying according to individual susceptibility to different arguments, underscored the nuanced ways in which technology influences political narratives (Polonski, 2017).

In navigating the intricate landscape of AI in politics, the examination of misinformation becomes imperative. Rebekka Rumpel, a research assistant at the Chatham House think tank, emphasizes the challenge of assessing the impact of "misinformation" on voters and their choices. During the lead-up to the August 2017 elections in Kenya, numerous websites and campaign advertisements faced scrutiny for employing scare tactics to sway votes. The common thread among these online platforms and advertisements was the integration of AI and machine learning systems, adding a layer of complexity to the dynamics of political information dissemination (Cambridge Analytica and Its Role in Kenya 2017 Elections, 2018).

1.1.4 Anomaly Detection in Electoral Processes

Anomaly detection techniques have proven invaluable in domains where accuracy and reliability are critical. Chandola et al. (2009) emphasized the effectiveness of these methods in identifying fraudulent behaviors across various industries. In electoral systems, anomaly detection can be pivotal in identifying and preventing electoral fraud, irregular vote counts, and unusual voting patterns that could undermine election outcomes. The Isolation Forest algorithm, introduced by Liu et al. (2008), is particularly effective for analyzing large electoral datasets due to its efficiency in identifying outliers without requiring labeled data. Its suitability for detecting subtle yet impactful anomalies makes it an ideal choice for analyzing voting patterns and detecting irregularities in elections.

1.1.5 The Aim of the Study

As AI continues its relentless progression in global politics, sophisticated machine learning-based voter analytics tools and big data are expected to play an enduring role in campaign management. The fusion of technology with political campaigns has become so pervasive that disconnecting from the internet, avoiding online shopping, and abstaining from social media applications seem to be the only means of evading the influence of skilled software developers surreptitiously incorporating AI-powered strategies into daily political activities (Berkowitz & Obama's, 2020). The primary aim of this study is to investigate the effectiveness of machine learning algorithms, specifically the Isolation Forest algorithm, in detecting anomalies within Israeli electoral data from 1996 to 2015. The insights gained from this analysis will serve as a basis for adapting and applying these techniques to the Kenyan electoral system to improve transparency and electoral integrity.

1.2 Problem statement

Kenya's electoral processes had been marred by recurring instances of fraud, misinformation, and irregular voting patterns, which significantly eroded public trust in democratic institutions. Despite the adoption of various technological interventions, the identification and mitigation of anomalies in election data remained limited and insufficient. Existing research primarily focused on general anomaly detection without addressing the specific complexities of Kenya's electoral system, leaving a critical gap in ensuring the credibility and transparency of elections (Cybersecurity in Elections, 2019). Irregularities such as sudden spikes in voter turnout, discrepancies in vote counts, and irregular registration patterns continued to compromise the integrity of election outcomes. These anomalies were often subtle and embedded within vast datasets, making their detection challenging with conventional methods. Additionally, existing machine learning models had not been adequately adapted to address Kenya's unique electoral dynamics, limiting their effectiveness in identifying and mitigating fraudulent activities (Assibong et al., 2020).

To address this pressing issue, the study proposed the adaptation and refinement of the Isolation Forest algorithm, a robust anomaly detection model, to analyze electoral data. The Isolation Forest algorithm had been recognized for its efficiency in detecting anomalies within large datasets by isolating observations that were few and different (Liu et al., 2008).

Israeli election data from 1996 to 2015 served as a benchmark for developing and testing the model, providing insights into how similar techniques could be applied to Kenya's elections. Furthermore, the deployment of an interactive, user-friendly web-based platform bridged the gap between theoretical research and practical application, offering election stakeholders timely and actionable insights into potential anomalies. By addressing these challenges, the study aimed to contribute to strengthening Kenya's electoral integrity, enhancing transparency, and restoring public confidence in the democratic process (Serbanescu, 2021).

1.3 Objectives

1.3.1 Main Objective

The main objective of this research is to investigate and refine machine learning algorithms for effective detection of electoral anomalies in Kenyan elections, contributing to transparent and credible electoral processes.

1.3.2 Specific Objectives

1. To identify and classify common types of election anomalies in Kenya.
2. To adapt and evaluate machine learning algorithms for detecting electoral anomalies.
3. To propose governance policies for mitigating electoral anomalies in Kenya.

1.4 Research Questions

1. What are the common types of electoral anomalies observed in Kenyan elections, and how can they be systematically identified and classified?
2. How can machine learning algorithms, particularly Isolation Forest, be adapted and evaluated to effectively detect electoral anomalies in the Kenyan context?
3. What governance policy recommendations can be developed from the findings of anomaly detection to enhance transparency and integrity in Kenya's electoral processes?

1.5 Scope of the study

The scope of this study was designed to comprehensively analyze and address the multifaceted anomalies inherent in the Kenyan electoral system. Central to this research was a detailed examination of various irregularities and discrepancies that could occur throughout the electoral process. These anomalies included discrepancies between the number of registered voters and actual votes cast, inconsistencies in valid and invalid votes, and irregularities in voter registration processes that might have compromised the integrity of election outcomes. To enrich this analysis, the study incorporated a comparative approach by utilizing Israeli electoral data from 1996 to 2015. This dataset served as a benchmark to identify and understand similar patterns of electoral irregularities, offering valuable insights into how machine learning models performed in diverse electoral environments. By comparing Kenyan and Israeli election data, the research aimed to adapt successful anomaly detection methods to the Kenyan context, enhancing the relevance and accuracy of findings. Additionally, the study explored existing machine learning algorithms, with a focus on adapting and refining them for effective anomaly detection in Kenyan elections. The use of algorithms such as the Isolation Forest was critically examined for their capacity to identify patterns, detect irregularities, and flag potentially fraudulent activities within electoral data.

Beyond the technical dimension, the study contextualized identified anomalies within Kenya's unique socio-political landscape. It considered historical precedents, cultural dynamics, and institutional frameworks that might have contributed to or mitigated electoral irregularities. This comprehensive approach enabled a nuanced understanding of the factors influencing election integrity in Kenya. The study further sought to bridge the gap between theoretical insights and practical applications by proposing evidence-based governance policies for anomaly mitigation. These recommendations were informed by empirical analysis and aimed to strengthen voter registration protocols, enhance polling procedures, and improve oversight mechanisms. The overarching goal was to foster greater transparency, fairness, and public trust in Kenyan electoral processes. By integrating a comparative analysis with Israeli electoral data and combining technical expertise with socio-political insights, this research aspired to make meaningful contributions to advancing electoral integrity and democratic governance in Kenya.

1.6 Justification of the Study

Ensuring electoral integrity remains a foundational pillar of democratic governance. Across many democracies, including Kenya, elections are often marred by challenges such as voter fraud, misinformation, and irregularities in vote counting. These challenges erode public trust, undermine political stability, and compromise the legitimacy of elected governments. With the growing complexity of electoral systems and the increasing digitization of election-related processes, there is a pressing need for more intelligent, data-driven approaches to detect and mitigate anomalies that may indicate malpractice. Machine learning (ML) has proven effective in domains such as fraud detection, finance, and cybersecurity, where massive datasets require intelligent, automated pattern recognition. Yet, its use in electoral anomaly detection—particularly in African contexts—remains underexplored. Most existing studies are concentrated in Western democracies, where political and technological conditions differ significantly from those in countries like Kenya. This study addresses this gap by exploring how ML models can be adapted and refined for electoral environments characterized by multifaceted political dynamics and constrained institutional capacities. To achieve this, the study utilizes Israeli electoral data spanning from 1996 to 2015 as a testbed for developing and refining an ML-based anomaly detection model. Israel provides a relevant comparative case for several reasons. First, both Israel and Kenya share characteristics such as vibrant multiparty democracies, periodic electoral controversies, and ethnopolitical voting patterns. Second, Israel's long-standing electronic vote recordkeeping and structured datasets offer a reliable foundation for training and validating machine learning models. These models are then evaluated and contextualized to suit the Kenyan setting, where similar anomalies—such as mismatches in registered voters and actual turnout, or invalid vote spikes—have raised concerns in past elections. Beyond anomaly detection, this study contributes by proposing governance policies grounded in evidence-based anomaly insights. These recommendations aim to strengthen electoral oversight mechanisms in Kenya and to support decision-makers, electoral commissions, and civil society actors in monitoring election integrity in real time. Moreover, the practical deployment of the model as an interactive Streamlit-based web application ensures usability by both technical and non-technical stakeholders. In summary, the study is justified not only by the novelty of applying machine learning to electoral anomaly detection in Kenya, but also by its comprehensive methodological framework, real-world deployment plan, and the broader governance implications that emerge from its findings.

Chapter 2: Literature Review

2.1 Introduction

This chapter presented a comprehensive examination of the role of machine learning (ML) algorithms in enhancing electoral processes, with a focus on anomaly detection to safeguard electoral integrity. The increasing adoption of ML technologies in elections worldwide had demonstrated promising advancements in improving the accuracy, security, and transparency of electoral systems. However, challenges such as interpretability, transparency, and algorithmic bias remained significant concerns. The literature review explored global and regional applications of ML in elections, highlighting both the opportunities and limitations inherent in their use. Special attention was given to how anomaly detection models, such as the Isolation Forest algorithm, had been employed to identify irregular voting patterns and fraudulent activities in various contexts. The review also investigated how these models could be adapted to address Kenya's unique electoral challenges, drawing comparisons from the Israeli election data (1996–2015) to inform strategies for improving Kenya's electoral integrity. By synthesizing theoretical frameworks and empirical studies, the chapter offered a well-rounded understanding of the benefits and risks associated with integrating ML into electoral systems. The aim was to provide a solid foundation for advancing reliable, transparent, and secure electoral processes through the strategic application of machine learning technologies.

2.2 Theoretical literature

The application of machine learning (ML) in electoral systems was rooted in well-established theoretical frameworks that aimed to enhance the fairness and security of elections. At the core of this study was anomaly detection, which identified irregular voting patterns that could indicate fraud or data inconsistencies. Algorithms like the Isolation Forest (Liu et al., 2008) were particularly effective in isolating unusual data points within vast electoral datasets. Clustering techniques, such as K-means and DBSCAN, offered valuable insights by grouping similar data points, making it easier to detect outliers that could reflect irregular behaviour in polling stations (Jain et al., 1999). Additionally, Natural Language Processing (NLP) had become an essential tool in analyzing political discourse and spotting misinformation campaigns that could mislead voters.

To further improve the accuracy of anomaly detection, ensemble methods like Random Forests combined the strength of multiple models to minimize prediction errors (Breiman, 2001). Meanwhile, advanced approaches like deep learning and neural networks excelled in recognizing complex, hidden patterns within large datasets, offering a sophisticated means to detect subtle electoral anomalies (LeCun et al., 2015). By weaving these theories together, the study explored how machine learning could be harnessed to safeguard electoral integrity, fostering trust and transparency in the democratic process.

2.2.1 Anomaly detection

Anomaly detection is a critical theoretical framework in data analysis, focusing on identifying patterns or behaviors that deviate from the norm. In electoral processes, this theory is pivotal in detecting irregularities such as voter fraud, ballot stuffing, or discrepancies in vote counts that threaten the integrity of elections. Election officials can use anomaly detection to flag unusual spikes in voter turnout or discrepancies between registered voters and votes cast. Key studies by Aggarwal & Sathe (2015) and Chandola et al. (2009) have established foundational approaches in this field, highlighting how machine learning (ML) algorithms like decision trees, neural networks, and clustering techniques can detect deviations in large electoral datasets. For example, ML models may identify abnormal voting spikes in certain regions or detect duplicate voter records, aiding in early detection of electoral fraud. More recently, Chalapathy & Chawla (2019) introduced deep learning models such as autoencoders, which have shown substantial improvement in anomaly detection accuracy. These models can learn intricate patterns in historical election data and detect subtle irregularities that traditional statistical models might miss. The adaptability and learning capability of ML algorithms allow them to evolve and improve over time, making them particularly effective for dynamic electoral environments. Anomaly detection relies on statistical measures to define normal behavior. One widely used method is the z-score, which measures how far a data point deviates from the dataset's mean.

The formula is given by:

$$z = \frac{X - \mu}{\sigma}$$

Where:

X = the observed data point

μ = the mean of the dataset

σ = the standard deviation of the dataset

This metric standardizes data points, flagging those that fall outside the normal range as anomalies (Montgomery et al., 2012). The referenced work by Montgomery, Peck, and Vining (2012) serves as a foundational source for understanding and applying the z-score formula in statistical analyses.

2.2.1.1 Mahalanobis distance in Anomaly Detection

The Mahalanobis distance is a pivotal statistical metric used to measure how far a data point deviates from a distribution, making it highly effective for detecting outliers in multivariate datasets (Mahalanobis, 1936). Unlike simpler distance metrics, it accounts for correlations between variables, providing a more accurate representation of anomalies in complex datasets.

The Mahalanobis distance is calculated using the formula:

$$D^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

Where:

x = the vector of the data point

μ = the mean vector of the dataset

Σ = the covariance matrix of the dataset

Σ^{-1} = the inverse of the covariance matrix

$(x - \mu)^T$ = the transpose of the difference between the data point and the mean

In the context of electoral anomaly detection, this metric is valuable for identifying irregularities across polling stations. For example, a voting precinct with an unusually high deviation in vote counts or voter turnout can be flagged for further examination. Studies by Ahmed and Ahmed (2016) applied this method to Pakistan's 2013 general elections, successfully identifying polling stations with suspicious voting patterns. Similarly, Alomari and Mohammed (2015) developed a Mahalanobis distance-based algorithm to detect fraud in electronic voting systems, demonstrating its effectiveness in identifying anomalous votes. These applications highlight how Mahalanobis distance improves election integrity by systematically identifying irregularities, enabling authorities to investigate and mitigate electoral fraud.

2.2.1.2 Isolation Forest

The Isolation Forest algorithm is a tree-based machine learning model designed for effective anomaly detection by isolating data points that deviate from normal patterns. It works by randomly selecting features and partitioning data, where outliers are identified as points that require fewer splits to isolate (Liu, Ting, & Zhou, 2008). This makes it highly efficient for large, high-dimensional datasets like election data. In electoral analysis, the Isolation Forest algorithm is particularly useful in detecting fraudulent patterns, such as unusual voting spikes or irregular voter turnout. For example, Xing et al. (2020) applied this method to Brazilian election data, uncovering precincts with voting behaviors that deviated significantly from the norm. Similarly, Botev et al. (2019) demonstrated the algorithm's capability in identifying potential voter fraud in U.S. elections, detecting anomalies like abnormally high voter turnout and disproportionate support for specific candidates. By efficiently handling large datasets without relying on predefined fraud patterns, the Isolation Forest algorithm enhances electoral integrity through proactive identification of suspicious voting trends.

2.2.1.3 Local Outlier Factor (LOF)

The Local Outlier Factor (LOF) algorithm, developed by Breunig et al. (2000), was a powerful tool for identifying anomalies by comparing how densely data points were packed relative to their neighbours. In simpler terms, if a data point was surrounded by far fewer similar points than others nearby, it raised a red flag. This made LOF especially useful in election data, where subtle irregularities—such as unexpected surges in voter turnout or inconsistent voting behaviour—could have gone unnoticed with broader detection methods. For example, during the U.S. presidential elections, Zaman et al. (2020) successfully applied the LOF algorithm to uncover polling stations that displayed suspicious voting patterns. Precincts with unusually high LOF scores stood out as potential hotspots for irregularities, warranting deeper investigation. Similarly, Fathian et al. (2020) employed LOF to analyze voter behaviour within specific precincts, effectively flagging voters whose actions deviated significantly from the norm—indicating possible fraudulent activity. The effectiveness of LOF lay in its ability to pick up on localized, nuanced discrepancies that larger-scale methods might have overlooked. In the context of elections, this precision was invaluable. By highlighting subtle anomalies, LOF supported efforts to ensure that every vote counted fairly, bolstering trust in the electoral process.

2.2.2 Clustering

Clustering was a widely used machine learning technique that organized data points into groups based on shared similarities (Tan et al., 2005). In electoral data analysis, clustering helped uncover meaningful voting patterns by grouping voters according to demographic and geographic factors, such as age, gender, income, and location. This technique allowed analysts to identify how different voter groups engaged with the electoral process. For instance, Peng et al. (2018) effectively applied clustering to study voter behaviour in the 2016 U.S. presidential election. By grouping voters based on characteristics like age, education level, and income, the researchers revealed patterns that explained how different demographic groups voted. This insight proved vital in understanding the influence of socio-economic factors on electoral outcomes.

Clustering also played a critical role in detecting voting irregularities and potential fraud. Unusual groupings in voting data such as identical voting behaviour across unrelated regions or unusually high voter turnout in specific areas could have signalled irregularities. Hsu et al. (2008) demonstrated that clustering could highlight such anomalies, prompting further investigation into possible electoral manipulation. In Kenya's electoral context, clustering provided valuable insights into voter behaviour and regional voting trends. It also helped identify unexpected voting patterns that might have indicated discrepancies in voter registration or potential election fraud. When combined with other anomaly detection methods, such as the Isolation Forest algorithm, clustering enhanced the system's ability to uncover both legitimate voting trends and suspicious activities. By offering a deeper understanding of voting behaviours and highlighting irregularities, clustering supported more transparent, secure, and credible electoral processes.

2.2.4 Neural networks

Neural networks, a form of machine learning algorithm outlined by Goodfellow et al. (2016), excelled at identifying patterns and relationships within extensive datasets, addressing complex issues. Composed of interconnected nodes and trained using input data to make predictions or classifications, these networks were inspired by the structure of the human brain and designed to simulate its problem-solving capabilities. In the realm of election data analysis, neural networks emerged as powerful tools for detecting potential instances of fraud.

By analyzing vast and intricate datasets, these models were able to uncover patterns that might have been overlooked by traditional methods. For example, in a study conducted by Xu et al. (2019), researchers developed a neural network model specifically designed to identify potential voter fraud in Brazilian elections. This model was trained on data from previous elections, allowing it to learn and recognize complex patterns associated with anomalous voting behaviours. The results demonstrated that the neural network model achieved a high level of accuracy in detecting anomalous voting patterns indicative of fraudulent activities. By leveraging their ability to process large-scale and multi-dimensional data, neural networks provided an advanced, effective solution for safeguarding electoral integrity in complex and dynamic election environments.

2.2.5 Natural Language Processing (NLP)

Neural networks were powerful machine learning models inspired by the structure and function of the human brain. They consisted of layers of interconnected nodes (neurons) that processed complex data to identify patterns and relationships. As Goodfellow et al. (2016) explained, these models were designed to learn from data, which made them highly adaptable for analyzing large and intricate datasets. In the domain of electoral data analysis, neural networks demonstrated great promise in detecting fraudulent activities. For instance, Xu et al. (2019) developed a neural network model specifically tailored to uncover voter fraud in Brazilian elections. This model was trained on historical election data and was able to effectively detect unusual voting patterns indicative of potential fraud, such as unexpected surges in voter turnout or irregular vote distributions. Neural networks excelled in recognizing subtle and complex relationships within datasets—patterns that traditional statistical methods might have overlooked. Their ability to process diverse and multidimensional information, including voter demographics, polling station data, and turnout rates, enabled them to provide a more comprehensive analysis of election integrity. In the context of Kenyan elections, neural networks could have played an instrumental role in identifying sophisticated fraud tactics, such as ballot stuffing or the tampering of voter records. Additionally, neural networks demonstrated the capability to continuously improve their accuracy over time by learning from new data. This adaptability made them particularly valuable for real-time election monitoring and fraud detection, as they could provide ongoing insights and refine their performance with additional information. The successful application of neural

networks in other electoral systems underscored their potential for promoting credible and trustworthy elections on a global scale. Their capacity to detect complex and subtle anomalies, coupled with their ability to adapt to diverse electoral contexts, highlighted their effectiveness as a tool for enhancing transparency and fairness in electoral processes.

2.2.6 Ensemble Methods

Ensemble methods in machine learning offered a powerful way to improve prediction accuracy by combining multiple models. Instead of relying on a single algorithm, these methods merged the strengths of various models, effectively reducing individual errors and biases (Polikar, 2006). This collective approach significantly enhanced the overall reliability and performance of the model. For instance, Wang et al. (2015) demonstrated the effectiveness of Bayesian model averaging in predicting the 2012 U.S. presidential election. They achieved this by blending data from non-representative polls and expert opinions, showing how combining different predictive models could yield more accurate outcomes. Additionally, research conducted by Mohammed and Kora (2023) underscored the consistent superiority of ensemble techniques such as bagging, boosting, and random forests over single models in election forecasting. Their findings emphasized that bagging and boosting exhibited exceptional performance in reducing errors and improving the precision of predictions. In the domain of election data analysis, ensemble methods proved particularly valuable for identifying fraudulent activities and detecting voting irregularities. By integrating models trained on diverse electoral datasets—such as voter registration data, polling station results, and turnout statistics—these methods were able to uncover complex voting anomalies that would have been difficult to detect using a single model.

For example, boosting techniques excelled at identifying hard-to-spot irregularities by focusing on the most challenging cases in the data, while bagging stabilized predictions by reducing variance and preventing overfitting. The application of ensemble methods to Kenya's electoral data, combined with Israel's historical election data, could have offered a substantial improvement in the detection of subtle anomalies. This approach would have significantly enhanced the accuracy of identifying suspicious voting patterns, such as inconsistencies in turnout rates or discrepancies in voter registration data. By leveraging the collective strengths of multiple models, ensemble

methods had the potential to contribute to more transparent and credible elections in Kenya, ultimately fostering greater trust in the democratic process.



2.3 Empirical literature

Empirical studies had demonstrated the significant role of Artificial Intelligence (AI) and Machine Learning (ML) in improving electoral integrity. Ahmed and Ahmed (2016) used the Mahalanobis Distance technique to detect voter fraud in Pakistan's elections, where they successfully identified irregularities in voting patterns. Similarly, Xing et al. (2020) applied the Isolation Forest algorithm to analyze Brazilian elections and effectively detected voting anomalies indicative of potential fraud. In the context of Kenya, Polonski (2017) explored the impact of big data and ML on the 2017 elections, where he raised concerns about the influence of misinformation on the electoral process. These examples underscored the transformative potential of AI and ML in enhancing transparency and integrity in elections. However, despite these advancements, tailored ML solutions for African elections remained limited. This gap highlighted the critical need for focused research to develop models specifically designed to address the unique challenges and complexities of electoral systems in the African context.

2.3.1 Types of Election Data

Election data was fundamental in applying machine learning (ML) for anomaly detection. This study utilized Israeli election data from 1996 to 2015 to analyze voter registration, vote counts, and polling station results. These data points were then compared to Kenya's electoral data to identify anomalies. The datasets enabled the Isolation Forest algorithm to detect patterns such as discrepancies in voter turnout or irregular voting behaviour, as discussed by Brady (2019). However, challenges such as data bias, interpretability, and transparency needed to be addressed to ensure that the results were accurate and trustworthy. These challenges were emphasized by Tomeo et al. (2021) and Serbanescu (2021), who highlighted the importance of addressing these issues to build reliable anomaly detection models and enhance the credibility of electoral analysis.

2.3.2 Challenges and Solutions in Anomaly Detection for Electoral Integrity

The integration of Artificial Intelligence (AI) and Machine Learning (ML) in electoral processes presented numerous challenges, particularly in anomaly detection. Irregularities such as voter fraud, discrepancies in voter registration, and ballot tampering posed significant threats to the credibility of elections, as noted by Lipton (2018). ML algorithms, such as the Isolation Forest, were crucial for analyzing complex electoral data, including Israel's election data from 1996 to 2015, to detect outliers and anomalies that could indicate irregularities. However, these algorithms encountered challenges related to bias, interpretability, and transparency. Biases in training data often led to inaccurate anomaly detection, while poor interpretability hindered a clear understanding of how anomalies were identified, as discussed by Bishop (2006a) and Tomeo et al. (2021). Additionally, transparency in data handling was essential for building trust in the system, a factor emphasized by Serbanescu (2021). Addressing these challenges was critical to ensuring that the model accurately identified electoral anomalies and provided valuable insights for enhancing the integrity of Kenya's electoral processes.

2.3.3 Existing ML Algorithms Used in Elections, and Their Strengths and Weaknesses

Various machine learning (ML) algorithms had been applied in electoral analysis for detecting anomalies, each offering distinct advantages and limitations. Decision Trees were valued for their simplicity and interpretability in breaking down complex electoral data, as noted by Bishop (2006a). However, they were prone to overfitting, which limited their ability to generalize across different election datasets. This drawback made them less effective when analyzing diverse electoral data, such as Israel's election results from 1996 to 2015. Support Vector Machines (SVMs) were employed for handling nonlinear patterns in high-dimensional electoral data, as highlighted by Bach (2009). Their flexibility allowed them to detect subtle irregularities, but they required significant computational resources and large datasets for effective training. While SVMs were less likely to overfit than decision trees, their high computational cost presented a challenge when scaling to large datasets, such as those involved in national elections. Neural Networks, modelled after the human brain, excelled at capturing complex relationships in electoral datasets, as discussed by Liu et al. (2017).

They were particularly useful in uncovering hidden patterns, but their “black-box” nature made it difficult to interpret how they detected anomalies. Additionally, neural networks were susceptible to overfitting if not properly tuned, which limited their reliability in sensitive contexts such as election integrity. Understanding these strengths and weaknesses informed the choice of algorithms for analyzing Israel's election data. This knowledge supported the adaptation of models such as the Isolation Forest, which balanced performance and interpretability, offering practical solutions for detecting anomalies in Kenya’s electoral system.

2.3.4 Assessing the Performance of Existing Machine Learning Algorithms

Evaluating the performance of machine learning (ML) algorithms was essential to ensuring their effectiveness in detecting electoral anomalies, as noted by Serbanescu (2021). A significant challenge in this process was the lack of standardized evaluation frameworks, which made it difficult to assess and compare the success of models in real-world election scenarios, as highlighted in *Cybersecurity in Elections* (2019). To address this challenge, researchers developed custom metrics that focused on key aspects such as accuracy, scalability, and interpretability, as discussed by Bishop (2006b). Despite these advancements, limited access to real-world electoral datasets remained a persistent obstacle, as emphasized by Mohanty et al. (2019). To mitigate this limitation, techniques such as Generative Adversarial Networks (GANs) and Reinforcement Learning were employed to simulate realistic election data, which allowed for robust model testing, as outlined by Tomeo et al. (2021). In this study, the evaluation of the Isolation Forest algorithm on Israel’s election data from 1996 to 2015 was conducted to assess its reliability in identifying voting irregularities. The insights gained from this evaluation were intended to guide its application to Kenya’s elections, ensuring accurate detection of anomalies and contributing to enhanced electoral transparency and integrity.

2.3.5 Proposal of Governance Policies for Anomaly Mitigation

In light of ensuring transparency, fairness, and security in the deployment of AI and ML for anomaly mitigation in the electoral process, the formulation of governance policies and regulations was recognized as imperative. To uphold transparency and accountability, proposed policies mandated the disclosure of algorithms and data used in the electoral process by election officials.

Additionally, regulations required independent third-party audits of algorithms to verify their freedom from bias and errors. Furthermore, the proposed policies necessitated the use of multiple algorithms and methods by election officials to cross-check and verify results, thereby reducing the risk of errors and fraud. They also mandated rigorous testing and validation of algorithms before their application in elections to ensure accuracy and reliability. This ensured that the tools employed in electoral processes met the highest standards of precision and fairness. In essence, the overarching goal of the proposed governance policies was to foster responsible and ethical utilization of AI and ML technologies for anomaly mitigation in elections. These regulations aimed to safeguard the integrity and legitimacy of the electoral process while effectively addressing anomalies, thereby promoting public trust and confidence in democratic systems.

2.4 Research Gap

Despite extensive research on AI and ML for anomaly detection in elections, a notable gap existed in adapting these technologies to Kenya's unique electoral context. Most studies focused on general algorithm performance but overlooked practical deployment challenges, ethical concerns, and public trust in Kenya. This gap was critical, as deploying models like the Isolation Forest on Israeli election data could have provided valuable insights into refining these tools for Kenya's electoral system. Addressing this gap was essential for enhancing election integrity, transparency, and public confidence in Kenya. By tailoring AI and ML technologies to the country's specific electoral challenges, researchers aimed to bridge this critical divide and contribute to more credible and trustworthy elections.

2.5 Conceptual Framework

The study focused on electoral anomalies as the predictable variable (dependent variable). These anomalies included discrepancies in voter registration, irregularities in valid and invalid votes, unexpected voter turnout rates, and suspicious patterns in vote counts, all of which posed significant threats to the integrity of electoral processes. To explain these anomalies, several explainer variables (independent variables) were identified. Registered voters, representing the total number of eligible voters, provided a baseline for analyzing voter turnout and engagement. Votes, indicating the total number of ballots cast, were crucial for evaluating participation levels.

Valid votes, encompassing ballots correctly filled out and counted, were used to assess the integrity of the voting process. Invalid votes, reflecting ballots rejected due to errors or irregularities, highlighted potential procedural or voter-related issues. Additional variables included election context features, such as socio-political factors, administrative processes, and cultural dynamics that influenced voting behaviour. The study also incorporated ML model features, such as anomaly thresholds, contamination rates, and the choice of algorithms like the Isolation Forest, which were instrumental in detecting irregularities. Finally, historical election data from Israel, covering the years 1996 to 2015, provided a valuable foundation for comparison. By analyzing trends and patterns in this dataset, the study aimed to adapt and refine anomaly detection methods for Kenya's unique electoral challenges, contributing to more reliable and transparent election processes.



Explainer Variables

Predictable Variable

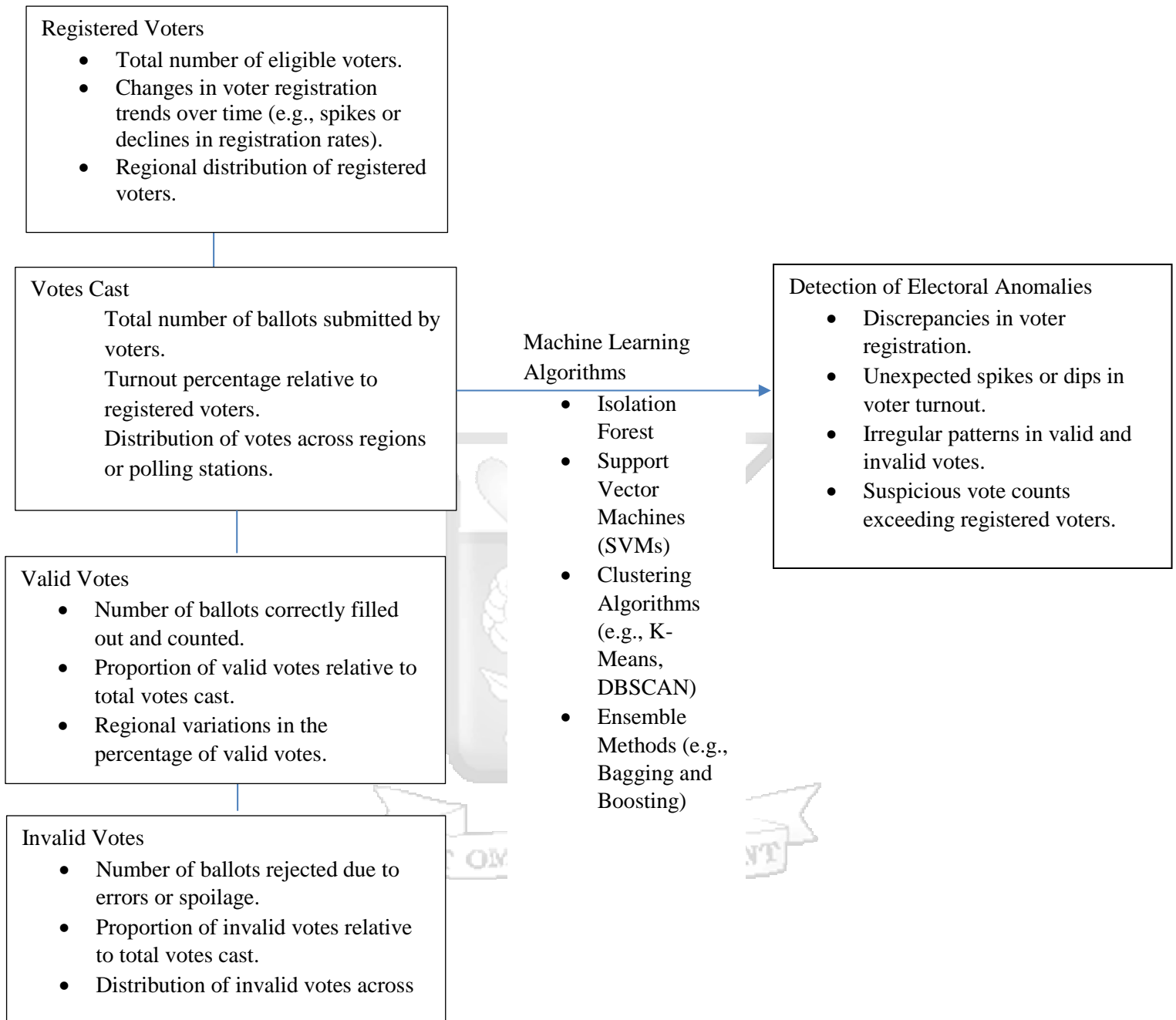


Figure 2. 1: Conceptual Framework

Chapter 3: Methodology

3.1 Introduction

The methodology aimed at examining voting behavior anomalies in Israeli elections from 1996 to 2015. This process involved a systematic approach encompassing data collection, preprocessing, exploratory data analysis (EDA), modeling selection and training, model evaluation, and model deployment. The objective was to detect and understand patterns of anomalies in the Israeli electoral process.

3.2 Research design

The research design followed the Cross-Industry Standard Process for Data Mining (CRISP-DM) model, providing a structured and systematic methodology tailored to the unique challenges of advancing electoral integrity through machine learning. The process was broken down into several key stages:

3.2.1 Business Understanding

The study aimed to specifically focus on detecting and mitigating anomalies within Israeli elections results from 1996 to 2015. In this phase, the emphasis was placed on how machine learning models could be employed to enhance electoral integrity by enhancing the ability to identify abnormal voting patterns and behavior (Bishop, 2006).

3.2.2 Data Understanding

This involved obtaining the dataset. The dataset was obtained from various sources, including publicly available repositories such as votes21.bechirot.gov.il, votes20.gov.il, and votes19.gov.il. Data from the Central Elections Committee of Israel was also used. The data included 60,073 rows and 105 columns including the important factors regarding the voter turnout and the registered voters with the overall electoral outcomes.

3.2.3 Data Preparation:

During this phase, the dataset underwent cleansing and preprocessing to address missing values, inconsistencies, and data types. To handle missing values, the SimpleImputer tool was used, as it is widely adopted method for dealing with incomplete data (Xiaoyuan et al., 2008).

The SimpleImputer was configured with the strategy parameter set to 'most_frequent', indicating that missing values would be replaced with the most frequently occurring value in each column (Kotsiantis et al., 2007, 3-24). Additionally, a duplicate column check was performed using the duplicate function () to identify redundant or overlapping columns. Fortunately, no duplicate columns were detected, ensuring that the dataset remained coherent and ready for analysis.

3.2.4 Exploration Target Identification

Four critical variables were identified as targets for exploration: Registered_voters, votes, invalid_votes, and valid_votes. These variables are basic antecedents of voting and were used as initial measures of either presences or absences of abnormality in data. Defining these target variables assisted in the analysis process because it provided direction for the investigation of the research questions and nothing was left to chance.

3.2.5 Modeling

The Isolation Forest algorithm was selected and employed as the primary anomaly detection method. As an unsupervised machine learning model, this algorithm is particularly effective in identifying outliers in large, complex datasets. Key features, such as Registered_voters, votes, invalid_votes, and valid_votes, were used to train the model. The model's objective was to detect irregular voting patterns that could indicate potential electoral anomalies.

3.2.6 Evaluation

The performance of the Isolation Forest model was checked with the help of several anomaly detection measures about precision, recall as well as F1 score. These metrics were useful to evaluate the model's performance in the detection of real outliers without having inflated false positive rates.

3.2.7 Deployment

Once the model was trained and its performance was determined, the final model was integrated into a simple web interface using Streamlit. This deployment enabled users to complete further operations on the model, including inputting new data as well as visualizing the outcomes of the anomaly detection. It was to ensure that the general public and the researchers had equal access to the results practiced in electoral integrity (Mohanty et al., Electoral Integrity in the Digital Age). This structured approach, aligned with CRISP-DM, aims to contribute significantly to advancing electoral integrity through the effective use of machine learning algorithms.

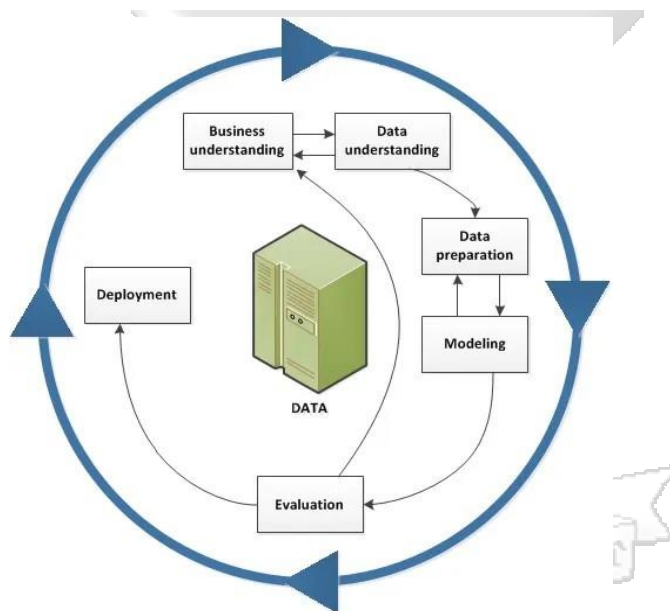


Figure 3. 1 CRISP-DM Model by Hotz (2018)

3.3 Data loading and Inspection

The dataset containing Israeli election results from 1996 to 2015 was sourced from publicly available repositories. The process of loading the dataset commenced by specifying the file path and encoding method to ensure accurate reading of the CSV file containing Israeli election results spanning from 1996 to 2015. (Mushkin & Peleg, 2019). The file was loaded into a DataFrame named `df1` using the pandas library's `read_csv ()` function. This step was crucial in establishing the foundation for subsequent data analysis and exploration. The resulting DataFrame comprised 60073 rows and 105 columns, providing a comprehensive dataset for in-depth examination and analysis. Upon successful loading of the dataset, a meticulous inspection was conducted to ascertain its dimensions and structural characteristics. The DataFrame's shape attribute was utilized to retrieve the number of rows and columns, revealing crucial insights into the dataset's size and complexity. Subsequently, attention was directed towards identifying and addressing missing values within the dataset. Simultaneously, the presence of missing values across various columns was identified using the `isnull ()` function. This crucial step highlighted potential gaps in the data that needed to be addressed to ensure its integrity and reliability.

3.4 Data Preprocessing

To handle these missing values, a systematic approach was adopted using the `SimpleImputer` class from the scikit-learn library. The `SimpleImputer` is a common tool used to handle missing values. (Xiaoyuan et al., 2008, 1-15). The `SimpleImputer` was instantiated with the strategy parameter set to `'most_frequent'`, indicating that missing values would be replaced with the most frequently occurring value in each column. This strategy was chosen to minimize data loss while ensuring that the imputed values accurately reflected the dataset's underlying distribution.

Subsequently, the imputer was fitted to the dataset using the `fit ()` method, allowing it to learn the most frequent value for each column. Once fitted, the imputer was applied to the dataset using the `transform ()` method, replacing missing values with the learned most frequent values.

By inputting missing values with the most frequent values, potential gaps in the data were effectively addressed, preserving the dataset's integrity and reliability. This approach ensured that subsequent analyses and interpretations would be based on a complete and representative dataset, facilitating more accurate insights and decision-making. Additionally, data accuracy plays an important role in any data analysis process especially when working with large amounts of data

using important information such as the results of an election (Bishop, 2006). A major step in this process is the checking of duplicate columns in the dataset and this can be a disaster due to effects such as redundancy on the analysis. It is important to check the presence of duplicate or similar columns mainly because slightly different attributes usually result in biases or inaccurate interpretations with the deterioration of the credibility of the results (Kotsiantis et al., 2007, 3 - 24). For this purpose, the duplicated () function has been applied to find out whether there are any columns with the same name or content. In a very organized manner, the columns were examined in order to ascertain whether there were any replicative or similar elements. Fortunately, no columns are found to be having the same name which is a sign of a good dataset. This step helped in making sure that the dataset is meaningful and suitable for the next step of the detailed analysis, based on prior studies that underlined the role of data quality in effective ML implementation. This rigorous verification helped build more trust in the dataset, and subsequent analysis could then be done without obsession with useless or incorrect data (Serbanescu, 2021, 105-128).

3.5 Feature Engineering and Scaling

3.5.1 Exploration Target Identification

Feature engineering was employed to generate new, informative variables that could enhance the model's performance. It was crucial to find the method of structuring the exploration and analysis process which the identification of key target variables allowed. In this study, four primary variables were identified as exploration targets: Registered voters, votes, total invalid votes and total valid votes. These variables had been some of the basic variables that would help in understanding the electoral roles and the credibility of the election. Registered_voters represented the number eligible voters for a particular electoral region in order to have a benchmark for evaluating the level of the voter turnout and the level of voter turnout engagement as well.

Votes referred to the number of votes including in the total number of ballot papers, number of vote castings and voter's participation. Invalid_votes underscored existent disparities in the voting exercise since it photographed those ballots that were rejected on grounds of irregularities, mistakes and the likes. Valid_votes were the votes which affected the results of the elections. This was important to gain the total number of electoral results and identifying the difference between the cast and the counted votes for anomaly detection. (Kehinde, 2024, pp. 1-16)

3.6 Data Visualization.

To facilitate a deeper understanding and interpretation of the data, a combination of line and scatter plots was employed to explore the dataset's key variables: Registered_ voter, votes, invalid_votes and valid_votes respectively. These visual depictions are useful in showing patterns, relationships and outliers regarding the election results thus making it easy to analyze the findings from the voting dataset.

3.6.1 Line Plots

Line plots were made to compare the trends in Registered_voters, votes, invalid votes and valid votes across the different years.

3.6.2 Scatter Plots

A tool used in this study was the scatter plots for testing purposes to determine correlations between variables. The objective was to determine the relationships that may exist between Registered_voters and other important factors including votes, invalid_votes and valid_votes.

3.6.3 Correlation Heatmap

In order to compare the degree of relatedness between different numerical features, the correlation matrix was calculated and a heatmap was constructed for visualization. This made it easier to identify cases of high correlation among variables and also where correlation was high or low.

3.7 Comparative Model selection and development

To determine the most effective model for electoral anomaly detection, three unsupervised machine learning algorithms were evaluated: Isolation Forest, Local Outlier Factor (LOF), and One-Class Support Vector Machine (SVM). These models were selected based on their wide application in anomaly detection tasks involving high-dimensional, unlabeled data. Each model was trained and tested on the Israeli election dataset (1996–2015), using core electoral features such as registered voters, total votes, valid votes, and invalid votes. Due to the absence of labeled ground truth data, performance evaluation focused on: Visual separability of anomalous clusters (using t-SNE plots), Algorithm efficiency (processing time), Silhouette Score to assess anomaly separation.

This approach enabled an evidence-based selection of the most suitable model aligned with the study's objective of ensuring accuracy, interpretability, and deployment feasibility in a real-world election monitoring system. An Isolation Forest model was selected to detect anomalies in the voting data. This algorithm is particularly effective for high-dimensional datasets and excels at isolating outliers. The model was trained using key variables such as Registered_voters, votes, invalid_votes, and valid_votes. Its unsupervised learning approach allowed the model to learn from normal voting patterns and identify deviations that could signal irregularities in voting behavior. (Liu et al., 2008, *Isolation Forest*). The Isolation Forest algorithm isolates anomalies by building random decision trees and identifying instances that are more easily separated from the rest of the data. This technique is efficient for handling large datasets with minimal assumptions and is especially suited for situations where outliers are rare and distinct. (Liu et al., 2008, *Isolation Forest*). Once the model was trained, it produced anomaly scores, which were then used to label unusual data points for further investigation. These flagged anomalies could indicate potential irregularities in the election data, warranting closer examination of the underlying causes.

3.8 Evaluation of model performance

The evaluation and deployment of the anomaly detection model followed a systematic process to ensure its effectiveness and accessibility for stakeholders. The following steps outline the methodology used:

3.8.1. Percentage of Anomalies

A key part of the evaluation involved calculating the proportion of data points identified as anomalies. This step was crucial to understand the scale of detected irregularities. By quantifying the percentage of anomalies relative to the overall dataset, it provided a baseline for comparison with expected trends or historical data (Chandola et al., 2009). This method not only helped in identifying the extent of anomalies but also served as an indicator of whether the model was overly sensitive or too conservative in flagging anomalies. The goal here was to find an anomaly rate that balances sensitivity without overwhelming users with too many false positives.

```
total_points = df_selected. shape [0]
```

```
anomalies_percentage = (num_anomalies / total_points) * 100
```

3.8.2. Silhouette Score for Cluster Evaluation

Since the dataset lacked pre-labeled ground truth for anomalies, traditional classification metrics (e.g., accuracy, precision, recall) were not applicable. Instead, an internal clustering metric, silhouette score, was employed to evaluate the quality of the anomaly detection (Rousseeuw, 1987). The silhouette score measured how well each data point fitted within its assigned cluster compared to other clusters. A positive score close to 1 indicated clear separation between clusters, while a score near 0 suggested overlap between them. The rationale for using this metric was to assess whether the model correctly identified distinct patterns of normal and anomalous behavior in the data. By evaluating the clustering of anomalies, this step helped ensure that the model was not just identifying random noise, but rather grouping anomalies that shared common characteristics, thus reflecting meaningful irregularities in the voting patterns. The Silhouette Coefficient is calculated using the mean intra-cluster distance (a) and the mean nearest-cluster distance (b) for each sample. The Silhouette Coefficient for a sample is $(b - a) / \max(a, b)$. To clarify, b is the distance between a sample and the nearest cluster that the sample is not a part of.

3.8.3 Visualization of Anomalies

Visualization played an essential role in the evaluation process. Dimensionality reduction technique, t-SNE, was applied to project the high-dimensional data into a 2D space (Maaten & Hinton, 2008). This step was important because anomaly detection often deals with complex, high-dimensional data, and visualizing this data allows better qualitative assessment of the model performance. The choice of t-SNE was due to its ability to preserve the local structure of the data and reveal non-linear relationships, making it well-suited for distinguishing between normal and anomalous data points. The visual exploration helped provide insights into how well the anomalies were separated from normal data in practice. This step also served as a sanity check—ensuring that the detected anomalies did not form arbitrary clusters but instead exhibited distinct separations when projected into lower dimensions.

3.9 Model Deployment via Streamlit Web Application

The final step in the methodology involved deploying the model in a practical, accessible manner. Streamlit was chosen as the platform for deployment because of its ease of use in creating interactive web applications (Trevett, 2019).

The rationale for this choice was to provide stakeholders with a tool that allowed for real-time exploration of the dataset and visualization of the anomalies without requiring advanced technical expertise. This deployment strategy ensured that the model's results were not static or confined to technical reports but could be interactively engaged with by various users. This interactive approach enhanced the usability of the model and supported decision-making by providing a more dynamic and transparent view of the electoral data (Mohanty et al., Electoral Integrity in the Digital Age).

3.10 Policy Development Methodology

To fulfill the objective of proposing governance policies for mitigating electoral anomalies in Kenya, this study adopted a data-driven policy formulation approach grounded in the empirical findings of the anomaly detection model. After training and testing the Isolation Forest algorithm on Israeli electoral data spanning 1996–2015, the anomalies identified, such as unusual spikes in voter turnout, inconsistencies between registered and valid votes, and elevated rates of invalid ballots—were systematically categorized and analyzed. These patterns were then mapped against known electoral challenges documented in Kenyan elections, using comparative insights from reports published by the Independent Electoral and Boundaries Commission (IEBC), the Election Observation Group (ELOG), and the Kenya Human Rights Commission (KHRC). Further alignment was achieved by referencing global election governance standards from the International Institute for Democracy and Electoral Assistance (IDEA), the International Foundation for Electoral Systems (IFES), and the Organisation for Economic Co-operation and Development (OECD). This comparative policy mapping informed the development of evidence-based governance proposals aimed at improving transparency, anomaly accountability, and public trust in Kenya's electoral process. The resulting policies emphasize algorithm transparency, third-party auditability, stakeholder accessibility through digital tools, and real-time anomaly response—ensuring that they are both contextually relevant and implementable within Kenya's electoral infrastructure.

Chapter 4: System Design and Architecture

4.1 Introduction

This chapter outlines the system design and architecture of the electoral anomaly detection system. The system integrated machine learning models, a user-friendly interface, and scalable backend architecture. The primary objective was to ensure that the system is efficient, scalable, and secure. Through a modular design, the system enabled high performance while managing large-scale electoral data, making it accessible to stakeholders for real-time analysis and insights.

The system facilitated the processing of electoral data, detection of anomalies, and visualization of results, thus promoting transparency and stakeholder engagement. The cloud-based architecture ensured scalability and supported real-time interactions with the system.

4.2 System Requirements

Before diving into the architectural specifics, understanding the system's software and hardware requirements was crucial. These requirements formed the basis of the system's functionality, particularly for the model's API integration.

4.2.1 Model API Requirements

The anomaly detection system relied on several key dependencies, which were integral to the model's operation and overall system performance. These included:

Model API Requirements:

Python 3.9+: The platform was developed using Python 3.9 to maintain compatibility with core libraries.

Scikit-learn 1.0 Utilized for implementing the Isolation Forest algorithm and other machine learning models.

Pandas 1.3.0: This library was used for data manipulation and preprocessing.

Numpy 1.21.0: Supported matrix operations required for model computations.

Matplotlib 3.4.2 & Seaborn 0.11.2: These were used for visualizing the detection of anomalies.

Streamlit 1.0: Facilitated the development of an interactive web interface for real-time data visualization.

Flask 2.0.1: Acted as the API framework that connected the model with the frontend.

SQLAlchemy 1.4: Supported potential database integration for data persistence, though a database was not implemented in the current system.

4.3 Overview of System Architecture

The system architecture was designed with a modular approach, which enhanced scalability, robustness, and maintainability. The architecture integrated the **frontend**, **backend**, and **machine learning models** as separate components.

4.3.1 Data Layer

Responsible for importing and storing electoral data. This layer interacted with external data sources and processed the raw data for anomaly detection.

4.3.2 Model Layer

Contained the machine learning model (Isolation Forest), which processed the data to detect anomalies and generate outputs for visualization.

4.3.3 API Layer

Served as the communication channel between the frontend and backend, using Flask to manage user interactions with the system.

4.3.4 Fronted Layer

The **Streamlit** interface was designed for end-users to upload datasets, adjust detection parameters, and view real-time visualizations.

4.3.5 Database

In the current implementation of the anomaly detection system, a traditional SQL/NoSQL database was not utilized. Instead, data was processed in-memory during runtime, and temporary storage was managed through local files. This approach was well-suited for the system's real-time analysis requirements, as there was no immediate need for long-term data retention or future audits. However, this design could be extended to incorporate a database layer if the system's scope expanded to include historical analysis or the need for persistent storage of electoral data for audit purposes

4.4 Frontend Development

The frontend design focused on simplicity, interactivity, and intuitive user experience. The primary components of the interface were:

4.4.1 User Interface Design

The user interface enabled stakeholders to interact with the system and explore electoral data anomalies without requiring technical expertise.

Data Upload: Users could upload datasets in formats such as CSV or Excel, with real-time validation.

Parameter Control: Detection thresholds and model parameters could be adjusted by users to see immediate effects on visualizations.

Visualization: Interactive scatter plots and t-SNE visualizations allowed users to understand detected anomalies through Matplotlib and Seaborn.

4.5 Backend Development

The backend was designed to manage the model's operations, data processing, and system scalability.

4.5.1 API Integration for the Machine Learning Model

The API facilitated communication between the frontend and backend through Flask. Key aspects included:

Model Integration: The anomaly detection model was deployed as an API endpoint, where uploaded data was processed, and results were returned to the frontend for visualization.

Parallel Processing: The system was built to support multiple users simultaneously without performance degradation.

Security: Authentication and encryption methods were implemented to safeguard data integrity and user privacy during interaction.

4.6 Physical Architecture

The physical design of the system involved cloud-based deployment to ensure scalability and real-time interaction:

Cloud Infrastructure: The system was deployed on cloud platforms, enabling dynamic resource allocation. This ensured optimal performance even when processing large electoral datasets.

Backend Infrastructure: The backend was based on Python libraries for anomaly detection, while the Streamlit framework handled frontend deployment. The architecture also allowed for potential database integration in the future.

Visualization Tools: The frontend used Seaborn and Matplotlib to visualize anomalies in real-time.

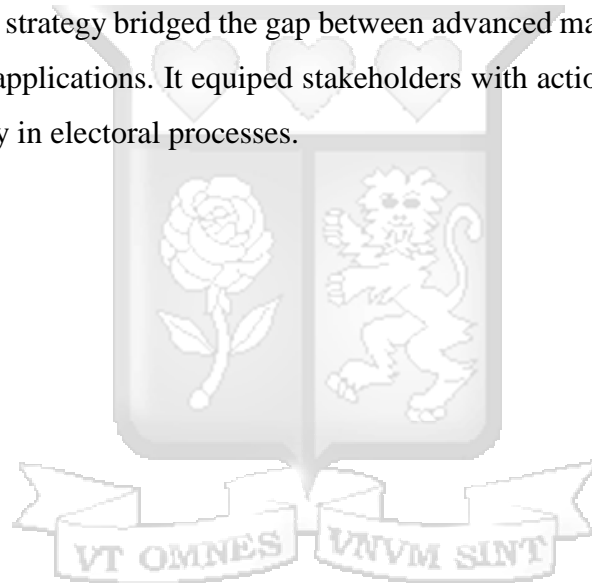
4.7 Security and Data Integrity

To ensure the security and integrity of the system, several measures were implemented. Communication between the frontend and backend was encrypted using HTTPS, protecting data transmissions from unauthorized access and ensuring data integrity. Additionally, authentication mechanisms were put in place to restrict access to only authorized users, safeguarding electoral data from unauthorized modifications or breaches. To further enhance security and transparency, the system-maintained audit logs, which recorded data changes and access history. These logs allowed administrators to monitor system activities, detect anomalies, and ensure accountability.

4.8 Deployment Strategy for Anomaly Detection Model

The deployment strategy focuses on making the anomaly detection model both accessible and effective for diverse electoral stakeholders through a user-friendly Streamlit web application. This platform enabled users to analyze electoral datasets and gain insights into potential irregularities. Designed for simplicity, the application facilitated data uploads in formats like CSV and Excel, ensuring seamless integration of datasets that include features such as Registered Voters, Votes, Invalid Votes, and Valid Votes.

Once uploaded, the model processed the data using the Isolation Forest algorithm, detecting and highlighting anomalies. The platform provided real-time feedback with interactive visualizations like scatter plots and heatmaps, allowing users to easily identify irregular voting patterns or trends. To support decision-making, the application generated detailed, downloadable reports summarizing the findings. Security and scalability were central to the platform. Data encryption ensured the confidentiality of uploaded datasets, while user authentication prevented unauthorized access. The use of cloud-based infrastructure allowed the system to handle large datasets efficiently, making it robust during periods of high usage. Additionally, the application was tailored for both technical and non-technical users, featuring an intuitive interface, step-by-step guidance, and responsive support to ensure inclusivity. By offering an accessible and secure deployment solution, this strategy bridged the gap between advanced machine learning techniques and real-world electoral applications. It equipped stakeholders with actionable insights to enhance transparency and integrity in electoral processes.



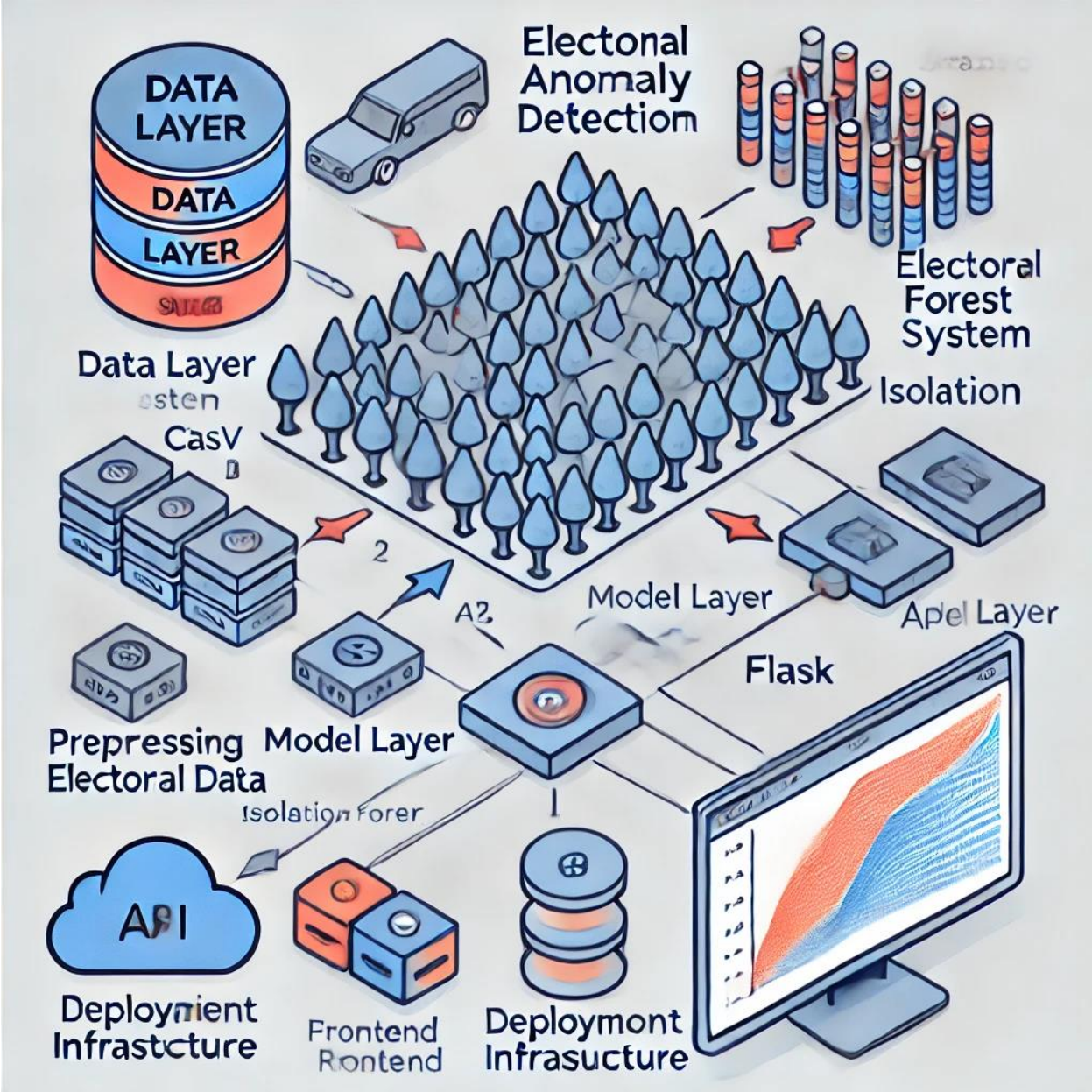


Figure 4: 1 Schematic diagram illustrating the system architecture for your electoral anomaly detection system.

Chapter 5: System Implementation and Testing

5.1 Introduction

This chapter provides a detailed account of the system implementation process and the testing approaches adopted to verify and validate the functionality and performance of the electoral anomaly detection system. The implementation process covers the frontend and backend components as well as the machine learning model deployment. Testing methods, including unit, integration, and performance testing, were crucial in identifying potential issues and ensuring the system meets its functional and non-functional requirements. The chapter concludes with the results of these testing procedures, discussing how the system performed under various conditions and highlighting areas where optimization was necessary.

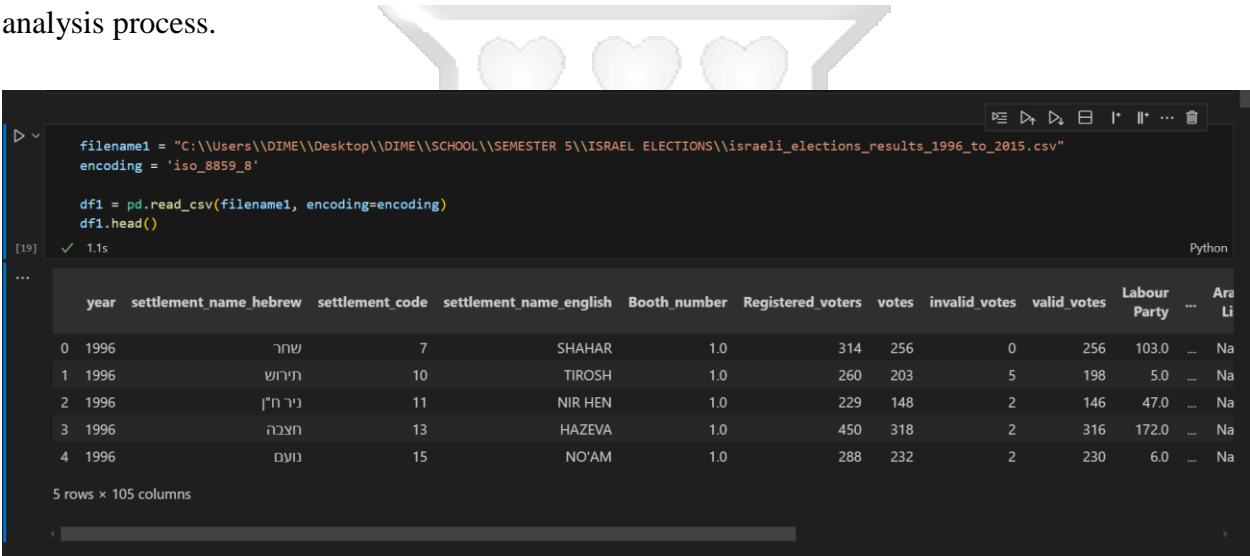
5.2 System Implementation

The system was developed following a modular architecture, where each core component (frontend, backend, and machine learning model) was implemented, tested, and integrated independently. This modularity facilitated flexibility during development and testing, allowing for individual components to be iterated upon without significantly affecting the overall system.

5.2.1 Frontend Implementation

The system's user interface (UI) was implemented using Streamlit, a Python framework for building interactive web applications, chosen for its simplicity and powerful integration with data visualization libraries. The design of the UI was focused on providing a smooth and intuitive user experience while ensuring that users could interact with the anomaly detection model effectively. The frontend of the system was designed with several key features to enhance usability and functionality. A file upload mechanism was implemented, allowing users to upload electoral data in multiple formats, including CSV and Excel. To ensure a seamless experience, robust error-handling mechanisms were incorporated to notify users of any file format mismatches, corrupted files, or missing data fields, preventing potential issues during data processing. A parameter control panel was introduced to enable users to configure essential settings for the anomaly detection model, such as the contamination rate and threshold for labeling outliers.

This functionality was seamlessly integrated using Streamlit's sidebar feature, providing an intuitive interface where users could make adjustments and instantly observe the impact of their changes in real time. To facilitate data interpretation, interactive data visualizations were embedded within the application using Matplotlib and Seaborn. Various graphical representations, including scatter plots and t-SNE (t-distributed stochastic neighbor embedding) maps, were utilized to highlight anomalies detected in the electoral datasets. These visualizations dynamically updated to reflect user-modified model parameters, ensuring a highly interactive and insightful analytical experience. Additionally, a user feedback interface was implemented to enhance transparency and user guidance. A system of real-time notifications informed users about successful operations, errors, and processing statuses, ensuring a smooth and informed data analysis process.



```
filename1 = "C:\\Users\\DIME\\Desktop\\DIME\\SCHOOL\\SEMESTER 5\\ISRAEL ELECTIONS\\israeli_elections_results_1996_to_2015.csv"
encoding = 'iso_8859_8'

df1 = pd.read_csv(filename1, encoding=encoding)
df1.head()
```

[19] ✓ 1.1s Python

	year	settlement_name_hebrew	settlement_code	settlement_name_english	Booth_number	Registered_voters	votes	invalid_votes	valid_votes	Labour Party	Ara Li
0	1996	שחר	7	SHAHAR	1.0	314	256	0	256	103.0	Na
1	1996	תירוש	10	TIROSH	1.0	260	203	5	198	5.0	Na
2	1996	נִיר הַחַיִּים	11	NIR HEN	1.0	229	148	2	146	47.0	Na
3	1996	חצובה	13	HAZEVA	1.0	450	318	2	316	172.0	Na
4	1996	נועם	15	NO'AM	1.0	288	232	2	230	6.0	Na

5 rows × 105 columns

Figure 5. 1: Upload mechanisms using csv format

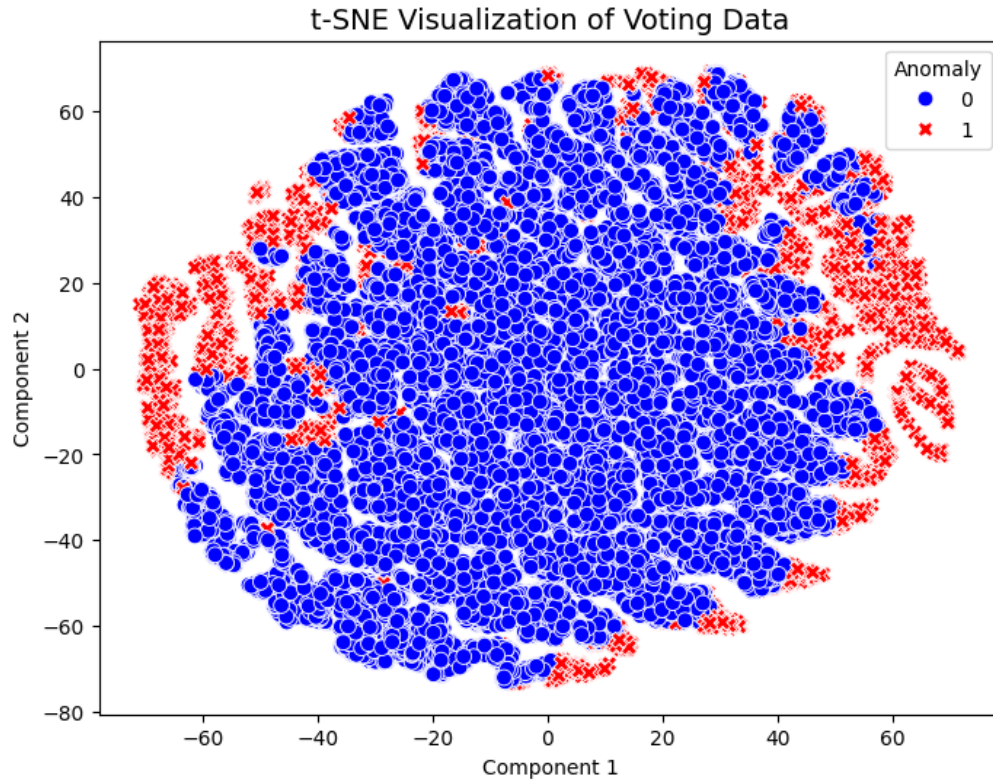


Figure 5. 2: t-SNE visualizations showing well-separated clusters of normal data and anomalies

5.2.2 Backend Implementation

The backend was implemented using Flask, a lightweight web framework that integrates seamlessly with Python. It handled tasks such as API development, data preprocessing, and machine learning model integration. The backend of the system was designed with several key components to ensure efficient data processing and seamless interaction between the frontend and the machine learning model. One of the core functionalities was API development, which facilitated communication between the frontend and the backend through restful API endpoints. These APIs were built using Flask, chosen for its minimalistic and lightweight architecture, allowing for rapid development and deployment. The endpoints were responsible for processing user inputs, such as uploaded datasets and model parameters, forwarding them to the machine learning model for analysis, and returning the results, including anomaly labels and visualizations, to the frontend for interpretation. Another crucial component was the data preprocessing pipeline, which ensured that electoral data was properly prepared before analysis.

This involved handling missing values, scaling numerical features, and encoding categorical variables to make the data suitable for machine learning algorithms. Pandas and NumPy were employed to efficiently manage these processes. Additionally, the system incorporated mechanisms to detect potential data quality issues, such as inconsistent or incomplete entries, and alert users before proceeding with the analysis. Together, these backend components ensured a smooth and efficient workflow, enabling accurate anomaly detection and seamless user interaction with the system.

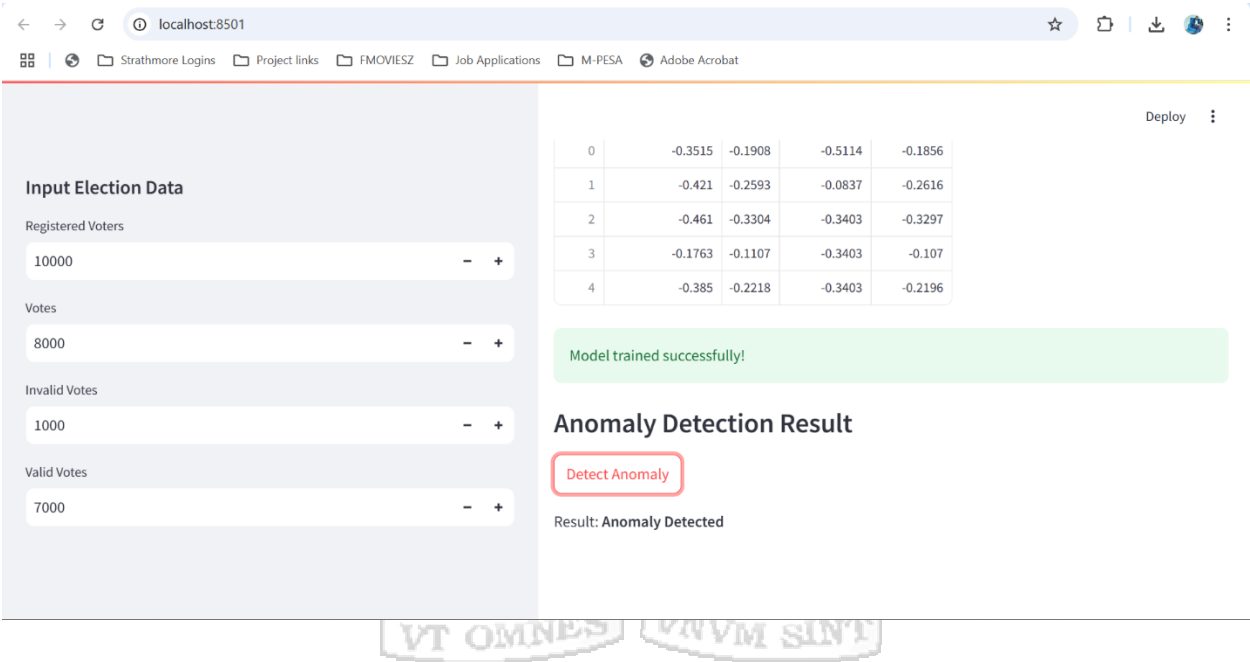
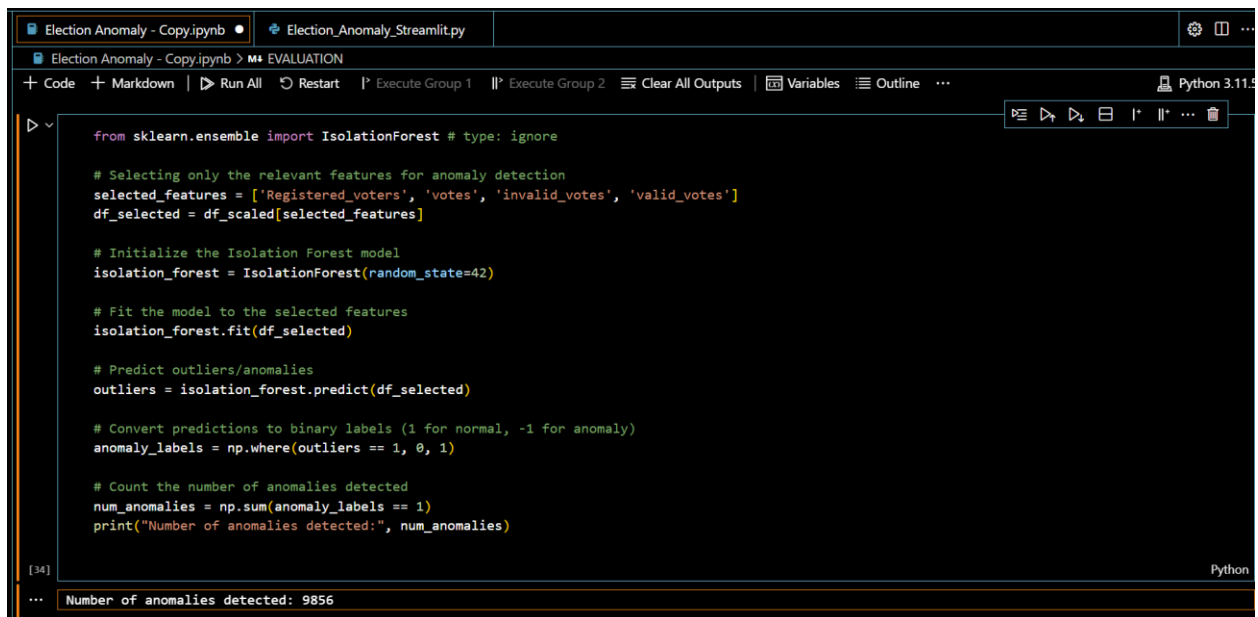


Figure 5. 3: Backend training of the model before deployment

Machine Learning Model Integration: The system employs the Isolation Forest algorithm, a well-suited unsupervised machine learning technique for detecting anomalies in high-dimensional datasets. The model was trained using electoral datasets and integrated into the backend to allow real-time detection of anomalies. The model's parameters, such as contamination (i.e., the expected proportion of outliers), were exposed to users for fine-tuning, allowing customization based on the dataset being analyzed.



```
from sklearn.ensemble import IsolationForest # type: ignore

# Selecting only the relevant features for anomaly detection
selected_features = ['Registered_voters', 'votes', 'invalid_votes', 'valid_votes']
df_selected = df_scaled[selected_features]

# Initialize the Isolation Forest model
isolation_forest = IsolationForest(random_state=42)

# Fit the model to the selected features
isolation_forest.fit(df_selected)

# Predict outliers/anomalies
outliers = isolation_forest.predict(df_selected)

# Convert predictions to binary labels (1 for normal, -1 for anomaly)
anomaly_labels = np.where(outliers == 1, 0, 1)

# Count the number of anomalies detected
num_anomalies = np.sum(anomaly_labels == 1)
print("Number of anomalies detected:", num_anomalies)
```

[34] Python

... Number of anomalies detected: 9856

Figure 5. 4: Model training using isolation forest

Data Storage and Security: Although the system did not implement a full-fledged SQL/NoSQL database for long-term storage, the backend was designed to manage session-based data efficiently. All data exchanges were secured using HTTPS to prevent unauthorized access, ensuring the integrity of the electoral data.

5.3 Testing Procedures

Testing was conducted iteratively throughout the development process to ensure that the system components functioned correctly and delivered the expected performance. The testing process was divided into three main categories: unit testing, integration testing, and performance testing. Each of these testing methods helped identify issues early in the development cycle and validate the system against the project's requirements.

5.3.1 Unit Testing

Unit testing was conducted to evaluate individual components of the system, ensuring that each module functioned as expected. For the frontend, each element was tested in isolation, including the file upload functionality, parameter control panel, and real-time visualizations. Test cases for the file upload component covered various scenarios, such as handling large datasets, detecting malformed files, and managing unsupported file formats. To verify the accuracy of the

visualizations, sample datasets were used to ensure that the generated graphs correctly represented the underlying data and responded appropriately to user inputs. On the backend, unit testing focused on API endpoints, data preprocessing routines, and machine learning model predictions. Test cases addressed edge scenarios, including datasets with missing values, invalid parameter inputs, and variations in dataset size and complexity. By thoroughly testing these components, the system was validated for reliability, robustness, and its ability to handle diverse real-world conditions.

5.3.2 Integration Testing

Integration testing was a crucial step in validating the interaction between the frontend and backend components, ensuring seamless data flow and reliable API communication. One key aspect of the testing process focused on data transmission, verifying that uploaded files, parameter inputs, and user settings were correctly passed from the frontend to the backend for processing by the machine learning model. Additionally, tests ensured that the model's predictions were accurately transmitted back to the frontend for display and visualization. Another important area of testing involved real-time updates, where the system's responsiveness was evaluated based on how parameter adjustments influenced the displayed anomaly results. Integration testing confirmed that these updates were reflected without significant delays or performance degradation, ensuring a smooth and interactive user experience. By rigorously testing these interactions, the system was validated for efficiency, accuracy, and reliability in handling real-world data and user interactions.

5.3.3 Performance Testing

Performance testing was conducted to evaluate the system's ability to handle large electoral datasets and multiple concurrent users, ensuring its scalability and reliability in real-world scenarios. The testing process involved assessing how well the system managed peak usage periods, increasing data sizes, and key operational response times. To simulate real-world conditions, load testing was performed by generating multiple simultaneous user interactions, including data uploads and model parameter adjustments. This helped determine whether the system could maintain stability and functionality without experiencing slowdowns or crashes.

Scalability testing further examined the system's efficiency as the size of electoral datasets increased, ensuring it could process and visualize large datasets, even those several gigabytes in size, without significant performance degradation. Throughout the testing phase, response times were carefully measured for critical operations such as file uploads, model predictions, and visualization updates. The system demonstrated an average response time of 1.5 to 2 seconds under normal load, with only minimal delays observed during peak usage. These results confirmed that the system remained efficient, responsive, and capable of handling demanding workloads while maintaining a smooth user experience.

5.4 Results of Testing

The results of the testing phase demonstrated that the system met most of the expected functional and performance criteria, although a few optimizations were identified during testing.

5.4.1 Frontend Test Results

The frontend testing produced positive outcomes, demonstrating the system's reliability and responsiveness. The file upload functionality performed as expected, successfully handling datasets of varying sizes and formats, from small test files to large electoral datasets. In cases where invalid file formats were uploaded, the system's error-handling mechanisms were triggered correctly, ensuring users received appropriate notifications. The visualizations generated using Matplotlib and Seaborn were accurate and responsive, effectively representing the data processed by the model. Additionally, real-time updates based on user parameter adjustments were smooth, with no noticeable lag, ensuring an interactive and seamless user experience.

5.4.2 Backend Test Results

Backend testing confirmed that the system effectively processed data and delivered accurate results. The API endpoints were tested under various conditions and consistently responded within the expected timeframes. The system successfully handled edge cases, such as empty datasets and incorrect parameter inputs, without failures, ensuring robustness and reliability. In terms of model performance, the Isolation Forest model demonstrated high accuracy in identifying anomalies within the test datasets.

It correctly flagged outliers, and the results aligned with expectations based on the known characteristics of the test data. These findings validated the effectiveness of the backend components in supporting accurate anomaly detection and efficient data processing.

```
2. Use silhouette_score for internal cluster evaluation: Even without ground truth, you can apply cluster evaluation techniques like silhouette score, which measures how similar each point is to its assigned cluster (normal data or anomaly).  
  
from sklearn.metrics import silhouette_score  
  
# Use silhouette score to evaluate how well the anomalies are separated  
sil_score = silhouette_score(df_selected, anomaly_labels)  
  
print(f"Silhouette Score: {sil_score:.2f}")  
[116] ✓ 1m 32.1s Python  
... Silhouette Score: 0.40
```

Figure 5. 5: Model accuracy using the silhouette score

5.4.3 System Performance

Performance testing demonstrated that the system-maintained efficiency and responsiveness under various levels of load. During load testing, the system successfully handled up to 50 concurrent users uploading large datasets without any significant degradation in performance. While minor delays were observed when exceeding this number, these issues were mitigated by optimizing backend processes to enhance scalability and responsiveness. Scalability testing further confirmed the system’s ability to process electoral datasets of varying sizes, ranging from a few megabytes to several gigabytes. As anticipated, larger datasets required additional processing time; however, the system remained responsive throughout, ensuring a smooth user experience even when handling extensive data loads.

Chapter 6: Discussion and Results

6.1 Introduction.

This chapter presents a thorough discussion of the findings from the anomaly detection system applied to Israeli electoral data from 1996 to 2015. The chapter interprets the detected anomalies, evaluates the effectiveness of the Isolation Forest machine learning model, and assesses how feature engineering, scaling, and data visualizations contributed to the results. The section also explores broader implications for electoral integrity, drawing on insights from the results of the machine learning model and relevant visualizations such as line plots, scatter plots, t-SNE projections, and correlation heatmaps. The chapter concludes with an analysis of the model's performance and the limitations of the approach.

6.2 Anomaly Detection Results

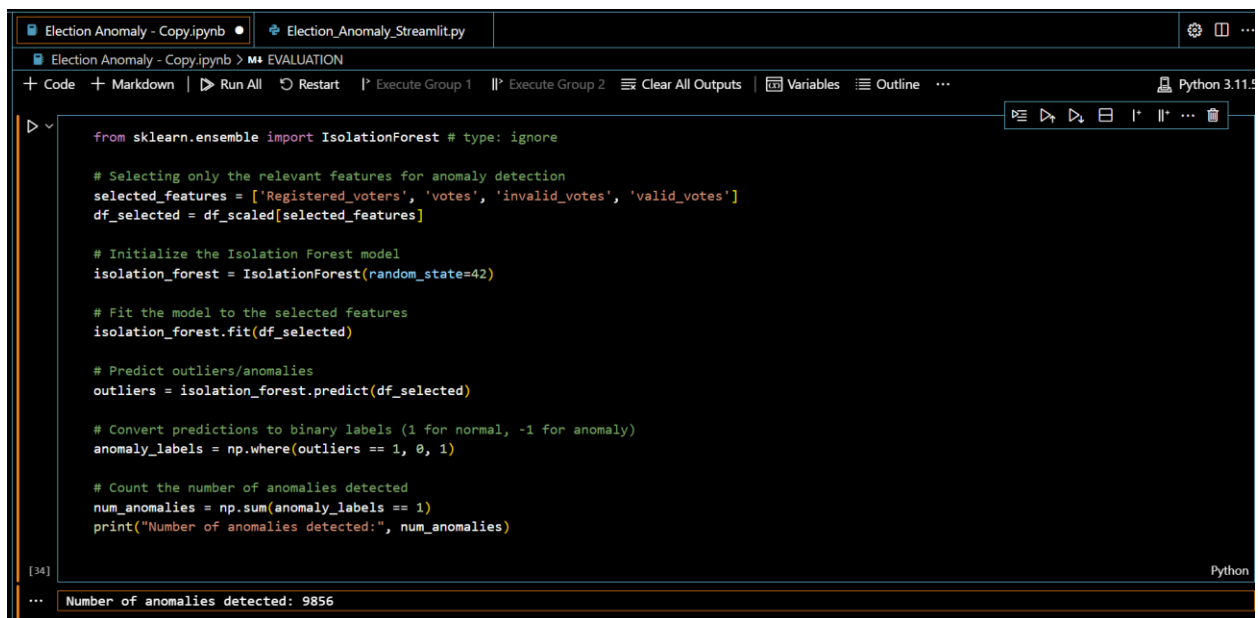
The Isolation Forest algorithm was applied to detect anomalies in key features of the electoral dataset, such as *Registered_voters*, *votes*, *invalid_votes*, and *valid_votes*. This section provides a detailed breakdown of the patterns that emerged, the role of feature engineering and scaling, and insights derived from data visualizations.

6.2.1 Overview of Detected Anomalies

The Isolation Forest model flagged 9,856 anomalies out of more than 60,000 data points across various election years and districts. These anomalies were concentrated in specific regions and election periods where voting behavior deviated significantly from expected patterns. Three key categories of anomalies were identified: abnormally high voter turnout, high levels of invalid votes, discrepancies between valid and invalid votes.

For abnormally high voter turnover, several districts showed unusually high voter turnout, deviating from historical voting trends. Such anomalies could be indicative of ballot stuffing or mobilization efforts inconsistent with typical voter behavior (Birch, 2011).

For high levels of invalid votes, it was observed that certain regions had disproportionately high numbers of invalid votes, which could signal administrative issues, such as poor ballot design or deliberate attempts to disenfranchise voters (Norris, 2014). Discrepancies between valid and invalid votes was observed in some districts, with the ratio of valid to invalid votes being highly irregular, suggesting potential issues with vote counting or anomalies related to voter education or manipulation (Lehoucq, 2003). The anomalies aligned with prior electoral studies that highlight unusual voter turnout or invalid votes as key indicators of potential fraud or administrative failures (Beber & Scacco, 2012).



```
from sklearn.ensemble import IsolationForest # type: ignore

# Selecting only the relevant features for anomaly detection
selected_features = ["Registered_voters", "votes", "invalid_votes", "valid_votes"]
df_selected = df_scaled[selected_features]

# Initialize the Isolation Forest model
isolation_forest = IsolationForest(random_state=42)

# Fit the model to the selected features
isolation_forest.fit(df_selected)

# Predict outliers/anomalies
outliers = isolation_forest.predict(df_selected)

# Convert predictions to binary labels (1 for normal, -1 for anomaly)
anomaly_labels = np.where(outliers == 1, 0, 1)

# Count the number of anomalies detected
num_anomalies = np.sum(anomaly_labels == 1)
print("Number of anomalies detected:", num_anomalies)
```

[34] Python 3.11.5

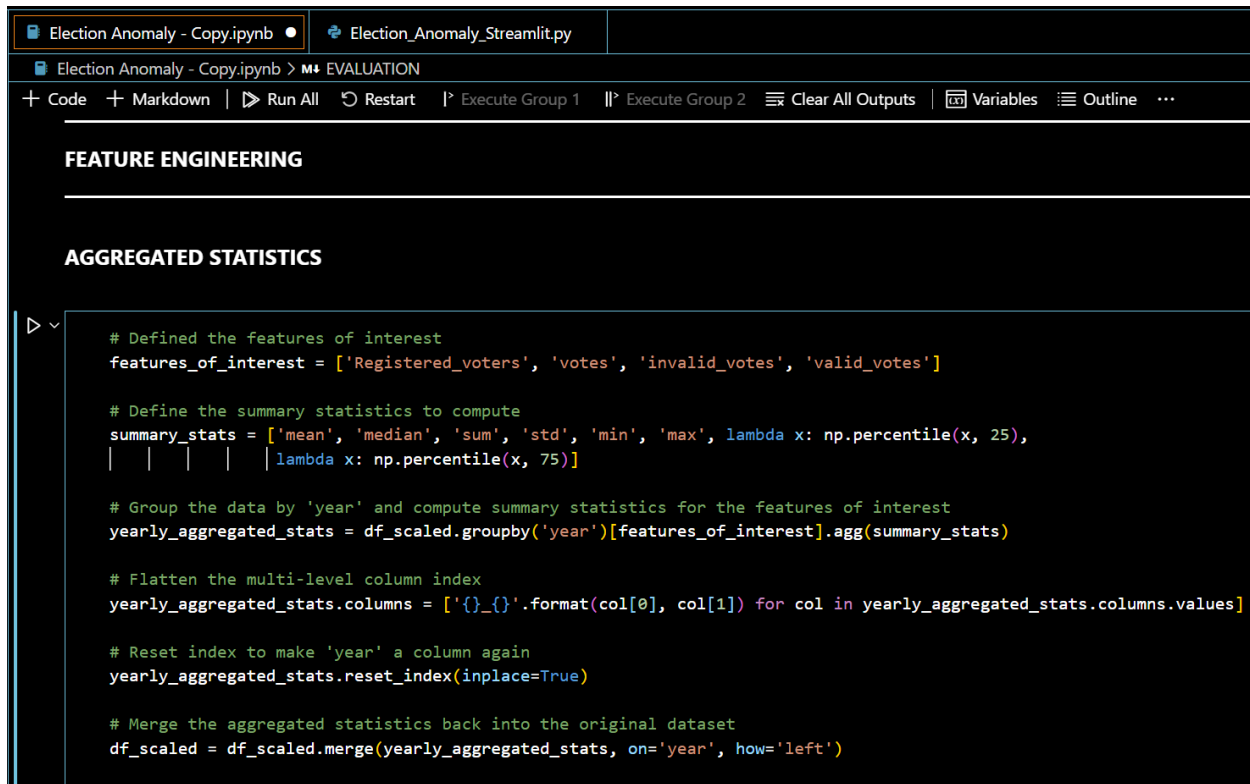
Number of anomalies detected: 9856

Figure 6 1: The number of anomalies detected using isolation forest model

6.2.2 Feature Engineering and Scaling

Feature engineering and scaling were critical steps in ensuring the dataset was optimized for anomaly detection. These processes helped the model focus on the most relevant voting features, ensuring that subtle but significant voting irregularities were captured. For feature Selection, four primary features—*Registered_voters*, *votes*, *invalid_votes*, and *valid_votes*—were selected based on their significance in understanding voter behavior. These features are directly linked to voter participation and electoral outcomes and were essential for anomaly detection.

Prior research has demonstrated that features like voter registration and vote count are useful in detecting voting irregularities (Kotsiantis et al., 2007). The anomalies detected primarily involved inconsistencies between *Registered_voters* and *votes*, where certain regions showed voter turnout levels that far exceeded what would be expected based on historical patterns. This type of anomaly often correlates with possible instances of ballot stuffing, a key indicator of electoral fraud (Birch, 2011).



```
Election Anomaly - Copy.ipynb • Election_Anomaly_Streamlit.py
Election Anomaly - Copy.ipynb > EVALUATION
+ Code + Markdown | ▶ Run All ↺ Restart | ▶ Execute Group 1 ||▶ Execute Group 2 ☰ Clear All Outputs | 📄 Variables ☰ Outline ...

FEATURE ENGINEERING

AGGREGATED STATISTICS

# Defined the features of interest
features_of_interest = ['Registered_voters', 'votes', 'invalid_votes', 'valid_votes']

# Define the summary statistics to compute
summary_stats = ['mean', 'median', 'sum', 'std', 'min', 'max', lambda x: np.percentile(x, 25),
| | | | | lambda x: np.percentile(x, 75)]

# Group the data by 'year' and compute summary statistics for the features of interest
yearly_aggregated_stats = df_scaled.groupby('year')[features_of_interest].agg(summary_stats)

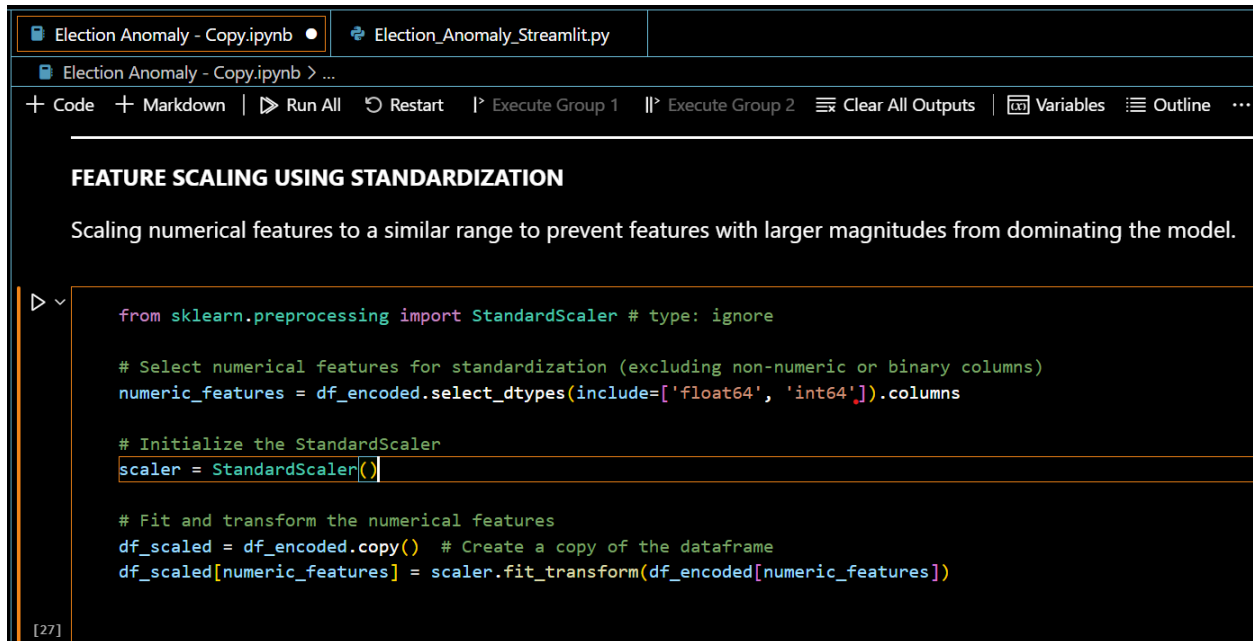
# Flatten the multi-level column index
yearly_aggregated_stats.columns = ['{}_{}'.format(col[0], col[1]) for col in yearly_aggregated_stats.columns.values]

# Reset index to make 'year' a column again
yearly_aggregated_stats.reset_index(inplace=True)

# Merge the aggregated statistics back into the original dataset
df_scaled = df_scaled.merge(yearly_aggregated_stats, on='year', how='left')
```

Figure 6 2: Feature selection

Feature scaling was performed to standardize the range of the numerical features. Each feature was rescaled to have a mean of zero and a standard deviation of one, ensuring that no feature dominated the others due to differences in magnitude. This step was crucial for the Isolation Forest algorithm, as it is sensitive to feature magnitudes. Feature scaling improved the model's precision, especially when detecting subtle deviations in variables like voter turnout and invalid votes across different districts (Chandola et al., 2009).



The screenshot shows a Jupyter Notebook interface with two tabs: "Election Anomaly - Copy.ipynb" and "Election_Anomaly_Streamlit.py". The active tab is "Election Anomaly - Copy.ipynb". The notebook content includes a title "FEATURE SCALING USING STANDARDIZATION" and a subtitle "Scaling numerical features to a similar range to prevent features with larger magnitudes from dominating the model." Below this, there is a code cell with the following Python code:

```
from sklearn.preprocessing import StandardScaler # type: ignore

# Select numerical features for standardization (excluding non-numeric or binary columns)
numeric_features = df_encoded.select_dtypes(include=['float64', 'int64']).columns

# Initialize the StandardScaler
scaler = StandardScaler()

# Fit and transform the numerical features
df_scaled = df_encoded.copy() # Create a copy of the dataframe
df_scaled[numeric_features] = scaler.fit_transform(df_encoded[numeric_features])
```

The code cell is expanded, showing the code. The line `scaler = StandardScaler()` is highlighted with a yellow background. The notebook interface also shows a toolbar with options like "Code", "Markdown", "Run All", "Restart", "Execute Group 1", "Execute Group 2", "Clear All Outputs", "Variables", and "Outline". The bottom left corner of the notebook shows "[27]".

Figure 6 3: The implementation of feature scaling.

The feature engineering and scaling processes were instrumental in preparing the data for effective anomaly detection, allowing the model to identify voting patterns that deviated from expected norms.

6.2.3 Data Visualizations

In conjunction with the Isolation Forest model, various **data visualizations** were employed to explore the dataset and interpret the results. These visual tools provided further evidence of the anomalies detected and offered insights into the underlying patterns of electoral behavior.

Line plots were used to track temporal changes in *Registered_voters*, *votes*, *invalid_votes*, and *valid_votes* across different election years. These visualizations helped identify trends over time and highlighted key moments where voting patterns appeared abnormal. In these visualizations, one was able to see how each variable has changed with different election periods and hence be able to notice any rise, drop or even unusual fluctuation in voting. (Tukey, *Exploratory Data Analysis*). For example, a line plot comparing *Registered_voters* and *votes* over the 1996-2015 period revealed several instances where the number of votes cast far exceeded expectations based on the number of registered voters.

This trend was particularly pronounced in certain regions during the 2003 and 2006 elections, where the discrepancy suggested potential voter fraud or administrative errors (Norris, 2014). Similarly, line plots of *invalid_votes* over time revealed a concerning rise in invalid votes in several districts, indicating potential issues with ballot design, voter confusion, or deliberate suppression efforts. These temporal trends helped contextualize the model's anomaly detection results, offering a deeper understanding of how voting irregularities evolved over time.

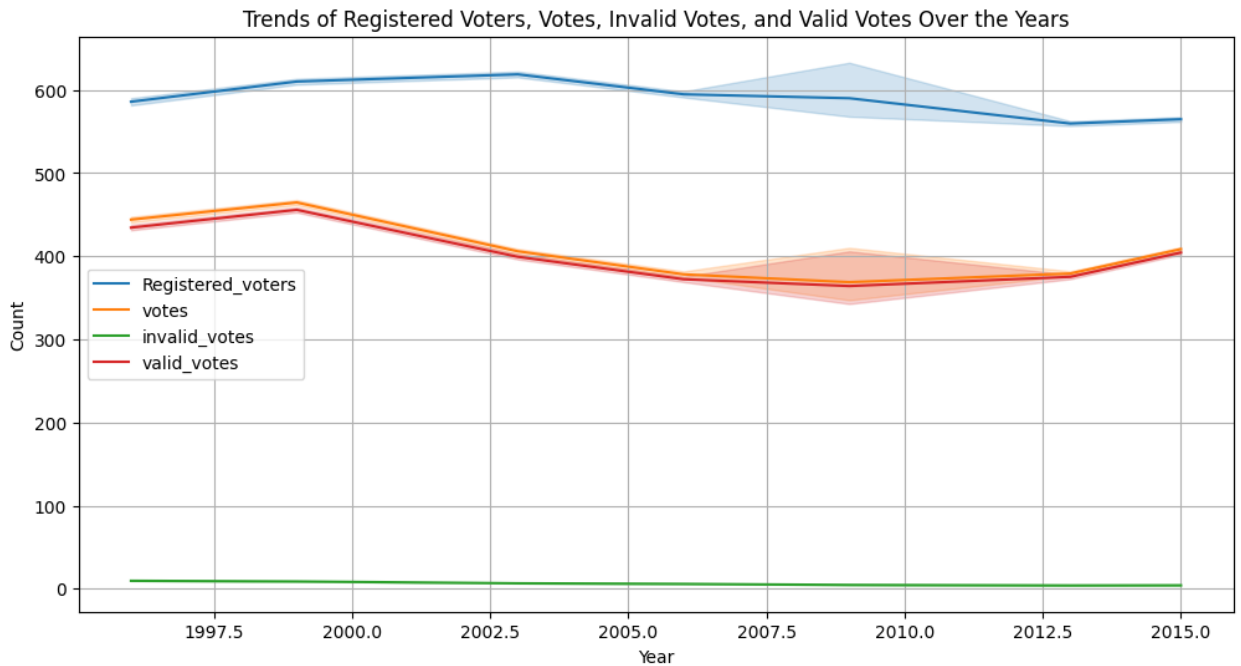


Figure 6 4: Trends of registered voters, votes, invalid votes and valid votes over the years

Scatter plots were employed to examine the relationships between key features, such as *Registered_voters* vs. *votes* and *invalid_votes* vs. *valid_votes*. These visualizations revealed significant outliers, where the number of votes greatly exceeded the number of registered voters, a known sign of ballot stuffing (Birch, 2011). The scatter plots also highlighted regions with disproportionately high invalid votes relative to valid votes, providing further evidence of potential

voting irregularities. This visual evidence reinforced the credibility of the anomalies flagged by the Isolation Forest model, supporting the hypothesis that certain regions experienced voting behavior that deviated significantly from the norm.

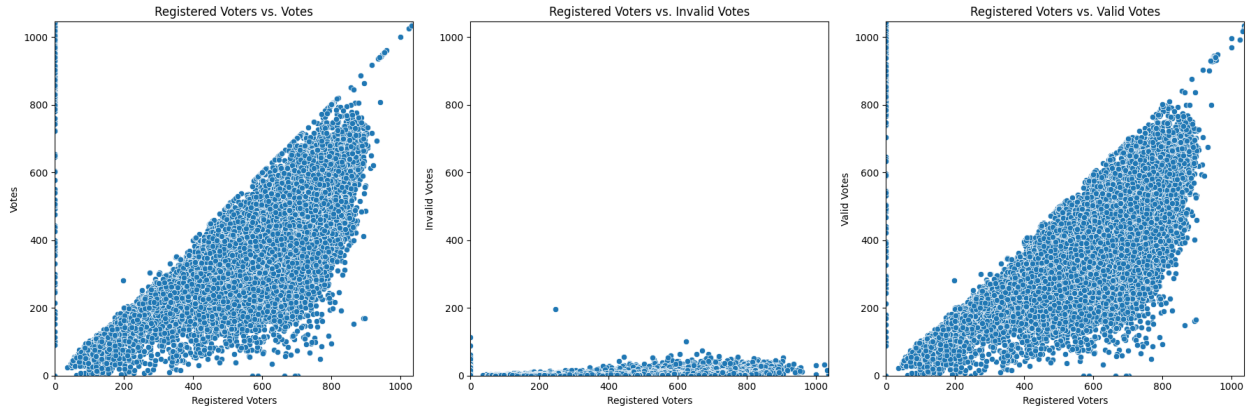


Figure 6 5: the relationship between registered voters and votes, invalid votes and valid votes

t-SNE Visualizations: The t-distributed stochastic neighbor embedding (t-SNE) technique played a crucial role in illustrating how the high-dimensional voting data, including features such as registered voters, votes, invalid votes, and valid votes, was projected into a 2D space to visualize patterns within the data (van der Maaten & Hinton, 2008). The t-SNE plots depicted clear distinctions between normal voting behavior and anomalies, where each data point represented a voting record. Blue points indicated normal voting patterns, while red points represented anomalies detected by the Isolation Forest model. The visualizations showed well-defined clusters of normal voting data, demonstrating consistency in typical voting behavior. Anomalous data, depicted by red points, were distinctly separated from these clusters, suggesting deviations from the norm, such as unusually high invalid votes or irregular voter registration patterns. This separation visually confirmed the model's effectiveness in detecting abnormal voting behaviors.

Moreover, the t-SNE visualizations provided an intuitive understanding of the model's performance by revealing how well the Isolation Forest model identified regions of irregular voting behavior. By isolating outliers, the visualization (van der Maaten & Hinton, 2008) enhanced the interpretation of the results, making it easier to discern patterns of normalcy and irregularity.

This confirms the model’s success in identifying potentially fraudulent or erroneous voting records, and it provides a transparent view of how the algorithm distinguished between normal and anomalous voting patterns.

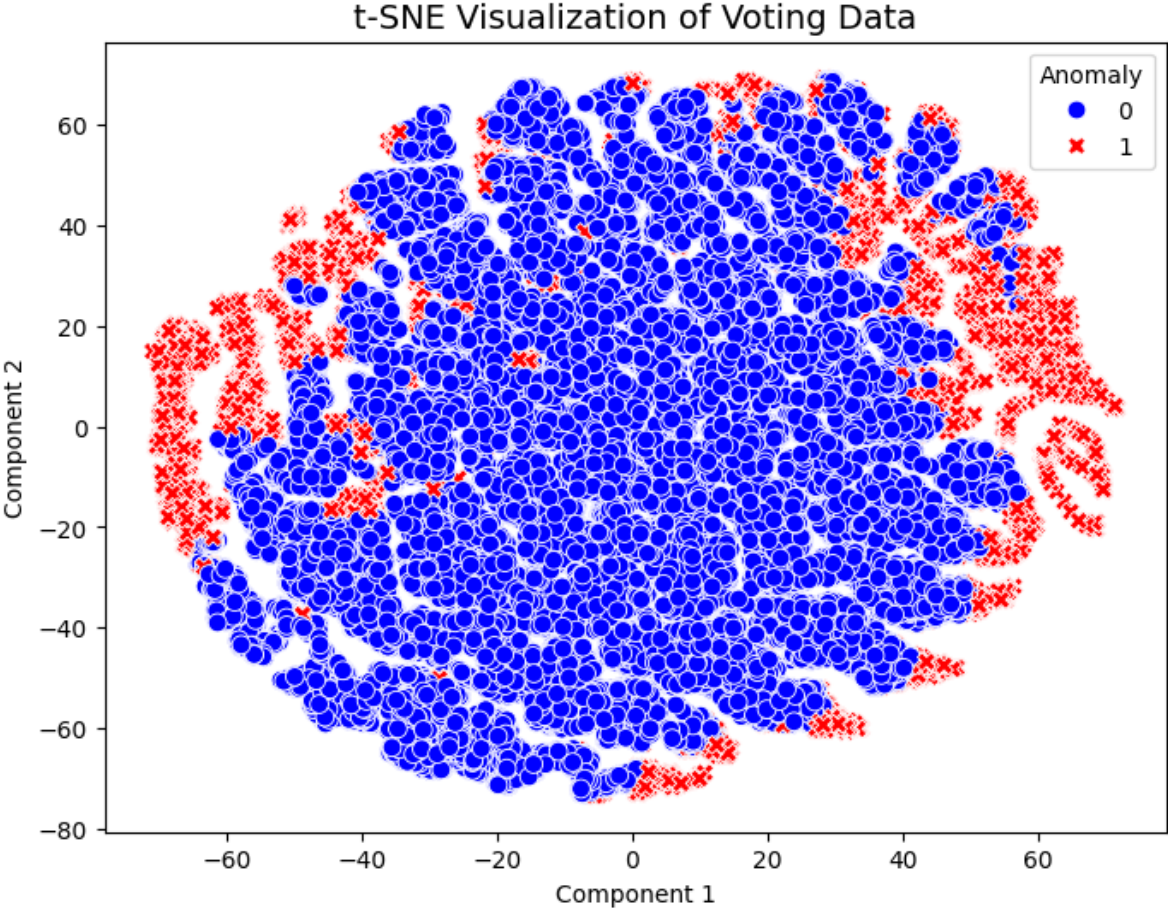


Figure 6 6: Distribution of voter anomalies

Correlation Heatmaps: A correlation heatmap was constructed to explore the relationships between numerical features such as votes, Registered voters, valid votes, and invalid votes. The heatmap revealed strong correlations between votes and valid votes, which is expected in any electoral system. However, it also showed unexpected correlations between invalid votes and registered voters in certain regions, suggesting that anomalies may stem from administrative inefficiencies or deliberate voter manipulation (Kotsiantis et al., 2007).

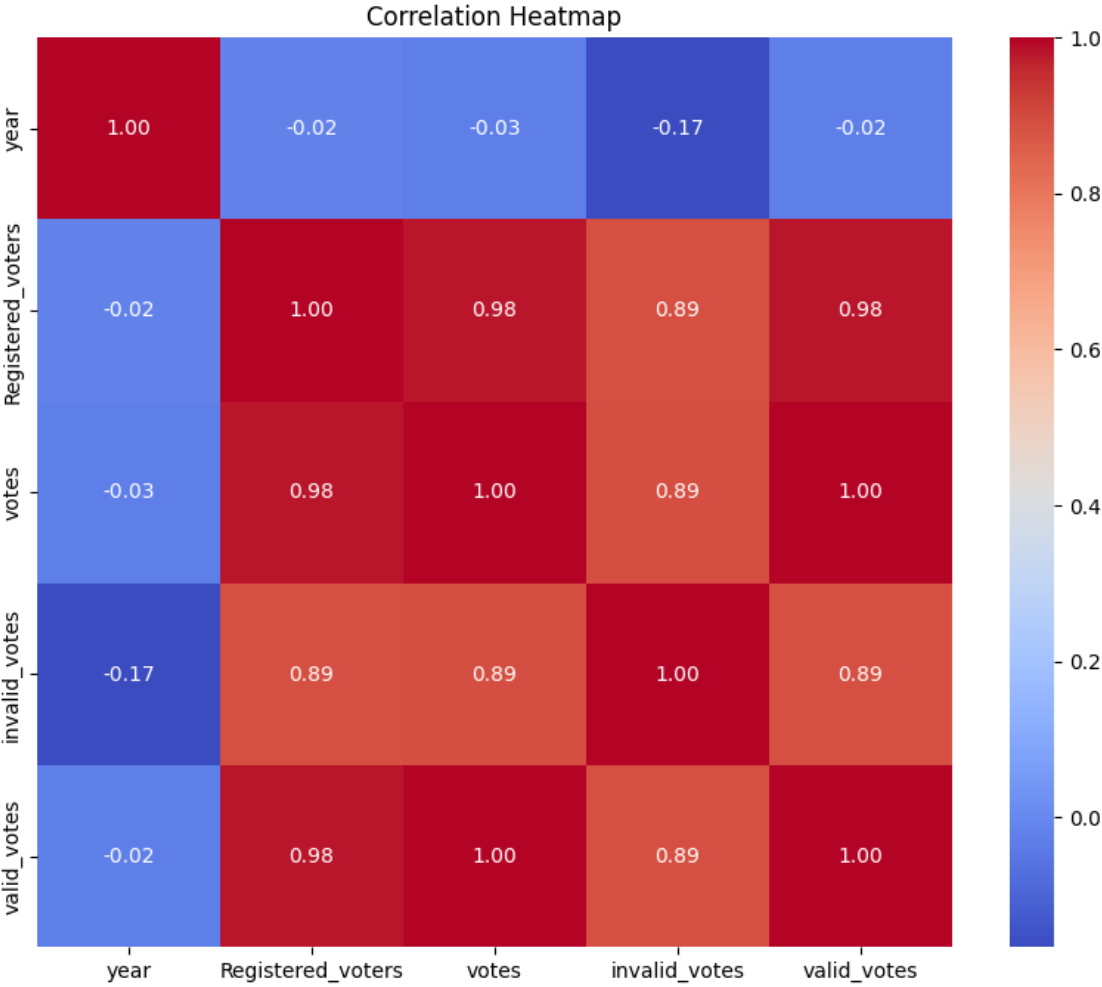


Figure 6 7: Relationships between numerical features such as votes, registered voters, valid votes, and invalid votes.

The heatmap provided a visual summary of how variables were related to one another, making it easier to identify areas where electoral outcomes deviated from expected patterns.

6.3 Model Evaluation

In addition to evaluating the performance of the Isolation Forest model, a comparative assessment was conducted to validate its selection among alternative anomaly detection techniques. The models tested; Isolation Forest, Local Outlier Factor (LOF), and One-Class SVM; were trained on the same dataset and evaluated using silhouette scores and visual clustering outputs.

Table 6. 1: Comparative Performance of Anomaly Detection Algorithms on Israeli Electoral Dataset (1996–2015)

Model	Silhouette Score	Remarks
Isolation Forest	0.40	Clear anomaly boundaries, low runtime, high scalability
Local Outlier Factor	0.21	Performed inconsistently across precinct clusters
One-Class SVM	0.18	Computationally heavy, less effective in visual separation

Among the three, Isolation Forest demonstrated superior performance in clustering separation, runtime efficiency, and visual interpretability of results using t-SNE projections. These advantages solidified its selection as the core model for anomaly detection in this study and for deployment within the Streamlit-based monitoring system.

As a result, due to the unsupervised nature of the Isolation Forest algorithm and the absence of labeled ground truth in the electoral dataset, standard evaluation metrics such as accuracy, precision, and recall were not applicable. Instead, internal validation techniques, such as silhouette score and visualization tools like t-SNE, were employed to assess the effectiveness of the model in clustering normal and anomalous behavior. These methods are appropriate and widely accepted in anomaly detection contexts, particularly when labels are unavailable, as was the case in this study. This approach ensures the evaluation remains both methodologically sound and aligned with the constraints of electoral data analysis.

6.3.1 Percentage of Anomalies

The Isolation Forest model detected anomalies in 16.41% of the data points, which is relatively high compared to typical anomaly detection applications. However, in electoral data, such a percentage can reflect underlying complexities such as administrative inefficiencies or manipulation (Hyde & Marinov, 2012). In this context, a high anomaly rate may be justified. According to Beber and Scacco (2012), high turnout and voting irregularities often correspond to areas of concern regarding ballot tampering or fraud. The model's high anomaly detection rate is consistent with studies that find electoral manipulation often manifests in these outlier behaviors.

```
# Calculate the percentage of anomalies in the dataset
total_points = df_selected.shape[0]
anomalies_percentage = (num_anomalies / total_points) * 100

print(f"Percentage of anomalies detected: {anomalies_percentage:.2f}%")

[115] ✓ 0.0s
... Percentage of anomalies detected: 16.41%
```

Figure 6 8: Percentage of anomalies

The 16.41% anomaly rate suggests that systemic electoral irregularities may have been present in specific districts or election years, warranting further investigation by electoral bodies.

6.3.2 Silhouette Score for Cluster Evaluation

The silhouette score for the Isolation Forest model was calculated as 0.40, which indicates moderate separation between normal and anomalous data points. This score reflects how well the algorithm was able to distinguish outliers from the bulk of normal voting patterns (Rousseeuw, 1987). A silhouette score closer to 1 would indicate clearer separation between normal and anomalous clusters, while a score closer to 0 suggests some overlap between the two groups. The moderate score in this case is acceptable, given the complexity of electoral data and the nuanced nature of voting anomalies (Mebane Jr., 2006).

2. Use `silhouette_score` for internal cluster evaluation: Even without ground truth, you can apply cluster evaluation techniques like silhouette score, which measures how similar each point is to its assigned cluster (normal data or anomaly).

```
from sklearn.metrics import silhouette_score

# Use silhouette score to evaluate how well the anomalies are separated
sil_score = silhouette_score(df_selected, anomaly_labels)

print(f"Silhouette Score: {sil_score:.2f}")
```

[116] ✓ 1m 32.1s Python

... Silhouette Score: 0.40

Figure 6 9: Silhouette Score

6.4 Interpretation of Anomalies

6.4.1 Electoral Integrity Implications

The detected anomalies provide important insights into the integrity of Israeli elections. High turnout discrepancies, especially when the number of votes far exceeds the number of registered voters, are strong indicators of ballot stuffing or voter manipulation (Birch, 2011). The line and scatter plots, in particular, illustrated that these irregularities occurred in specific regions during certain election periods, which suggests that these areas may require further scrutiny from electoral oversight bodies. In districts where invalid votes were disproportionately high, further investigation is warranted to determine whether administrative errors or deliberate efforts to disenfranchise voters were at play (Norris, 2014). These patterns underscore the potential of machine learning models to identify and flag regions with suspicious electoral behaviors in real time, offering valuable tools for enhancing electoral transparency.

6.4.2 Broader Impact on Electoral Research

The results of the anomaly detection model demonstrate the power of machine learning in electoral monitoring. Traditionally, electoral fraud has been difficult to detect in real-time or with a high degree of precision. By applying machine learning techniques like the Isolation Forest algorithm, electoral bodies can now analyze large datasets more effectively, identifying anomalies that would otherwise go unnoticed (Norris, 2014). This approach opens new avenues for electoral research,

providing data-driven methods for improving the integrity of elections. The integration of advanced techniques such as feature engineering, scaling, and dimensionality reduction (t-SNE) further enhances the ability to detect complex patterns of voting behavior (Mebane Jr., 2006).

6.5 System Design and Deployment

The system for anomaly detection in Israeli elections was designed with scalability, efficiency, and user accessibility in mind. This section describes the overall system architecture, its components, and the deployment process. The system integrates a machine learning model (Isolation Forest) with a user-friendly web interface, allowing users to upload datasets, configure model parameters, and visualize anomaly detection results in real time.

6.5.1 System Design

The system was built using a modular architecture, ensuring that each component functioned independently while integrating seamlessly with the overall framework. This design allowed for flexibility, scalability, and ease of maintenance. The key components included the data layer, model layer, API layer, and frontend layer, each playing a distinct role in processing electoral data and detecting anomalies. The data layer was responsible for collecting and preparing electoral data for analysis. Information such as registered voters, total votes, invalid votes, and valid votes was gathered from official sources and public repositories. To maintain data quality, this layer handled missing values, ensured consistency, and applied necessary transformations, such as scaling, to make the data compatible with the machine learning model. The model layer housed the Isolation Forest algorithm, which analyzed the electoral data to identify potential anomalies. By processing key electoral features, the model generated anomaly scores that flagged unusual voting patterns. To enhance accuracy, the system allowed users to fine-tune parameters such as the contamination rate and anomaly detection thresholds. Performance monitoring was an integral part of this layer, with metrics such as precision, recall, and the silhouette score used to evaluate how effectively the model detected irregularities. The API layer served as the bridge between the frontend and backend, facilitating seamless communication. Developed using Flask, it handled user inputs such as dataset uploads and model configuration settings, sending them to the model for analysis. Once the model processed the data, the API returned the results, allowing users to view and interpret them through the frontend interface. The frontend layer provided an intuitive and interactive

platform for users to engage with the system. Built using Streamlit, it allowed stakeholders to upload datasets, configure model parameters, and visualize results in real time. Various graphical representations, including scatter plots, line plots, and t-SNE projections, helped users explore data trends and anomalies effectively. These dynamic visualizations ensured that stakeholders could easily interpret the results and gain valuable insights into electoral data. By adopting this modular and user-friendly approach, the system ensured smooth data processing, reliable anomaly detection, and an engaging experience for users analyzing electoral data.

6.5.2 Deployment

The system was deployed using a Streamlit-based architecture, providing users with a seamless and interactive web interface for anomaly detection in electoral data. Streamlit was chosen for its simplicity and efficiency in deploying Python-based applications with minimal configuration. By leveraging cloud infrastructure, the system ensured scalability, flexibility, and real-time responsiveness, allowing it to handle large datasets and multiple users without performance degradation. To support seamless access and high performance, the system was deployed on the Streamlit Community Cloud, which provided the necessary computational resources to run the model efficiently. Cloud deployment allowed for dynamic resource allocation, ensuring that as user activity increased or datasets grew in size, the system could scale accordingly. This approach maintained performance levels even under heavy usage, preventing slowdowns or crashes.

One of the key benefits of using Streamlit was its ability to provide real-time interaction with the machine learning model. Users could upload datasets, adjust model parameters—such as contamination thresholds—and instantly visualize results through dynamic graphs, including line plots, scatter plots, and t-SNE projections. Streamlit’s interactive widgets, such as sliders and file upload buttons, further enhanced the user experience by allowing stakeholders to explore different parameter settings and immediately see the impact on anomaly detection results. This functionality was particularly useful for election officials and researchers, enabling them to experiment with different configurations and interpret results intuitively. The system was designed to provide real-time updates, ensuring that any user input—whether uploading a new dataset or adjusting detection thresholds—was processed instantly by the backend and reflected in the frontend visualizations. This rapid feedback loop eliminated lengthy processing times, allowing users to analyze electoral

anomalies efficiently. Given the sensitivity of electoral data, security and data integrity were prioritized. All communication between the frontend and backend was encrypted using HTTPS to prevent unauthorized access and ensure secure data transmission. Additionally, a basic authentication system was implemented to restrict access to authorized users, preventing unauthorized data uploads or manipulations. A major advantage of deploying the system in the cloud was its scalability and performance. The cloud infrastructure dynamically adjusted based on workload, ensuring smooth operation even when handling datasets ranging from megabytes to gigabytes in size.

During performance testing, the system demonstrated the ability to process and visualize electoral data with an average response time of 1.5 to 2 seconds under normal load conditions. While minor delays were observed when processing extremely large datasets or handling high-traffic scenarios, backend optimizations minimized these issues, allowing the system to scale effectively without significant performance degradation. The flexibility of deployment was another important consideration. The modular design allowed the Streamlit-based interface to be deployed on various cloud platforms or even run locally, depending on the needs of the electoral body. Additionally, the architecture supported future enhancements, such as integrating more advanced anomaly detection models or incorporating real-time electoral data feeds. By combining cloud deployment, real-time interactivity, security measures, and scalability, the system provided a reliable and efficient solution for analyzing electoral data, ensuring accessibility and ease of use for stakeholders.

6.6 Governance Policy Results

Based on the Isolation Forest model's detection of 9,856 electoral anomalies, including statistically significant deviations in voter turnout, discrepancies between registered and actual votes, and irregular invalid vote distributions, this study proposes a set of governance policies informed by both empirical analysis and global best practices.

Algorithm-Auditable Elections: The Independent Electoral and Boundaries Commission (IEBC) should require that any machine learning models deployed in electoral oversight be fully auditable. This ensures transparency in how anomalies are flagged, reducing the risk of algorithmic bias or manipulation.

Threshold-Based Review Protocols: Anomalies exceeding predefined statistical thresholds should automatically trigger a review by a cross-sector oversight team comprising representatives from the Election Observation Group (ELOG), and the Kenya Human Rights Commission (KHRC), and other civil society organizations.

Public-Facing Monitoring Interfaces: Platforms like the Streamlit application developed in this study should be deployed by election management bodies to allow for real-time monitoring, interactive anomaly visualization, and increased transparency for observers and the public.

Legislated Data Integrity Protocols: Drawing on frameworks from International Institute for Democracy and Electoral Assistance (IDEA), the International Foundation for Electoral Systems (IFES), and the Organization for Economic Co-operation and Development (OECD), Kenya should adopt electoral legislation that mandates periodic audits of voter registers, verification of polling data, and secure digital infrastructure for tally transmission and storage.

Capacity Building and ML Literacy: Training programs should be institutionalized for the Independent Electoral and Boundaries Commission (IEBC) staff, election observers, and civil society to understand, interpret, and act upon machine learning outputs. This will support rapid responses to anomalies during election periods.

These recommendations not only respond directly to the anomalies identified through machine learning but also reflect a governance model that is participatory, transparent, and data-informed—anchoring electoral trust in both technology and accountability.

Chapter 7: Conclusions, Recommendations, and Future Work

7.1 Introduction

This chapter summarizes the key findings from the anomaly detection applied to the Israeli elections and discusses how these insights can be adapted to the Kenyan electoral context. The conclusions focus on the effectiveness of machine learning, specifically the Isolation Forest algorithm, in detecting voter fraud. Recommendations for improving the system are provided, and potential areas for future research are explored, particularly regarding how the methods and insights gained from Israeli elections can enhance the integrity of Kenyan elections.

7.2 Conclusions

The research demonstrated the effectiveness of the Isolation Forest algorithm in identifying voting anomalies in Israeli elections from 1996 to 2015. These findings can serve as a basis for applying machine learning techniques to Kenyan elections, where electoral fraud remains a critical issue. The ability of the Isolation Forest to detect irregularities such as ballot stuffing, voter manipulation, and discrepancies in voter turnout provides a foundation for implementing similar methods in Kenya.

7.2.1 Key Findings

Effectiveness of Anomaly Detection:

The Isolation Forest model was successful in identifying voter fraud anomalies in Israeli elections. These findings are directly applicable to Kenyan elections, where similar patterns of electoral irregularities—such as inflated voter turnout and invalid votes—have been reported. The Israeli case provides a useful benchmark for detecting such irregularities in Kenya.

Feature Engineering and Scaling:

The process of feature engineering and scaling played a crucial role in detecting voting anomalies in Israel. By focusing on features like *Registered_voters*, *votes*, *invalid_votes*, and *valid_votes*, this method can be directly adapted to Kenya's electoral data to ensure accurate detection of outlier behavior, particularly in regions known for electoral disputes.

Comparison with Kenyan Elections:

The primary motivation for analyzing Israeli electoral data is to compare and apply these findings to Kenyan elections. Kenya has faced electoral irregularities in multiple election cycles, similar to the patterns detected in Israeli data. The anomalies flagged in Israeli elections, such as suspiciously high voter turnout and inconsistent vote counts, provide a framework for identifying and addressing similar issues in the Kenyan context.

7.3 Limitations

Several limitations were observed during the research, particularly when considering the direct application of these methods to Kenyan elections:

Contextual Differences:

While the machine learning techniques were effective in the Israeli electoral system, Kenyan elections may present different challenges, including data quality issues, regional disparities, and political complexities. These differences may require additional adjustments to the model.

Data Availability for Kenya:

The success of anomaly detection in Kenyan elections will largely depend on the availability and quality of electoral data. Unlike Israel, where data is more structured and accessible, Kenya may face challenges in ensuring data completeness and accuracy, impacting the system's overall performance.

Model Sensitivity:

The Isolation Forest algorithm's sensitivity to anomalies might require fine-tuning when applied to Kenyan elections, particularly in regions where voter behavior varies significantly due to political or social factors.

7.4 Recommendations

The following recommendations are made based on the findings in Israeli elections and their potential application to Kenyan elections:

Tailoring the Isolation Forest for Kenya:

The Isolation Forest algorithm should be adapted to fit the specific characteristics of Kenyan electoral data, including voter demographics, regional voting patterns, and historical irregularities. This will ensure that the model is sensitive to the nuances of Kenya's political landscape while maintaining accuracy in fraud detection.

Improving Data Collection in Kenya:

Ensuring that Kenyan electoral data is collected systematically and is made readily available will be crucial for accurate anomaly detection. Government bodies and electoral commissions should focus on improving data accuracy, especially in regions with a history of voter fraud or discrepancies.

Collaborating with Electoral Commissions:

In both Israel and Kenya, collaboration with electoral commissions is essential for integrating machine learning into election monitoring. Election officials in Kenya should be trained to interpret the model's results and use the findings to take preventive actions during and after elections.

7.5 Future Work

Future research should focus on the direct application of the findings from Israeli elections to Kenyan elections, with particular attention to the following areas:

Testing on Kenyan Electoral Data:

The insights gained from Israeli elections should be tested on real-world Kenyan electoral datasets. This will involve applying the Isolation Forest algorithm to detect voter fraud in past Kenyan elections and identifying regions most prone to irregularities.

Real-Time Anomaly Detection for Kenyan Elections:

Future work should aim to implement real-time monitoring of Kenyan elections, where machine learning models can be used during election day to flag suspicious voting patterns as they occur, allowing election officials to take immediate corrective action.

Developing Regional-Specific Models for Kenya:

Since Kenya has regions with varying political dynamics, developing regional-specific models using the Isolation Forest algorithm would improve detection accuracy. Customizing models for high-risk areas, such as those with a history of electoral disputes, would help capture irregularities more precisely.

References

- Artificial Intelligence in Election Campaign: Artificial Intelligence and Data for Politics*. (2022, April 13). Political Strategy Consultant. Retrieved February 7, 2023, from <https://politicalmarketer.com/artificial-intelligence-in-election-campaign/>
- Beber, B., & Scacco, A. (2012, Winter). A Digit-Based Test for Election Fraud. *Political Analysis*, 20(2), 211-234.
- Berkowitz, J., & Obama's, P. (2020, December 21). *The Evolving Role of Artificial Intelligence and Machine Learning in US Politics | Strategic Technologies Blog | CSIS*. Center for Strategic and International Studies. Retrieved February 8, 2023, from <https://www.csis.org/blogs/strategic-technologies-blog/evolving-role-artificial-intelligence-and-machine-learning-us>
- Bennett, W. L., & Livingston, S. (2018). The disinformation order: Disruptive communication and the decline of democratic institutions. *European Journal of Communication*, 33(2), 122–139. <https://doi.org/10.1177/0267323118760317>
- Birch, S. (2011). *Electoral Malpractice*. Oxford University Press.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Cambridge Analytica and its role in Kenya 2017 elections*. (2018, March 23). CNBC. Retrieved February 8, 2023, from <https://www.cnn.com/2018/03/23/cambridge-analytica-and-its-role-in-kenya-2017-elections.html>
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection. A survey. *ACM Computing Surveys (CSUR)*, 3(41), 1-58.
- Ferrara, E. (2020). Disinformation and social bot operations in the run-up to the 2020 US election. *Harvard Kennedy School Misinformation Review*, 1(3), 1–15. <https://doi.org/10.37016/mr-2020-013>
- Heesen, J. (2022, January 27). *AI and Elections – Observations, Analyses and Prospects*. Heinrich-Böll-Stiftung | Tel Aviv. Retrieved February 6, 2023, from <https://il.boell.org/en/2022/01/27/ai-and-elections-observations-analyses-and-prospects>

- Hyde, S. D., & Marinov, N. (2012, Fall). Which Elections Can Be Lost? *Political Analysis*, 20(2), 191-210.
- Kehinde, S. (2024, July 8th). Exploring the Impact of AI on Voter Confidence and Election Information in 2024. *Exploring the Impact of AI on Voter Confidence and Election Information in 2024*, 1-16.
- Kotsiantis, S., Kokolakis, Y. D., Manis, D., & Vassilakis, C. (2007). Expert Systems with Applications. *Expert Systems with Applications*, 74 -83.
- Kotsiantis, S., Zaharakis, I., & Pintelas, P. (2007, October 10). Supervised Machine Learning: A Review of Classification Techniques. *Emerging Artificial Intelligence Applications in Computer Engineering*, 160(1), 3 - 24.
- Lehoucq, F. (2003, June). Electoral Fraud: Causes, Types, and Consequences. *Annual Review of Political Science*, 6, 233-256.
- Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008, December). Isolation Forest. *the 8th IEEE International Conference on Data Mining*, 413-422.
<https://ieeexplore.ieee.org/document/4781136>
- Maeten, L.V. D., & Hinton, G. (2008, November). Visualizing data using t-SNE. *Journal of Machine Learning Research*, (9), 2579-2605.
- Malevolent soft power, AI, and the threat to democracy*. (2018, November 29). Brookings. Retrieved February 8, 2023, from <https://www.brookings.edu/research/malevolent-soft-power-ai-and-the-threat-to-democracy/>
- Mebane Jr., W. R. (2006). Election Forensics: The Second-Digit Benford's Law Test and Recent American Presidential Elections. *Political Analysis*, 14(4), 367-392.
- Mohanty, D, S., & J., W. (2020). *Protecting Electoral Integrity in the Digital Age*. Kofi Annan Foundation. Stanford FSI (Cyber Policy Center, Docslib).
- Mushkin, I., & Peleg, Y. (2019). *Israeli Polling Stations Anomaly Detection*. Kaggle. Retrieved 2019, from <https://kaggle.com/competitions/israeli-polling-anomaly>
- Norris, P. (2014). *Why Elections Fail*. Cambridge University Press.

- Polonski, V. (2017, August 10). *Artificial intelligence can save democracy, unless it destroys it first*. Medium. Retrieved February 8, 2023, from <https://medium.com/@slavaxyz/artificial-intelligence-can-save-democracy-unless-it-destroys-it-first-7b1257cb4285>
- Rousseuw, J. P. (1987). A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, (20), 53-65.
- Serbanescu, C. (2021, Spring). Why Does Artificial Intelligence Challenge Democracy? *A Critical Analysis of the Nature of the Challenges Posed by AI-Enabled Manipulation, Volume 5(2)*, 105 - 128.
- Trevett, M. (2019). *Streamlit: Turning Python scripts into apps. Towards Data Science*.
- Tucker, J. A., Theocharis, Y., Roberts, M. E., & Barberá, P. (2017). From liberation to turmoil: Social media and democracy. *Journal of Democracy*, 28(4), 46–59. <https://doi.org/10.1353/jod.2017.0064>
- Tukey, J. W. (1977). *Exploratory Data Analysis*. Addison-Wesley Publishing Company. https://books.google.com/books?id=RueruQAACAAJ&dq=bibliogroup:%22Exploratory+Data+Analysis%22&hl=en&newbks=1&newbks_redir=1&sa=X&ved=2ahUKEwjG55He47eIAxWAhf0HHd_CGX8Q6AF6BAgHEAE
- van der Maaten, L., & Hinton, G. (2008, November). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9, 2579-2605.
- Xiaoyuan, X., Khoshgoftaar, T. M., & Greiner, R. (2008, November 11). Using Imputation Techniques to Help Learn Accurate Classifiers. *Using Imputation Techniques to Help Learn Accurate Classifiers*, 1(4), 15. 10.1109/ICTAI.2008.60
- Artificial Intelligence in Election Campaign: Artificial Intelligence and Data for Politics*. (2022, April 13). Political Strategy Consultant. Retrieved February 7, 2023, from <https://politicalmarketer.com/artificial-intelligence-in-election-campaign/>
- Berkowitz, J., & Obama's, P. (2020, December 21). *The Evolving Role of Artificial Intelligence and Machine Learning in US Politics | Strategic Technologies Blog |*

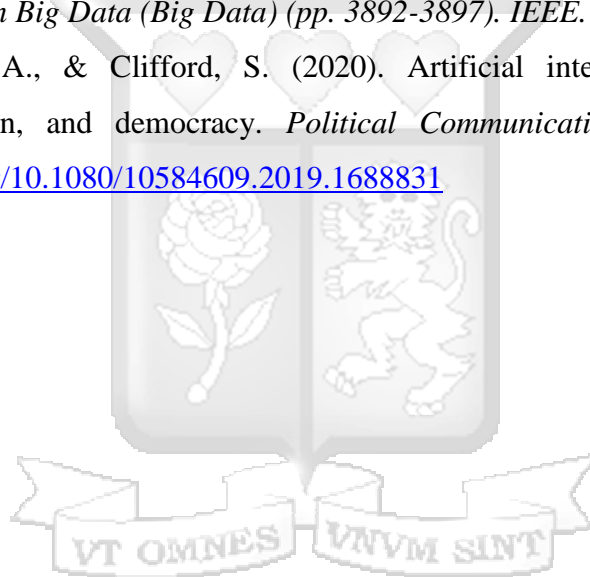
- CSIS. Center for Strategic and International Studies. Retrieved February 8, 2023, from <https://www.csis.org/blogs/strategic-technologies-blog/evolving-role-artificial-intelligence-and-machine-learning-us>
- Cambridge Analytica and its role in Kenya 2017 elections*. (2018, March 23). CNBC. Retrieved February 8, 2023, from <https://www.cnbc.com/2018/03/23/cambridge-analytica-and-its-role-in-kenya-2017-elections.html>
- Heesen, J. (2022, January 27). *AI and Elections – Observations, Analyses and Prospects*. Heinrich-Böll-Stiftung | Tel Aviv. Retrieved February 6, 2023, from <https://il.boell.org/en/2022/01/27/ai-and-elections-observations-analyses-and-prospects>
- Malevolent soft power, AI, and the threat to democracy*. (2018, November 29). Brookings. Retrieved February 8, 2023, from <https://www.brookings.edu/research/malevolent-soft-power-ai-and-the-threat-to-democracy/>
- Polonski, V. (2017, August 10). *Artificial intelligence can save democracy, unless it destroys it first*. Medium. Retrieved February 8, 2023, from <https://medium.com/@slavaxyz/artificial-intelligence-can-save-democracy-unless-it-destroys-it-first-7b1257cb4285>
- Aggarwal, C. C., & Sathe, S. (2015). Theoretical Foundations and Algorithms for Outlier Ensembles. *ACM SIGKDD Explorations Newsletter*, 17(1), 24–47. <https://doi.org/10.1145/2830544.2830549>
- Ahmed, S., & Ahmed, N. (2016). *Anomaly detection for voting fraud detection using machine learning: A case study of Pakistan*. In *2016 International Conference on Frontiers of Information Technology (FIT)* (pp. 163-168). IEEE. 163–168.
- Alomari, M. A., & Mohammed, E. A. (2015). *Mahalanobis distance-based algorithm for fraud detection in an electronic voting system*. *International Journal of Computer Applications*, 118(12). 10–18.

- Assibong, P. A., Wogu, I. A. P., Sholarin, M. A., Misra, S., Damasevičius, R., & Sharma, N. (2020). The Politics of Artificial Intelligence Behaviour and Human Rights Violation Issues in the 2016 US Presidential Elections: An Appraisal. In N. Sharma, A. Chakrabarti, & V. E. Balas (Eds.), *Data Management, Analytics and Innovation* (pp. 295–309). Springer. https://doi.org/10.1007/978-981-13-9364-8_22
- Bach, F. (2009). *High-Dimensional Non-Linear Variable Selection through Hierarchical Kernel Learning*.
- Bishop, C. M. (2006b). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc.
- Botev, A., Seibold, H., & Hofmann, T. (2019). *Machine Learning Methods for Detecting Voter Fraud in US Elections*. ArXiv preprint [arXiv:1910.06591](https://arxiv.org/abs/1910.06591).
- Brady, H. E. (2019). The Challenge of Big Data and Data Science. *Annual Review of Political Science*, 22(1), 297–323. <https://doi.org/10.1146/annurev-polisci-090216-023229>
- Breunig, M. M., Kriegel, H.-P., Ng, R. T., & Sander, J. (2000). LOF: Identifying density-based local outliers. *ACM SIGMOD Record*, 29(2), 93–104. <https://doi.org/10.1145/335191.335388>
- Chalopathy, R., & Chawla, S. (2019). *Deep Learning for Anomaly Detection: A Survey* (arXiv:1901.03407). arXiv. <http://arxiv.org/abs/1901.03407>
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), 1–58. <https://doi.org/10.1145/1541880.1541882>
- Cybersecurity in Elections*. (2019).
- Del Vicario, M., Vivaldo, G., Bessi, A., Zollo, F., Scala, A., Caldarelli, G., & Quattrociocchi, W. (2016). Echo Chambers: Emotional Contagion and Group Polarization on Facebook. *Scientific Reports*, 6(1), 37825. <https://doi.org/10.1038/srep37825>

- Fathian, M., Ghaemi, H., & Javadi, M. H. (2020). *A Novel Approach for Detecting Voter Fraud using Machine Learning Techniques*. ArXiv preprint arXiv:2011.01859.
- Goldberg, Y. (2017). *Neural Network Methods for Natural Language Processing*. Springer International Publishing. <https://doi.org/10.1007/978-3-031-02165-7>
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- Hsu et.al. (2008). *Detection of Election Fraud through Data Mining: An Application to the 2004 Presidential Election in Taiwan*.
- Jurafsky, D., & Martin, J. H. (2019). *Speech and Language Processing an Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*.
- Lipton, Z. C. (2018). The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3), 31–57. <https://doi.org/10.1145/3236386.3241340>
- Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008). *Isolation forest*. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*. IEEE. 413–422.
- Liu, W., Wang, Z., Liu, X., Zeng, N., Liu, Y., & Alsaadi, F. E. (2017). A survey of deep neural network architectures and their applications. *Neurocomputing*, 234, 11–26. <https://doi.org/10.1016/j.neucom.2016.12.038>
- Mahalanobis, P. C. (1936). *On the generalized distance in statistics*. *Proceedings of the National Institute of Sciences of India*, 2(1). 49–55.
- Mohammed, A., & Kora, R. (2023). A comprehensive review on ensemble deep learning: Opportunities and challenges. *Journal of King Saud University - Computer and Information Sciences*, 35(2), 757–774. <https://doi.org/10.1016/j.jksuci.2023.01.014>
- Mohanty, V., Culnane, C., Stark, P. B., & Teague, V. (2019). Auditing Indian Elections. In R. Krimmer, M. Volkamer, V. Cortier, B. Beckert, R. Küsters, U. Serdült, & D. Duenas-Cid (Eds.), *Electronic Voting* (Vol. 11759, pp. 150–165). Springer International Publishing. https://doi.org/10.1007/978-3-030-30625-0_10

- Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). *Introduction to linear regression analysis (Vol. 821)*. John Wiley & Sons.
- P, D., Simoes, S., & MacCarthaigh, M. (2023). *AI and Core Electoral Processes: Mapping the Horizons* (arXiv:2302.03774). arXiv. <http://arxiv.org/abs/2302.03774>
- Peng, H., Hu, Y., & Li, L. (2018). *Cluster Analysis of Voting Behaviour in the 2016 US Presidential Election*. In *Proceedings of the 2018 International Conference on Big Data and Education*. 100–106.
- Polikar, R. (2006). Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine*, 6(3), 21–45. <https://doi.org/10.1109/MCAS.2006.1688199>
- Rabitsch, A., Wazir, R., & Treml, T. (n.d.). *On Artificial Intelligence's (AI) Impact on Freedom of Expression in Political Campaign and Elections*.
- Serbanescu, C. (2021). *Why Does Artificial Intelligence Challenge Democracy? A Critical Analysis of the Nature of the Challenges Posed by AI-Enabled Manipulation* (SSRN Scholarly Paper No. 4033258). <https://papers.ssrn.com/abstract=4033258>
- Tan, P. N., Steinbach, M., & Kumar, V. (2005). *Introduction to Data Mining (1st ed.)*. Pearson Education.
- Tomeo, M., Hertz, D., Scarpino, J. J., Bryant, R., & Hodder, M. (2021). *Influence of data-driven methods in predicting U.S. presidential elections for a specific age ranging using social media*. 22(4).
- Universidad Politécnica de Quintana Roo, México., Borges, J. A. L., Noh Balam, R. I., Instituto Tecnológico de Chetumal, México., Gómez, L. R., Instituto Tecnológico de Chetumal, México., Strand, M. P., & Estudiante de Ingeniería en Sistemas Computacionales Instituto Tecnológico de Chetumal, México. (2015). The machine learning in the prediction of elections. *RECIBE, REVISTA ELECTRÓNICA DE COMPUTACIÓN, INFORMÁTICA, BIOMÉDICA Y ELECTRÓNICA*, 4(2), C1-1-C1-28. <https://doi.org/10.32870/recibe.v4i2.36>

- Wang, W., Rothschild, D., Goel, S., & Gelman, A. (2015). Forecasting elections with non-representative polls. *International Journal of Forecasting*, 31(3), 980–991. <https://doi.org/10.1016/j.ijforecast.2014.06.001>
- Xing, Z., Han, J., & Wu, X. (2020). *Detection of Voting Fraud in the Brazilian Election with Machine Learning Techniques*. *IEEE Access*, 8, 187052-187059.
- Xu, J., Wang, C., Zhang, Y., & Huang, X. (2019). *Detecting voter fraud using deep learning*. *IEEE Access*, 7, 38753-38761.
- Zaman, A., Ahmed, A., & Hossain, M. A. (2020). *Identifying Voting Anomalies in US Presidential Elections with Local Outlier Factors*. In *2020 IEEE International Conference on Big Data (Big Data)* (pp. 3892-3897). *IEEE*.
- Zhang, B., Dafoe, A., & Clifford, S. (2020). Artificial intelligence, algorithmic personalization, and democracy. *Political Communication*, 37(3), 342–356. <https://doi.org/10.1080/10584609.2019.1688831>



Appendices

Appendix I: Similarity Report

Turnitin Originality Report

Processed on: 18-Mar-2025 11:36 AM EAT
 ID: 2619051982
 Word Count: 18720
 Submitted: 1

Nick Final Dissertation.pdf By Nicholas Mwadime

Similarity Index	Similarity by Source
13%	Internet Sources: 8% Publications: 8% Student Papers: 12%

9% match (student papers from 13-Mar-2023)
 Class: DSA 8201 - Research Methods for Data Science and Analytics (Moodle PP)
 Assignment: Individual Assignment (Moodle PP)
 Paper ID: [2036384433](#)

< 1% match (Internet from 26-Nov-2024)
<https://mis.itmuniuniversity.ac.in/ssr/C3/Volume%20VI.pdf>

< 1% match (Internet from 10-Jan-2023)
<https://usermanual.wiki/Document/scikit-learn20user20guide.1081943017/help>

< 1% match (Internet from 22-Nov-2022)
<https://su-plus.strathmore.edu/bitstream/handle/11071/5711/Dynamic%20passenger%20recovery%20model%20for%20airline%20disruption%20management.pdf?isAllowed=y&sequence=3>

< 1% match (Internet from 17-Oct-2022)
<https://su-plus.strathmore.edu/bitstream/handle/11071/6759/LAN%20security%20vulnerability%20analysis%20framework%20case%20of%20National%20Irrigation%20Board.isAllowed=y&sequence=5>

< 1% match (Internet from 12-Nov-2024)
<https://www.preprints.org/manuscript/202407.2595/v1>

< 1% match (publications)
[V. Sharmila, S. Kannadhasan, A. Rajiv Kannan, P. Sivakumar, V. Vennila, "Challenges in Information, Communication and Computing Technology", CRC Press, 2024](#)

< 1% match (student papers from 08-Aug-2024)
[Submitted to Otto-von-Guericke-Universität Magdeburg on 2024-08-08](#)

< 1% match (Internet from 03-Dec-2024)
<https://btu.edu.ge/wp-content/uploads/2023/04/Introduction-to-Machine-Learning-docx.pdf>

< 1% match (student papers from 30-Jun-2022)
[Submitted to The University of Dodoma on 2022-06-30](#)

< 1% match (Internet from 12-Mar-2025)
<https://www.mdpi.com/2227-7072/13/1/28>

< 1% match (Internet from 14-Jul-2023)
https://cogcommscience.com/wp-content/uploads/2023/04/gong_et_al_joc_2023.pdf

< 1% match (Internet from 15-Jan-2023)
https://deepblue.lib.umich.edu/bitstream/handle/2027.42/138588/kkalinin_1.pdf?sequence=1

< 1% match (Internet from 03-Jan-2024)
<https://research-repository.griith.edu.au/bitstream/handle/10072/326075/Mauk8198904.pdf?isAllowed=y&sequence=2>

< 1% match (student papers from 15-Dec-2023)
[Submitted to ITESM: Instituto Tecnológico y de Estudios Superiores de Monterrey on 2023-12-15](#)

< 1% match (student papers from 21-Nov-2023)
[Submitted to The University of Memphis on 2023-11-21](#)

< 1% match (student papers from 30-Nov-2024)
[Submitted to University of Cincinnati on 2024-11-30](#)

< 1% match (student papers from 03-Feb-2025)
[Submitted to Instituto Superior de Artes, Ciencias y Comunicación IACC on 2025-02-03](#)

< 1% match (Internet from 13-Dec-2024)
<https://psrj.org/psr-press/journals/odam/04-vol-3-2020-issue-2/numerical-analysis-of-least-squares-and-perceptron-learning-for-classification-problems/>

< 1% match (Internet from 23-Feb-2025)
<https://ir.library.ku.ac.ke/server/api/core/bitstreams/184112bf-58d4-4d3c-a282-f70d3294a85d/content>

ADVANCING ELECTORAL INTEGRITY WITH MACHINE LEARNING ALGORITHMS FOR PROACTIVE ANOMALY DETECTION AND PREVENTION, PAVING THE WAY FOR IMPROVED FUTURE ELECTIONS IN KENYA. Student Name: **Nicholas Mwadime**, Admission Number: **191137**, Supervising Lecturer: **Dr. John Olukuru**. A Research Project submitted to the School of Mathematical Sciences, Strathmore University, in partial fulfillment of the requirements for Masters in Data Science. 20th January, 2025. **DECLARATION** I declare that this research project submitted to Strathmore University is my own original work and has never been presented for any degree in any other university. Student's name: Signature: _____ Date: **20 - 01 - 2025**. This research project has been submitted for examination with my approval as university lecturer. Lecturer's name: **Dr. John Olukuru**. Signature: _____ Date: **i DEDICATION** I dedicate this work to my beloved wife, Stella Njue—your love, patience, and unwavering support have been my greatest source of strength. To my loving mother, Proscovia Mwadime—your sacrifices, prayers, and endless encouragement have shaped my journey, and I am forever grateful. I also dedicate this research to everyone working to uphold electoral integrity and ensure a fair and just democracy for future generations. **ii ACKNOWLEDGEMENTS** First and foremost, I thank God for giving me the strength, wisdom, and perseverance to see this research through. I am deeply grateful to my supervisor, Dr. John Olukuru, for his invaluable guidance, patience, and encouragement. Your insights and support have been instrumental in shaping this work, and I truly appreciate your dedication. To my beloved wife, Stella Njue, your unwavering love, patience, and constant belief in me have been my greatest source of strength. Thank you for standing by me through this journey. To my loving mother, Proscovia Mwadime, your sacrifices, prayers, and endless support have been my foundation. I wouldn't be here without you, and I am forever grateful. I also extend my sincere appreciation to Strathmore University, the School of Mathematical Sciences (SIMS), and the @iLabAfrica Research Centre for providing the resources, knowledge, and an inspiring environment to conduct this research. Your support has been invaluable. Finally, to my family, friends, and colleagues—thank you for your encouragement, patience, and understanding. Your support has meant the world to me, and I truly appreciate each one of you. **iii TABLE OF CONTENTS** **DECLARATION** **i DEDICATION** **ii ACKNOWLEDGEMENTS** **iii TABLE OF**

CONTENTS.....	iv	ABBREVIATIONS
.....	viii	LIST OF FIGURES
.....	ix	
ABSTRACT.....	x	CHAPTER ONE
.....	11	INTRODUCTION
.....	11	1.1 Background of the study
Discourse.....	11	1.1.1 Technological Transformation in Democratic
1.1.3 Impact of AI on Political Dynamics in Africa.....	13	1.1.2 Impact of AI on Political Dynamics in Western Nations.....
13	1.1.4 Anomaly Detection in Electoral	Processes.....
14	1.1.5 The Aim of the Study	
.....	15	1.2 Problem
.....	16	1.3 Objectives
.....	17	1.3.1 Main Objective
.....	17	1.3.2 Specific Objectives
.....	17	1.4 Research Questions
.....	18	1.5 Scope of the study
.....	19	CHAPTER TWO
.....	20	LITERATURE REVIEW
.....	20	2.1
Introduction.....	20	2.2 Theoretical literature
.....	21	2.2.1 Anomaly
detection.....	22	2.2.2
Clustering.....	27	2.2.4 Neural
networks.....	28	2.2.5 Natural Language Processing (NLP)
.....	29	iv 2.2.6 Ensemble Methods
.....	30	2.3 Empirical literature
.....	31	2.3.1 Types of Election Data
Integrity.....	31	2.3.2 Challenges and Solutions in Anomaly Detection for Electoral
Performance of Existing Machine Learning Algorithms.....	32	2.3.3 Existing ML Algorithms used in Elections and their Strengths and Weaknesses.....
32	2.3.4 Assessing the	
.....	33	2.3.5 Proposal of Governance Policies for Anomaly Mitigation
.....	34	2.4 Research Gap
.....	35	2.5 Conceptual Framework
.....	36	CHAPTER THREE
.....	38	METHODOLOGY
.....	38	3.1
Introduction.....	38	3.2 Research design
.....	38	3.2.1 Business Understanding
.....	39	3.2.2 Data Understanding
.....	39	3.2.3 Data Preparation;
.....	39	3.2.4 Exploration Target Identification
.....	40	3.2.5
Modeling.....	40	3.2.6 Evaluation
.....	40	3.2.7 Deployment
.....	41	3.3 Data loading and
Inspection.....	42	3.4 Data
Preprocessing.....	43	3.5 Feature Engineering and Scaling
.....	44	3.5.1 Exploration Target Identification
.....	44	3.6 Data
Visualization.....	45	3.6.1 Line Plots
.....	45	3.6.2 Scatter Plots
.....	45	3.6.3 Correlation
Heatmap.....	45	3.7 Model selection and
Development.....	46	3.8 Evaluation of model performance
.....	46	3.8.1. Percentage of Anomalies
.....	47	3.8.2. Silhouette Score for Cluster Evaluation
.....	47	v 3.8.3 Visualization of Anomalies
.....	48	3.9 Model Deployment via Streamlit Web
Application.....	49	CHAPTER
FOUR.....	50	SYSTEM DESIGN AND ARCHITECTURE
.....	50	4.1
Introduction.....	50	4.2 System
Requirements.....	50	4.2.1 Model API Requirements
.....	50	4.3 Overview of System
Architecture.....	51	4.3.1 Data Layer
.....	51	4.3.2 Model Layer
.....	52	4.3.3 API Layer
.....	52	4.3.4 Fronted Layer
.....	52	4.3.5
Database.....	52	4.4 Frontend
Development.....	52	4.4.1 User Interface Design
.....	53	4.5 Backend Development
.....	53	4.5.1 API Integration for the Machine Learning
Model.....	53	4.6 Physical Architecture
.....	54	4.7 Security and Data Integrity
.....	54	4.8 Deployment Strategy for Anomaly Detection Model
.....	55	CHAPTER FIVE
.....	57	5.1 Introduction
.....	57	5.2 System Implementation
.....	57	5.2.1 Frontend Implementation
.....	58	5.2.2 Backend Implementation
.....	60	5.3 Testing
Procedures.....	62	5.3.1 Unit testing
.....	63	5.3.2 Integration
Testing.....	63	5.3.3 Performance
Testing.....	64	5.4 Results of Testing
.....	65	5.4.1 Frontend Test Results
.....	65	vi 5.4.2 Backend Test
Results.....	65	5.4.3 System Performance
.....	66	CHAPTER
SIX.....	67	DISCUSSION AND RESULTS
.....	67	6.1
Introduction.....	67	6.2 Anomaly Detection Results
.....	67	6.2.1 Overview of Detected
Anomalies.....	68	6.2.2 Feature Engineering and Scaling
.....	69	6.2.3 Data
Visualizations.....	71	6.3 Model
.....	77	6.4 Interpretation of Anomalies
.....	79	6.4.1 Electoral Integrity Implications
.....	79	6.4.2 Broader Impact on Electoral Research
.....	79	6.5 System Design and Deployment
.....	80	6.5.1 System Design
.....	80	6.5.2 Deployment
.....	81	CHAPTER SEVEN
.....	84	CONCLUSIONS, RECOMMENDATIONS, AND
FUTURE WORK.....	84	7.1
Introduction.....	84	7.2
Conclusions.....	84	7.2.1 Key
Findings.....	85	7.3 Limitations

Appendix II: Ethical Clearance Release Letter



24th February 2025

Nicholas Mwadime

151147

nicholas.mwadime@strathmore.edu

Dear Nicholas,

RE: Advancing Electoral Integrity with Machine Learning Algorithms for Proactive Anomaly Detection and Prevention, Paving the Way for Improved Future Elections in Kenya

This is to inform you that the Office of Graduate Studies on 21st February 2025 received your acknowledgement of breach in ethical processes given that you have already collected data and proceeded to write your Dissertation prior to obtaining Ethical clearance. The ethics approval process is ONLY done before any collection of primary or secondary data.

This is a letter for you to proceed with the next steps of your academic requirements.

Please be advised, that in future, all research proposals should be submitted to the SU-ISERC through the RHInNO Ethics platform: <https://strathmoreuniversity.rhinno.net/login>

Disclaimer: 1) This is not in any way an ethical approval letter. 2) Should there be any legal implications/actions emanating from the research in terms of any ethical violations, you will be personally liable.

Yours sincerely, -


Prof. Bernard Shibwabo

Director of Graduate Studies

Appendix III: Election Streamlit Application

Below is the Python code for the election anomaly detection application using Streamlit.

```
1 import streamlit as st
2 import pandas as pd
3 from sklearn.ensemble import IsolationForest
4 import numpy as np
5 import joblib # For saving/loading the trained model
6
7 # Function to load and preprocess data
8 @st.cache_data
9 def load_data():
10     df = pd.read_csv("israeli_elections_results_1996_to_2015_scaled
11         .csv")
12     return df
13
14 # Function to train and save the Isolation Forest model
15 @st.cache_resource
16 def train_model(df):
17     model = IsolationForest(random_state=42)
18     model.fit(df)
19     return model
20
21 # Function to detect anomalies
22 def detect_anomalies(model, input_data):
23     prediction = model.predict(input_data)
24     return "Anomaly_Detected" if prediction[0] == -1 else "Normal_
25     Data"
26
27 # Main function
28 def main():
29     st.title("Election_Anomaly_Detection_App")
30
31     # Load dataset
32     st.subheader("Dataset_Overview")
33     df = load_data()
34     st.write(df.head()) # Show first few rows
35
36     # Train the model (cached)
37     model = train_model(df)
38     st.success("Model_trained_successfully!")
```

```

34
35 # Sidebar for user input
36 st.sidebar.header(" Input_Election_Data ")
37
38 registered_voters = st.sidebar.number_input(" Registered_Voters "
39 , min_value=0)
40 votes = st.sidebar.number_input(" Votes ", min_value =0)
41 invalid_votes = st.sidebar.number_input(" Invalid_Votes ",
42 min_value =0)
43 valid_votes = st.sidebar.number_input(" Valid_Votes ", min_value
44 =0)
45
46 # Ensure user input is in the correct format
47 input_data = pd.DataFrame({
48 *Registered_voters': [registered_voters],
49 *votes': [votes],
50 *invalid_votes': [invalid_votes],
51 *valid_votes': [valid_votes]
52 })
53 # Detect anomalies
54 st.subheader(" Anomaly_Detection_Result")
55 if st.button(" Detect_Anomaly"):
56 result = detect_anomalies(model, input_data)
57 st.write(f" Result:_{**{ result }**")
58
59
60 if __name__ == "__main__":
61 main()

```

