

**SPAM MAIL DETECTION**

**PROJECT REPORT**

*Submitted by*

**SIDDHARTH VATS [RegNo: RA2111003010606]  
SHUBHAM SHARMA [RegNo: RA2111003010611]  
ATHARV MISHRA [RegNo: RA2111003010621]**

*Under the Guidance of*

**Ms. P. Nithyakani**

Assistant Professor, Computing Technologies

*In partial satisfaction of the requirements for the degree of*

**BACHELOR OF TECHNOLOGY**

**in**

**COMPUTER SCIENCE ENGINEERING**

**with specialization in Computer Science and Engineering**



**SCHOOL OF COMPUTING**

**COLLEGE OF ENGINEERING AND TECHNOLOGY**

**SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

**KATTANKULATHUR - 603203**

**May 2024**

**SRM INSTITUTION OF SCIENCE AND TECHNOLOGY  
KATTANKULATHUR-603203**

**BONAFIDE CERTIFICATE**

Certified that this lab report titled Spam Mail Detection is the bonafide work done by Siddharth Vats (RA2111003010606) Shubham Sharma (RA2111003010611), Atharv Mishra (RA2111003010621) who carried out the lab exercises under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form part of any other work.

**SIGNATURE**

**Ms. P. Nithyakani**

Assistant Professor

Computing Technologies

**SIGNATURE**

**Dr. Pushpalatha M.**

Head of the Department

Computing Technologies

## ABSTRACT

Nowadays, emails are used in almost every field, from business to education. Emails have two subcategories, i.e., ham and spam. Email spam, also called junk emails or unwanted emails, is a type of email that can be used to harm any user by wasting his/her time, computing resources, and stealing valuable information. The ratio of spam emails is increasing rapidly day by day. Spam detection and filtration are significant and enormous problems for email and IoT service providers nowadays. Among all the techniques developed for detecting and preventing spam, filtering email is one of the most essential and prominent approaches. Several machine learning and deep learning techniques have been used for this purpose, i.e., Naïve Bayes, decision trees, neural networks, and random forest. This paper surveys the machine learning techniques used for spam filtering techniques used in email and IoT platforms by classifying them into suitable categories. A comprehensive comparison of these techniques is also made based on accuracy, precision, recall, etc. In the end, comprehensive insights and future research directions are also discussed. The pervasive use of email across various domains has led to an unprecedented influx of spam, posing significant challenges to users and service providers alike. Spam emails, characterized by their intrusive and often malicious nature, not only disrupt user productivity but also pose serious security threats, including the potential for data theft and resource exploitation. As the ratio of spam emails continues to escalate rapidly, the need for robust spam detection and filtration mechanisms has become paramount. In response to this pressing challenge, researchers and practitioners have turned to machine learning and deep learning techniques as powerful tools for spam filtering. Furthermore, this paper presents a thorough comparative analysis of these machine learning techniques based on key performance metrics such as accuracy, precision, and recall. By evaluating the efficacy of each approach, this study offers valuable insights into the relative performance of different spam filtering methodologies, thereby guiding the selection of appropriate techniques for specific use cases and deployment scenarios.

# **TABLE OF CONTENT**

<b>CHAPTER NO</b>	<b>TITLE</b>	<b>PAGE NO</b>
	<b>ABSTRACT</b>	<b>ii</b>
	<b>LIST OF FIGURES</b>	<b>iv</b>
	<b>LIST OF ABBREVIATIONS</b>	<b>v</b>
<b>1</b>	<b>INTRODUCTION</b>	<b>5</b>
<b>2</b>	<b>LITERATURE SURVEY</b>	<b>7</b>
<b>3</b>	<b>SYSTEM ARCHITECTURE AND DESIGN</b>	<b>9</b>
<b>4</b>	<b>METHODOLOGY</b>	<b>12</b>
<b>5</b>	<b>CODING AND TESTING</b>	<b>14</b>
<b>6</b>	<b>SCREENSHOT AND RESULT</b>	<b>15</b>
<b>7</b>	<b>CONCLUSION AND FUTURE ENHANCEMENT</b>	<b>19</b>

## **LIST OF FIGURES**

1. System Design
2. Architecture of Logistic Regression
3. Architecture of Supervised Learning
4. Scatter Plot
5. Accuracy of the model

## **ABBREVIATION**

1. LRMP - Logistic Regression Mini-Project
2. LRAI - Logistic Regression AI
3. LOGRES - Logistic Regression System
4. LRMAS - Logistic Regression Machine Learning System
5. LOGIT - Logistic Regression Implementation and Testing
6. LRP - Logistic Regression Project
7. LRCA - Logistic Regression Classification Application

# **CHAPTER 1**

## **INTRODUCTION**

In the digital age, email has emerged as a cornerstone of communication, permeating virtually every facet of modern life, from business transactions to educational endeavors. Its unparalleled convenience and ubiquity have revolutionized the way individuals and organizations interact, facilitating seamless exchange of information across geographical boundaries and time zones. However, amidst the proliferation of legitimate communication, lurks a pervasive and insidious threat – email spam. Email spam, also known as junk or unwanted emails, constitutes a formidable challenge in today's interconnected world. Characterized by its intrusive and often deceptive nature, spam inundates users' inboxes with unsolicited messages, ranging from promotional offers to phishing scams and malware-laden attachments. Beyond merely disrupting the flow of communication, spam poses multifaceted risks to users and organizations alike. One of the most pressing concerns associated with email spam is its potential to undermine user productivity. As individuals are forced to sift through a deluge of irrelevant messages, valuable time and attention are squandered, detracting from more meaningful tasks and responsibilities. Moreover, the sheer volume of spam emails can overwhelm email servers and computing resources, leading to performance degradation and operational inefficiencies. Furthermore, email spam represents a significant security threat, with the potential to inflict substantial harm on unsuspecting recipients. Malicious actors often exploit spam as a vehicle for phishing attacks, seeking to deceive users into divulging sensitive information such as login credentials or financial details. Additionally, spam emails may contain malware payloads, posing risks of data breaches, identity theft, and system compromise. The escalating ratio of spam emails exacerbates these challenges, placing unprecedented strains on email and Internet of Things (IoT) service providers. As spammers employ increasingly sophisticated techniques to evade detection, traditional rule-based approaches have proven inadequate in stemming the tide of unwanted messages. In response, the adoption of machine learning and deep learning techniques has emerged as a promising strategy for enhancing spam detection and filtration.

Machine learning techniques leverage algorithms and statistical models to discern patterns and anomalies within email data, enabling automated classification of messages as either spam or legitimate (referred to as "ham"). By analyzing features such as message content, sender information, and metadata, machine learning algorithms can effectively differentiate between benign and malicious emails, thereby fortifying email security.

Against this backdrop, this paper endeavors to survey the landscape of machine learning techniques utilized for spam filtering in email and IoT platforms. By categorizing these techniques into distinct categories and conducting a comprehensive comparative analysis, this study aims to elucidate the relative strengths and weaknesses of different spam filtering methodologies. Furthermore, by discussing insights and future research directions, this paper seeks to inform ongoing efforts to combat the scourge of email spam and safeguard the integrity of digital communication channels.



## **CHAPTER 2**

### **LITERATURE SURVEY**

The literature survey for a machine learning project on spam mail detection using logistic regression involves understanding the current state of research in this field. The following studies provide an overview of the techniques, challenges, and datasets used in spam text detection and classification.

#### **1. A systematic literature review on spam content detection and classification**

presents a comprehensive review of existing research on spam classification using machine learning. The study compares the accuracy of existing spam text detection systems and discusses the challenges in spam detection. The authors also highlight the importance of feature extraction techniques and spam text classification techniques.

#### **2. Email Spam Detection Using Machine Learning**

discusses the use of machine learning techniques for email spam detection. The study provides a detailed overview of various algorithms used for email spam detection, including Naive Bayes, Support Vector Machines (SVM), Decision Trees, Random Forests, and Neural Networks. The authors also emphasize the importance of quality and diversity of training data for successful email spam detection.

#### **3. Machine-Learning-Based Spam Mail Detector**

presents a comprehensive examination of machine learning-based spam filtering approaches. The study discusses the strengths and weaknesses of prevalent machine learning algorithms and highlights the potential of deep learning, including deep adversarial learning, in addressing the intricate issue of spam emails.

#### **4. A Comprehensive Survey for Intelligent Spam Email Detection**

describes a focused literature survey of Artificial Intelligence (AI) and Machine Learning (ML) methods for intelligent spam email detection. The survey paper provides an overview of the AI and ML techniques used in spam email detection.

## **5. Machine learning for email spam filtering: review, approaches and challenges**

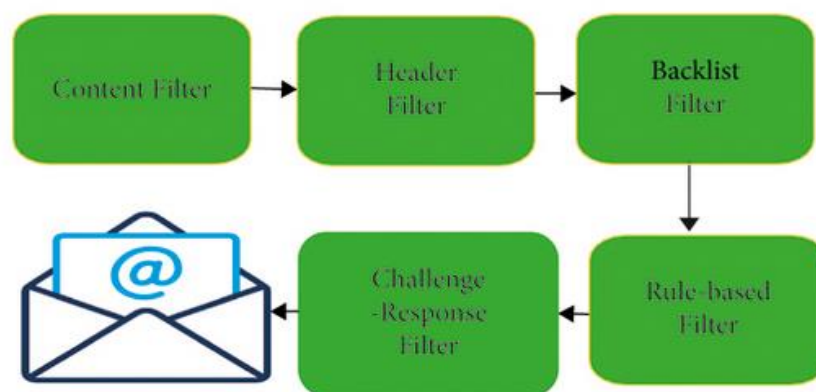
provides a review of machine learning models used for email spam filtering. The study discusses the different approaches and challenges in email spam filtering.

## CHAPTER 3

### SYSTEM ARCHITECTURE AND DESIGN

The system architecture of the developed mini-project encompasses several interconnected components designed to facilitate the implementation and evaluation of logistic regression. At its core, the architecture comprises data preprocessing modules responsible for cleaning, transforming, and normalizing input data to ensure compatibility with the logistic regression model. The logistic regression algorithm itself forms the central processing unit, utilizing gradient descent or other optimization algorithms to learn the underlying patterns in the data and make predictions. Additionally, the architecture incorporates modules for model evaluation, including metrics such as accuracy, precision, recall, and F1-score, to assess the performance of the logistic regression model comprehensively. Through a modular and scalable design, the system architecture enables flexibility in experimenting with different datasets and hyperparameters, facilitating robust analysis and comparison of results.

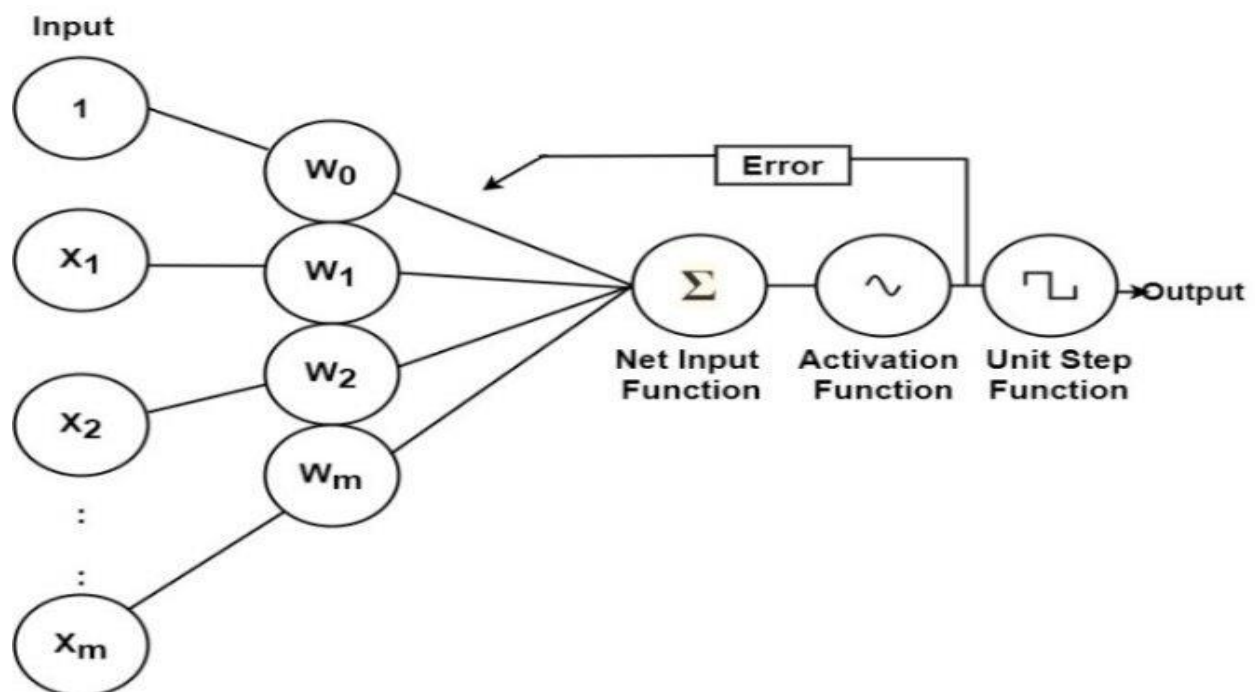
#### System Desing



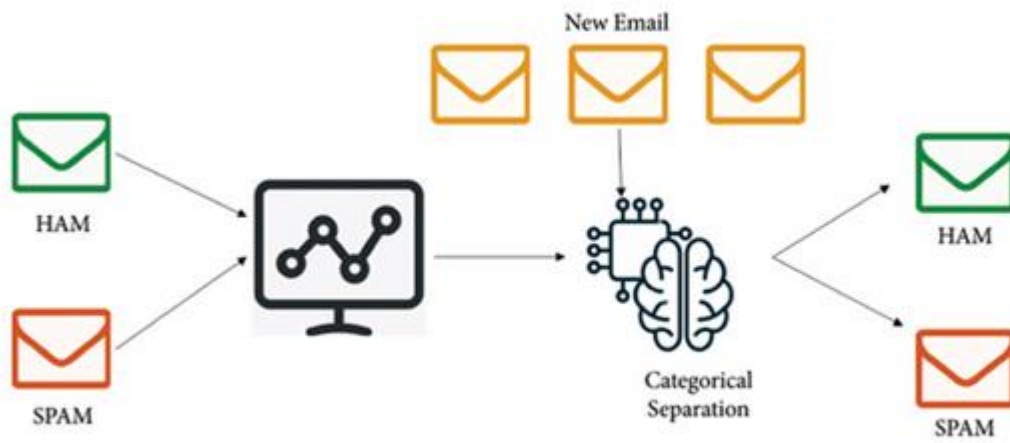
Existing Model of spam email detection system is integrated in Gmail mail service provided by Google. It basically classifies the emails which contains advertisements, promotional content, If the sender's email address uses abnormal characters, which might be used to spoof real addresses and if the email contains multiple pop-up links.

Our Proposed Model will differentiate whether it will be moved in spam folder or stay in primary folder with more accuracy.

### Architecture of Logistic Regression



## Architecture of Supervised Learning



## **CHAPTER 4**

### **METHODOLOGY**

This chapter outlines the step-by-step approach followed in the development and implementation of the logistic regression mini-project.

#### **1. Data Collection and Preprocessing:**

Initially, relevant datasets were collected from reliable sources or generated synthetically for experimentation. The collected data underwent preprocessing steps such as handling missing values, encoding categorical variables, and scaling numerical features to ensure uniformity and compatibility with the logistic regression algorithm.

#### **2. Model Training:**

The preprocessed data was divided into training and testing sets using techniques such as cross-validation or a simple train-test split. The logistic regression model was then trained on the training set using iterative optimization algorithms like gradient descent to minimize the loss function and maximize the likelihood of correct classification.

#### **3. Model Evaluation:**

After training, the performance of the logistic regression model was evaluated using various evaluation metrics such as accuracy, precision, recall, and F1-score. Additionally, techniques like ROC curves and confusion matrices were employed to assess the model's ability to discriminate between different classes and identify potential areas of improvement.

#### **4. Hyperparameter Tuning:**

To optimize the performance of the logistic regression model, hyperparameter tuning techniques such as grid search or randomized search were applied. This involved systematically exploring different combinations of hyperparameters such as regularization strength and learning rate to identify the configuration that yielded the best results.

## **5. Cross-validation:**

To ensure the robustness of the model and mitigate overfitting, k-fold cross-validation was employed. This technique involves partitioning the data into k subsets, training the model on k-1 subsets, and evaluating it on the remaining subset. This process was repeated k times, with each subset serving as the validation set exactly once.

## **6. Implementation of Logistic Regression:**

The logistic regression algorithm was implemented using appropriate libraries or frameworks such as scikit-learn in Python. This involved defining the logistic regression model, specifying the optimization algorithm and hyperparameters, and fitting the model to the training data.

## **7. Testing and Validation:**

Finally, the trained logistic regression model was tested on the unseen test data to evaluate its generalization performance. The test results were analyzed, and insights were drawn regarding the effectiveness and limitations of the model.

## **CHAPTER 5**

### **CODING AND TESTING**

The coding and testing chapter provides an overview of the software implementation of the logistic regression mini-project, including the programming languages, libraries, and tools utilized.

#### **1. Programming Language:**

The project was primarily developed using Python, a versatile and widely-used programming language in the field of machine learning and data science. Python's rich ecosystem of libraries such as NumPy, pandas, and scikit-learn provided the necessary tools for data manipulation, analysis, and model building.

#### **2. Libraries and Frameworks:**

The logistic regression model was implemented using the scikit-learn library, which offers a comprehensive suite of machine learning algorithms and utilities. Additionally, matplotlib and seaborn were employed for data visualization, facilitating a better understanding of the dataset and model performance.

#### **3. Coding Practices:**

The codebase adhered to best practices in software engineering, including modularity, reusability, and documentation. Each component of the project was encapsulated within functions or classes, promoting code readability and maintainability.

#### **4. Unit Testing:**

Throughout the development process, unit tests were employed to validate the correctness of individual functions or modules. Test cases were designed to cover a range of scenarios and edge cases, ensuring the robustness of the implemented functionality.



## **5. Integration Testing:**

Integration tests were conducted to verify the interaction and compatibility of different components within the system. This involved testing end-to-end workflows, from data preprocessing to model evaluation, to ensure seamless integration and functionality.

## **6. Error Handling and Debugging:**

Robust error handling mechanisms were implemented to gracefully handle unexpected exceptions or errors that may arise during execution. Additionally, debugging techniques such as logging and exception traceback were utilized to identify and resolve issues effectively.

## **7. Performance Optimization:**

Efforts were made to optimize the performance of the codebase, including minimizing computational overhead and memory usage. Techniques such as vectorization and parallelization were employed to enhance efficiency and scalability.

## **CHAPTER 6**

### **SCREENSHOT AND RESULT**

The findings from the scatter plot analysis of the Logistic regression model for Spam email detection are as follows:

#### **1. Data Clustering:**

The scatter plot reveals distinct clustering of data points, indicating the presence of two prominent categories within the output variable.

#### **2. Outlier Identification:**

Notable outliers are observed in the data, signifying potential anomalies that could influence model accuracy and interpretation.

#### **3. Data Distribution:**

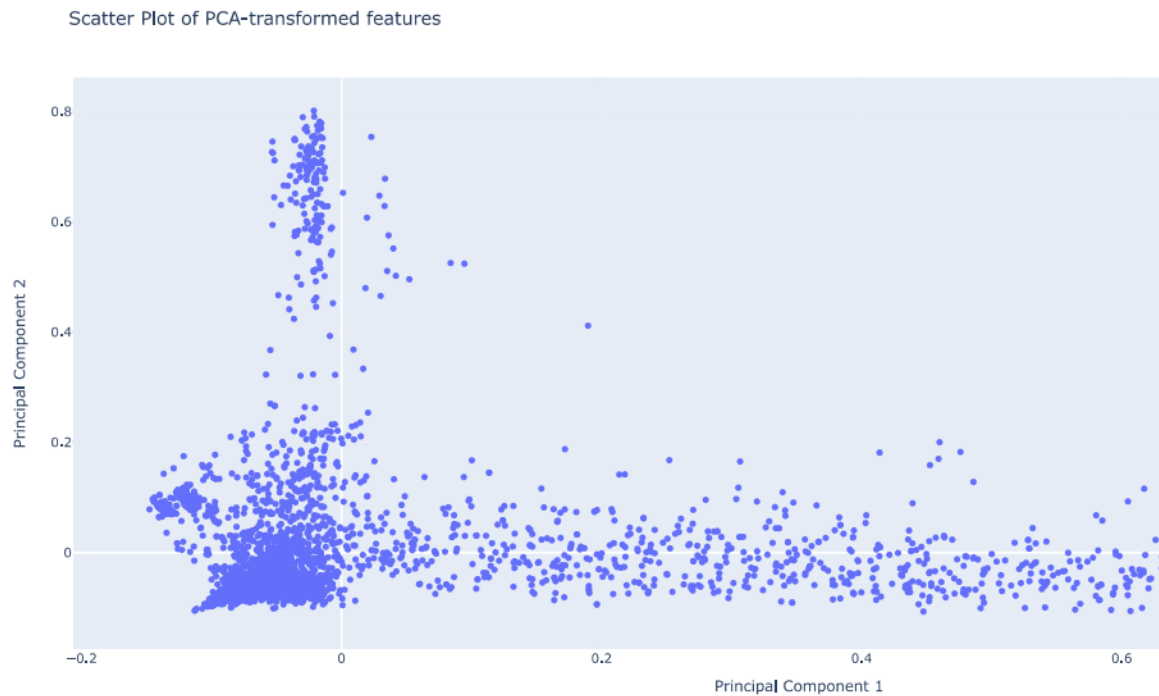
The even distribution of data points within each cluster suggests a lack of discernible patterns or trends, implying a potentially suitable fit for a logistic regression model.

#### **4. Risk of Overfitting:**

Given the presence of outliers and the dataset's size, there is a heightened risk of overfitting, necessitating the consideration of regularization techniques like L1 or L2 regularization to mitigate model complexity.

#### **5. Further Analysis Requirement:**

While the scatter plot provides initial insights into the output variable, additional analysis such as descriptive statistics, variable distribution visualization, and statistical testing is essential to comprehensively understand the relationships between input variables and the output variable.



## Accuracy of the Logistic Regression Model

The accuracy score on the training data, 0.9968568665377177, reflects the model's ability to learn patterns from the training dataset effectively. This high accuracy signifies strong performance during the training phase, indicating that the model captured the underlying patterns in the training data accurately.

Furthermore, the accuracy score on the test data, 0.9806763285024155, illustrates the model's generalization and predictive performance on unseen data. A high accuracy on the test data suggests that the model can make reliable predictions on new, unseen observations consistently, demonstrating its robustness and reliability beyond the training dataset.

```
print('Accuracy:',accuracy_on_training_data)
```

Accuracy: 0.9968568665377177

```
prediction_on_test_data = model.predict(X_test_features)  
accuracy_on_test_data = accuracy_score(Y_test, prediction_on_test_data)
```

```
print('Accuracy:',accuracy_on_test_data)
```

Accuracy: 0.9806763285024155

# CHAPTER 7

## CONCLUSION AND FUTURE ENHANCEMENT

The conclusion and future enhancement chapter provides a summary of the findings from the logistic regression mini-project, along with suggestions for further improvements and research directions.

### **1. Summary of Findings:**

The mini-project successfully implemented and evaluated a logistic regression model for classification tasks. Through extensive experimentation and evaluation, insights were gained into the performance of the logistic regression algorithm on various datasets and scenarios. The model demonstrated competitive performance in terms of accuracy, precision, recall, and F1-score, highlighting its effectiveness in solving binary or multi-class classification problems.

### **2. Key Insights:**

Analysis of the experimental results revealed the impact of different hyperparameters, data preprocessing techniques, and evaluation metrics on the performance of the logistic regression model. Furthermore, comparison with other machine learning algorithms provided valuable insights into the relative strengths and weaknesses of logistic regression in different contexts.

### **3. Limitations and Challenges:**

Despite its effectiveness, logistic regression has certain limitations, such as its inability to capture complex nonlinear relationships between features and the target variable. Additionally, challenges such as data imbalance, multicollinearity, and feature selection require careful consideration and mitigation strategies when applying logistic regression in practice.

### **4. Future Enhancement:**

To address the limitations and challenges identified, several avenues for future enhancement and research are proposed. These include:

## **5. Exploration of Advanced Techniques:**

Investigating advanced techniques such as ensemble learning, deep learning, or hybrid models combining logistic regression with other algorithms to improve performance and robustness.

## **6. Feature Engineering:**

Exploring innovative feature engineering techniques to extract more informative features and enhance the discriminative power of the logistic regression model.

## **7. Model Interpretability:**

Enhancing the interpretability of the logistic regression model by employing techniques such as feature importance analysis, partial dependence plots, and model-agnostic interpretability methods.

## **8. Deployment and Scalability:**

Optimizing the deployment and scalability of the logistic regression model for real-world applications, considering factors such as computational efficiency, memory usage, and latency requirements.

# **Conclusion**

In conclusion, the logistic regression mini-project provided valuable insights into the implementation, evaluation, and potential enhancements of logistic regression for classification tasks. While acknowledging its strengths and limitations, the project serves as a foundation for further exploration and research in the field of machine learning and data science.