

# NAAN MUDHALVAN IBM PROJECT REPORT

**Domain : Applied Data Science**

PROJECT TITLE

## **Audit AI: A Machine Learning for Detecting Fraud in Audit Data**

<b><i>TEAM ID</i></b> <i>NM2023TMID01937</i>
<b><u><i>TEAM MEMBERS</i></u></b>
<i>Chandrasekaran S (Team Leader) - 420420104301</i>
<i>Ajith R. (Team Member 1)- 420420104001</i>
<i>Tamilselvan T ( Team Member 2) - 420420104309</i>
<i>Vijay Kumar E ( Team Member 3) - 420420104311</i>

## CONTENTS

### 1 INTRODUCTION

#### 1.1.1 Project Overview

#### 1.1.2 Purpose

### 2 LITERATURE SURVEY

### 3 IDEATION & PROPOSED SOLUTION

#### 3.1.1 Problem Statement Definition

#### 3.1.2 Empathy Map Canvas

#### 3.1.3 Ideation & Brainstorming

#### 3.1.4 Proposed Solution

### 4 REQUIREMENT ANALYSIS

#### 4.1.1 Functional requirement

#### 4.1.2 Non-Functional requirements

#### 4.1.3 Technical architecture

### 5 PROJECT DESIGN

#### 5.1.1 Data Flow Diagrams

#### 5.1.2 Solution & Technical Architecture

#### 5.1.3 User Stories

### 6 CODING & SOLUTIONING

#### 6.1.1 Feature 1

#### 6.1.2 Feature 2

### 7 RESULTS

#### 7.1.1 Performance Metrics

### 8 ADVANTAGES & DISADVANTAGES

### 9 CONCLUSION

### 10 FUTURE SCOPE

### 11 APPENDIX

.

# **CHAPTER-1**

## **INTRODUCTION**

## **1.INTRODUCTION:**

Fraudulent activities pose significant risks to organizations, leading to financial losses, reputational damage, and legal consequences. In the field of auditing, traditional approaches to detecting fraud often rely on manual inspection and rule-based methods, which can be time-consuming, subjective, and limited in their ability to uncover complex fraud schemes.

To overcome these limitations, machine learning has emerged as a powerful tool for detecting fraud in audit data. The objective of this project is to develop a machine learning model for detecting fraud in audit data. By training the model on historical data and using it to analyze real-time transactions, we aim to identify potentially fraudulent activities and alert auditors or relevant stakeholders for further investigation.

### **1.1 Project Overview**

The goal of this project is to develop a machine learning model for detecting fraud in audit data. Fraudulent activities can have severe consequences for organizations, including financial losses and reputational damage. Traditional auditing methods often fall short in uncovering sophisticated fraud schemes, which is where machine learning can play a crucial role.

Throughout the project, document the methodologies, findings, and challenges faced. Maintain clear communication with stakeholders and subject matter experts to ensure the model aligns with their needs and incorporates their insights.

### **1.2 Purpose**

The purpose of using machine learning for detecting fraud in audit data is to enhance the effectiveness and efficiency of fraud detection processes in auditing.

Improved Detection Accuracy, Early Fraud Detection, Scalability and Efficiency, Adaptability to Evolving Fraud Patterns, Identification of Complex Fraud Networks, Augmentation of Human Expertise, Continuous Improvement.

**CHAPTER-2**  
**LITERATURE SURVEY**

## 2. LITERATURE SURVEY

### 2.1 Existing problem

The existing auditing processes in organizations often rely on manual inspection and rule-based methods, which can be time-consuming, subjective, and limited in their ability to detect complex fraud patterns and anomalies. This leads to potential risks of undetected fraudulent activities, resulting in financial losses, reputational damage, and legal consequences for the organizations

### 2.2 References

#### Paper 1

**Title:** "Intelligent Fraud Detection in Financial Statements Using Machine Learning and Data Mining: A Systematic Literature Review"

**Authors:** Matin N. Ashtiani and Bijan Raahemi

Published: July 13, 2021 (Current version as of July 15, 2022)

**Problem:** Detecting fraudulent financial statements is challenging due to complex and diverse data. The paper reviews machine learning and data mining approaches for efficient and accurate fraud detection.

**Methodology/Algorithm:** The authors used the Kitchenham methodology and explored regression, ensemble methods, and clustering techniques.

**Advantages:** Machine learning and data mining enable quick and accurate processing of large data volumes. They uncover complex patterns and improve efficiency in fraud detection, reducing investigation time and resources.

**Disadvantages:** Challenges include the need for high-quality data, potential lack of model transparency, and the possibility of some fraud cases being missed.

#### Paper 2

**Title:** "An Analysis on Financial Statement Fraud Detection for Chinese Listed Companies Using Deep Learning"

**Authors:** Wu Xiuguo and Du Shengyong

**Published:** February 22, 2022 (Current version as of March 4, 2022)

**Problem Definition:** The research aims to develop models using deep learning algorithms to detect financial statement fraud in Chinese listed companies' annual reports. Traditional approaches are ineffective due to companies' tactics, leading to significant losses for stakeholders. The objective is to build models with high classification performance and develop a classification framework that utilizes both textual and numerical data in annual reports.

**Methodology/Algorithm:** Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) are employed in the research.

**Advantages:** The research presents a framework that utilizes deep learning algorithms to detect financial statement fraud in Chinese listed companies' annual reports. Deep learning algorithms offer high classification performance that traditional approaches cannot achieve due to companies' tactics. The study also emphasizes the importance of textual analytics in detecting fraud in financial documents.

**Disadvantages:** The research has certain limitations, such as the limited sampling period of five years and the potential exclusion of delisted companies and incomplete annual reports, which could impact prediction results. Additionally, the data source only includes Chinese listed companies, limiting the applicability of the models to other markets, thus requiring further study.

### **Paper 3**

**Title:** "Machine Learning- and Evidence Theory-Based Fraud Risk Assessment of China's Box Office"

**Authors:** Shi Qiu and Hong-Qu He

**Published:** November 27, 2018 (Current version as of December 27, 2018)

**Problem Definition:** Box-office fraud in China poses a significant problem for the movie market, misleading consumers and investors, and potentially harming the developing motion picture industry and shadow movie market in the country. Accurate supervision and auditing are necessary to regulate the market and detect financial fraud.

**Methodology/Algorithm:** The framework includes an ordered logistic regression algorithm and iterative aggregation of different evidence to calculate the basic probability assignment.

**Advantages:** The proposed NFM method (based on evidence theory) serves as an effective complementary method to improve the efficiency of auditing and supervising box-office fraud in China. It utilizes publicly available data from various websites, including Web 2.0, to assess fraud risk. The evidence theory-based framework offers a convenient iterative assessment procedure, and the ordered logistic regression algorithm enhances the accuracy of fraud risk assessment.

**Disadvantages:** One potential drawback of the proposed NFM method is that it may not capture all aspects of box-office fraud, relying heavily on publicly available data that may not detect more covert fraudulent activities. Additionally, the algorithm used for fraud risk assessment may have limitations in certain situations, and the method may require continuous updates to effectively detect new fraud schemes.

## **Paper 4**

**Title:** "A Systematic Literature Review of Fraud Detection Metrics in Business Processes"

**Authors:** Badr Omair and Ahmad Alturki

**Published:** February 4, 2020 (Current version as of February 12, 2020)

**Problem Definition:** The literature review highlights the need for a comprehensive approach to detecting possible business process fraud (PBF). While some metrics for fraud detection have been studied, there is a lack of attention specifically given to PBF, which occurs when business processes deviate from standard operating procedures. Only a few detection metrics have been identified, and they do not fully address the different conceptual perspectives on business processes. Therefore, further research is needed to extend the current detection metrics for PBF and evaluate their effectiveness in exposing common fraud risks.

**Methodology/Algorithm:** The study employs a combination of Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), and ordered logistic regression algorithms.

**Advantages:** The content analysis methodology used in this study enables a comprehensive review of relevant literature



## Paper 5

**Title:** "Machine Learning for Financial Risk Management: A Survey"

**Authors:** Akib Mashrur, Wei Luo, Nayyar A. Zaidi, and Antonio Robles-Kelly

**Published:** November 5, 2020 (Current version as of November 19, 2020)

**Problem Definition:** Traditional statistical models for financial risk management have limitations, necessitating the use of advanced machine learning models to enhance accuracy and robustness. However, challenges such as long training times, the requirement for large amounts of data, and the lack of explainability and fairness in machine learning models need to be addressed

**Methodology/Algorithm:** The study employs GARCH or stochastic volatility models, Long Short-Term Memory (LSTM) networks, and Natural Language Processing (NLP) techniques.

**Advantages:** Advanced machine learning models, including deep learning approaches, offer significant improvements in accuracy and robustness for financial risk management compared to traditional statistical models. Machine learning models can also incorporate unstructured textual data, enhancing the accuracy of volatility forecasting. Additionally, federated learning systems enable private and secure learning utilizing sensitive financial data.

**Disadvantages:** Long training times and the need for large amounts of data in certain machine learning models, such as LSTMs, may limit their applicability in certain situations. Machine learning models also lack explainability and fairness, which can be a concern in financial risk management. Moreover, effectively utilizing unstructured textual data for volatility forecasting may require sophisticated NLP techniques to accurately extract relevant information.

### 2.3 Problem Statement Definition

As fraud detection becomes increasingly important in auditing, there is a need for an effective and efficient way to detect fraudulent activities in audit data. With the vast amount of data being generated by companies, it can be difficult for auditors to manually identify potential fraudulent transactions or patterns. Therefore, a machine learning model can be developed to assist auditors in detecting fraudulent activities in audit data.

The goal of this project is to build a machine learning model that can accurately identify instances of fraud in audit data, based on historical data of known fraudulent activities.

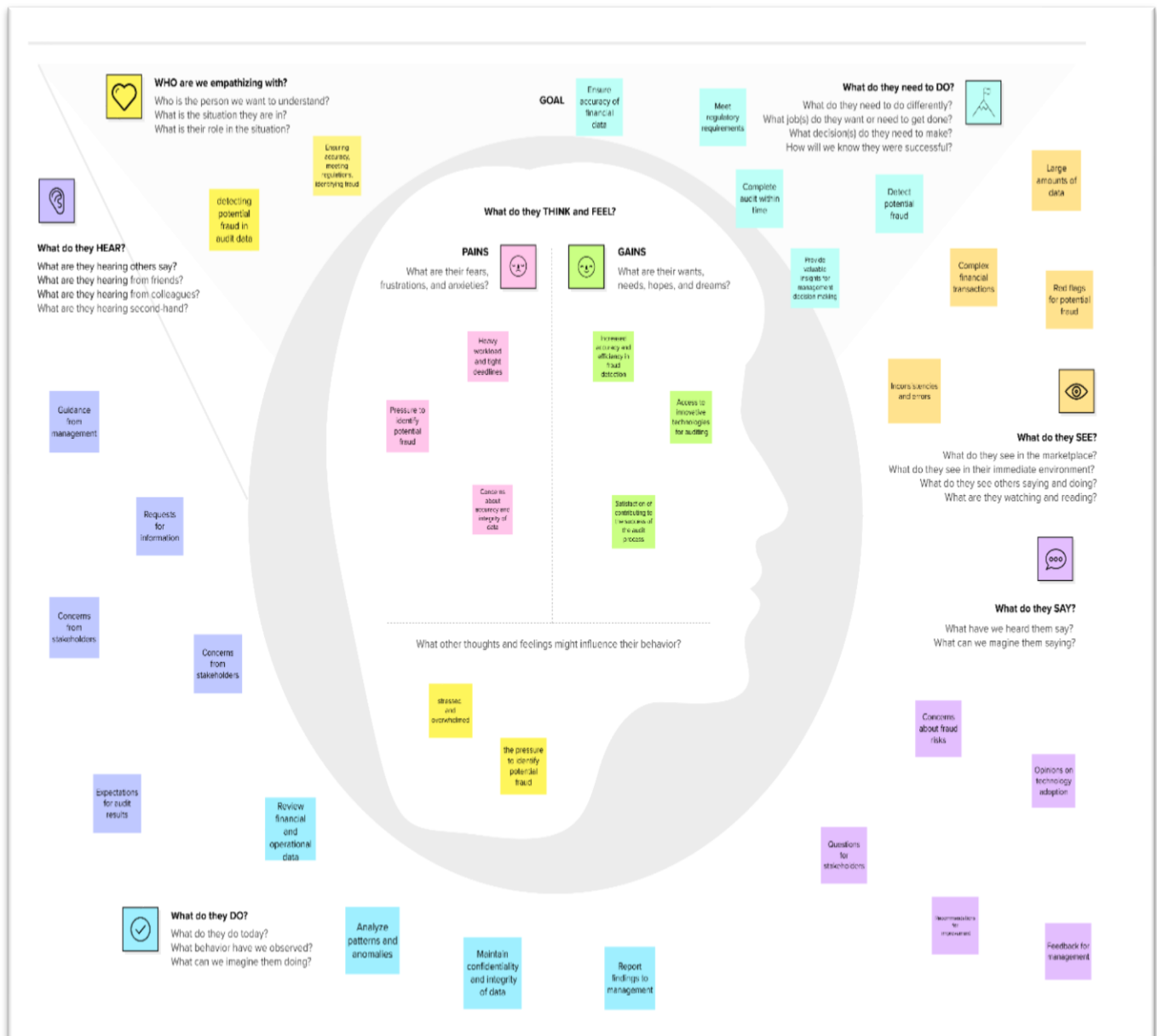
<b>I am</b>	Describe customer with 3-4 key characteristics - <i>who are they?</i>	Describe the customer and their attributes here
<b>I'm trying to</b>	List their outcome or "job" the care about - <i>what are they trying to achieve?</i>	List the thing they are trying to achieve here
<b>but</b>	Describe what problems or barriers stand in the way - <i>what bothers them most?</i>	Describe the problems or barriers that get in the way here
<b>because</b>	Enter the "root cause" of why the problem or barrier exists - <i>what needs to be solved?</i>	Describe the reason the problems or barriers exist
<b>which makes me feel</b>	Describe the emotions from the customer's point of view - <i>how does it impact them emotionally?</i>	Describe the emotions the result from experiencing the problems or barriers



**CHAPTER-3**  
**IDEATION & PROPOSED SOLUTION**

### 3. IDEATION AND PROPOSED SOLUTION

#### 3.1 Empathy Map Canvas



Reference: <https://app.mural.co/t/auditaiibm4517/m/auditaiibm4517/1682776371564/d60202f5f2c95fe3e0b3e8cfa879ec0a1794722?sender=u1f347ebc11e76f570f911985>

## 3.2 Ideation & Brainstorming

2

### Brainstorm

Write down any ideas that come to mind that address your problem statement.

10 minutes

#### TIP

You can select a sticky note and hit the pencil button to delete, don't start drawing!

Provide user-friendly UI

enable users to view their data in a more meaningful way

segment users with the help of machine learning

Provide the user with a personalized dashboard

CHANDRASEKHARAN S

AUTH R

Use machine learning to identify potential issues

use machine learning to identify potential issues

use machine learning to identify potential issues

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

TAMIL SELVAN T

VLAYAKUMAR E

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

3

### Group ideas

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you can break it up into smaller sub-groups.

20 minutes

#### TIP

Add a sentence-like label to each cluster of notes to describe the idea, process, capability, and category of the idea to help you group them better.

#### Services

User-friendly UI

Provide user with a personalized dashboard

Use machine learning to identify potential issues

provide user with a personalized dashboard

#### Additional Features

Use machine learning to identify potential issues

Use machine learning to identify potential issues

Use machine learning to identify potential issues

provide user with a personalized dashboard

#### Support

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

#### Benefits

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

provide user with a personalized dashboard

## Idea Prioritization

4

**Prioritize**

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

20 minutes

**Importance**

If each of these ideas could get done without any difficulty or cost, which would have the most positive impact?

**Feasibility**

Regardless of their importance, which tasks are more feasible than others? (Cost, time, effort, complexity, etc.)

**Tip**

Participants can use their cursor to point at where sticky notes should go on the grid. The facilitator can confirm the spot by using the laser pointer holding the H key on the keyboard.

**Options to enhance that service parameters**

**provide better personalized services**

**define documentation structure**

**offer a dedicated onboarding process to new users**

**provide better onboarding**

**prepare the use of interactive dashboards**

**enhance the overall user experience by creating a dedicated onboarding process**

**develop self-reports**

**use a tiered link approach**

**increase security and efficiency**

**address customer onboarding for new users**

**align with customer needs**

→

**After you collaborate**

You can export the mural as an image or pdf to share with members of your company who might find it helpful.

**Quick add-ons**

- Share the mural**  
 Share a view link to the mural with stakeholders to keep them in the loop about the outcomes of the session.
- Export the mural**  
 Export a copy of the mural as a PNG or PDF to attach to emails, include in slides, or save in your drive.

**Keep moving forward**

- Strategy blueprint**  
 Define the components of a new idea or strategy.  
[Open the template →](#)
- Customer experience journey map**  
 Understand customer needs, motivations, and obstacles for an experience.  
[Open the template →](#)
- Strengths, weaknesses, opportunities & threats**  
 Identify strengths, weaknesses, opportunities, and threats (SWOT) to develop a plan.  
[Open the template →](#)

[Share template feedback](#)

Reference:

<https://app.mural.co/t/auditaibm4517/m/auditaibm4517/1682816643232/850ca275d9fd5f53c216306f3d573caa19695f63?sender=u1f347ebc11e76f570f911985>

### 3.3 Proposed Solution

S.No.	Parameter	Description
1.	Problem Statement (Problem to be solved)	Auditors require an automated and reliable fraud detection system to enhance the security and trustworthiness.
2.	Idea / Solution description	<p>The ML model is used to detect fraud in audit data</p> <p>The model should analyse large volumes of data, identify patterns, and anomalies to provide auditors with actionable insights and learn from past audits to improve its accuracy over time</p>
3.	Novelty / Uniqueness	<p>It allows for the automation of a traditionally manual process, saving time and reducing errors.</p> <p>Its ability to learn from past audits and improve its accuracy over time</p>
4.	Social Impact / Customer Satisfaction	<p>The efficiency and effectiveness of audits, and its potential to prevent significant financial losses and reputational damage, can ultimately lead to increased customer satisfaction.</p> <p>It saves time and resources so the auditors to focus on more complex and value-adding tasks</p>
5.	Business Model (Revenue Model)	This could be done on a subscription or per-audit basis, with clients paying a fee for each audit analysed.
6.	Scalability of the Solution	<p>It must handle large volumes of data, with sufficient computing power and storage capacity.</p> <p>The machine learning algorithms used would need to be optimized to ensure that they can process data quickly and accurately, even as the dataset grows larger</p>

**CHAPTER – 4**  
**REQUIREMENT ANALYSIS**



## 4.1 Functional requirement

FR No.	Functional Requirement (Epic)	Sub Requirement (Story / Sub-Task)
FR-1	User Registration	Registration through Form Registration through Gmail Registration through LinkedIn
FR-2	User Confirmation	Confirmation via Email
FR-3	Fraud Detection by machine learning algorithms	System can detect unusual patterns and anomalies in financial and operational data  System can continuously learn from historical data to improve accuracy  System can generate alerts for auditors  System can provide detailed reports
FR-4	User-friendly interface for auditors	Auditors can easily view audit data  Auditors can filter and search for relevant data and activities Auditors can customize alerts and reports
FR-5	web-based interface for auditors	System is deployed on IBM Cloud platform to ensure scalability, availability, and security  System can integrate with other systems, such as accounting and financial systems, for seamless data exchange
FR-6	Performance System can quickly and accurately	System can provide real-time alerts and reports to auditors without delay or latency  System can handle high volumes of audit data without sacrificing performance or accuracy

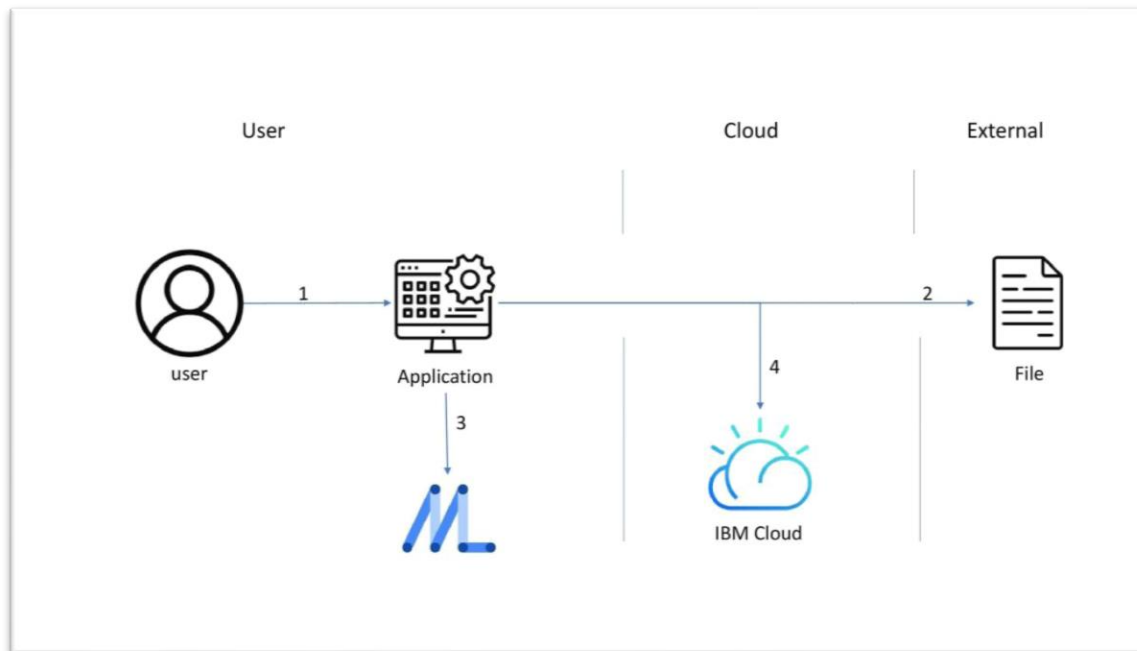
## 4.2 Non-Functional requirements

FR No.	Non-Functional Requirement	Description
NFR1	Usability	User-friendly interface with clear instructions and intuitive navigation.
NFR2	Security	Ensure data confidentiality, authorization, encryption, regular audits, and vulnerability assessments.
NFR3	Reliability	Consistent and error-free operation, quick recovery from errors, and avoidance of single points of failure.
NFR4	Performance	Timely analysis of large data volumes, fast response times, and ability to handle high loads without sacrificing performance.
NFR5	Availability	Minimal downtime or maintenance, avoidance of scheduled maintenance, and disaster recovery plan in place.
NFR6	Scalability	Designed to handle growth in data volume and user numbers, ability to scale up or down as needed, and handle peak loads without additional resources.

**CHAPTER-5**  
**PROJECT DESIGN**

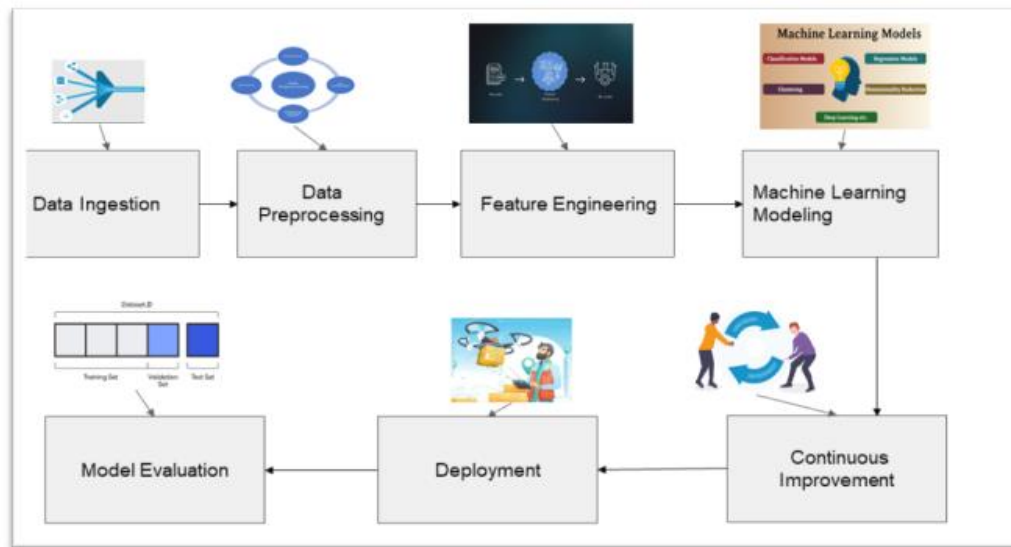
## 5. PROJECT DESIGN

### 5.1 Data Flow Diagram:

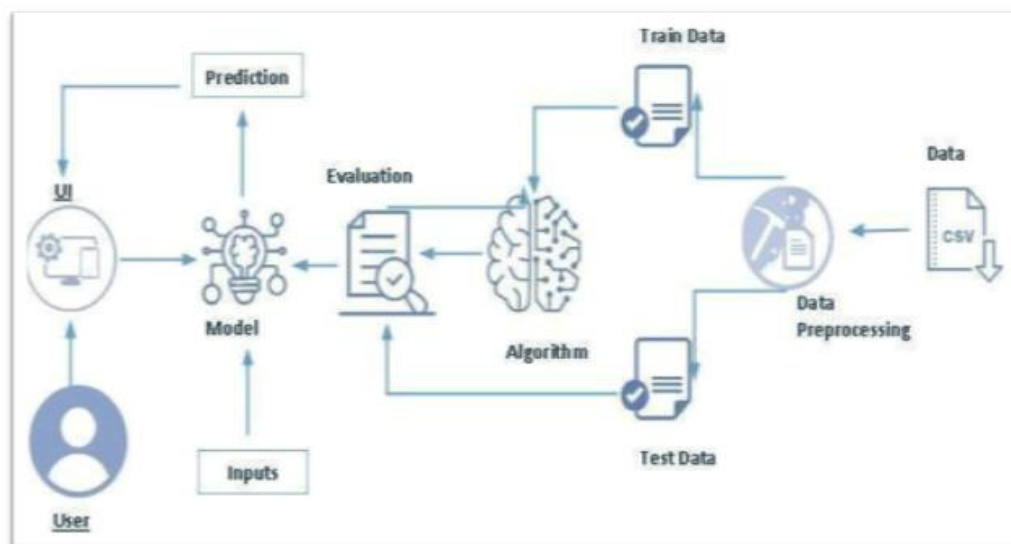


1. The user login into the webpage.
2. The User select and load the data needed to be processed and identified.
3. The input data is given to the machine learning model, and it process them with given algorithms and give the output data.
4. Dataset needed is stored in the IBM Cloud storage.

## 5.2 Solution & Technical Architecture



## Technical Architecture



### 5.3 User Stories

User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Team Member
Business Manager	Access to fraud detection reports	USN-1	I want to be able to access fraud detection reports so that I can monitor potential fraud within my organization.	Reports must be available in real-time and easily accessible through a user-friendly interface.	High	Chandrasekaran
Data Analyst	Access to raw audit data	USN-2	I want to be able to access raw audit data so that I can perform my own analysis and identify potential fraud	Data must be accessible through a secure and user-friendly interface, with appropriate access controls in place.	High	Ajith
IT Administrator	System monitoring and management	USN-3	I want to be able to monitor and manage the fraud detection system to ensure that it is running smoothly and efficiently.	System must be easily monitorable and manageable through a user-friendly interface, with alerts in place for potential issues or failures.	High	Vijay Kumar
Business User	View real-time dashboard of audit data	USN-4	As a business user, I want to be able to view a real-time dashboard of audit data so that I can quickly identify any potential fraud or anomalies.	The dashboard should display relevant metrics and KPIs, such as total transactions, average transaction amount, and percentage of flagged transactions	Medium	Tamilselvan

**CHAPTER-6**  
**CODING & SOLUTIONING**

## 6.1 Feature 1

### Python Flask

Python Flask is used to develop web integration with our model. Flask is mainly used to render and integrate the Audit AI application in the browser. By running the python application, the suitable server domain link is obtained and run in the browser.

### HTML

The HTML and CSS is used to design the overall Audit AI UI. HTML is used to add UI components and CSS is used to add style to those components

### Jupyter notebook

The jupyter notebook is used to create our model ,test and train with dataset we use many algorithms in our model like logistic regression ,random forest , xgboost and knn from this algorithm the best one chosen and we use knn algorithm in this project . pickle is used to store the model and integrate with the flask

#### *Build PYTHON FLASK Code:*

##### **app.py**

```
from flask import Flask, request, render_template
import pickle
```

```
app = Flask(__name__)
model = pickle.load(open('knn.pkl', 'rb'))
```

```
@app.route('/')
def home():
    return render_template('index.html')
```

```
@app.route('/project')
def project():
    return render_template('project.html')
```

```
@app.route('/riskpred', methods=['GET', 'POST'])
def predict_risk():
    if request.method == "POST":
        pred = [
            float(request.form['Sector_score']),
            float(request.form['PARA A']),
            float(request.form['Risk A']),
            float(request.form['PARA B']),
            float(request.form['Risk B']),
            float(request.form['TOTAL']),
```



```

        float(request.form['numbers']),
        float(request.form['Money_Value']),
        float(request.form['Score_MV']),
        float(request.form['District_Loss']),
        float(request.form['History']),
        float(request.form['Score']),
        float(request.form['Inherent Risk']),
        float(request.form['Audit_Risk'])
    ]
    output = model.predict([pred])[0]
    return render_template('result.html', predict="The Predicted Risk value is: " + str(output))
else:
    return render_template('main.html')

if __name__ == '__main__':
    app.run(debug=True)

```

## index.html

```

<!DOCTYPE html>
<html>
<head>
    <link rel="stylesheet" type="text/css" href="{{ url_for('static', filename='styles.css') }}">

    <title>Welcome to My Project</title>
</head>
<body>
    <h1>Welcome to Audit AI Home page</h1>
    <div class="image-container">
        
    </div>
    <p>Audit AI: A Machine Learning for Detecting Fraud in Audit Data</p>
    <br>"Description about our site
    click learn more"
    <div class="button-container">
        <a href="/project" class="button">Learn More</a>
        <a href="/riskpred" class="button">Start Prediction</a>
    </div>
</body>
</html>

```

## Main.html

```

<!DOCTYPE html>
<html>
<head>
    <link rel="stylesheet" type="text/css" href="{{ url_for('static', filename='styles.css') }}">
    <title>Main Page</title>
</head>
<body>
    <div class="container">

```

```
<h1>Detecting Fraud in Audit Data</h1>
<form action="/riskpred" method="POST">
  <label for="Sector_score">Sector Score:</label>
  <input type="number" name="Sector_score" step="any" required>
  <label for="PARA_A">PARA A:</label>
  <input type="number" name="PARA_A" step="any" required><br><br>

  <label for="Risk A">Risk A:</label>
  <input type="number" name="Risk A" step="any" required><br><br>

  <label for="PARA_B">PARA B:</label>
  <input type="number" name="PARA_B" step="any" required><br><br>

  <label for="Risk B">Risk B:</label>
  <input type="number" name="Risk B" step="any" required><br><br>

  <label for="TOTAL">TOTAL:</label>
  <input type="number" name="TOTAL" step="any" required><br><br>

  <label for="numbers">Numbers:</label>
  <input type="number" name="numbers" step="any" required><br><br>

  <label for="Money_Value">Money Value:</label>
  <input type="number" name="Money_Value" step="any" required><br><br>

  <label for="Score_MV">Score MV:</label>
  <input type="number" name="Score_MV" step="any" required><br><br>

  <label for="District_Loss">District Loss:</label>
  <input type="number" name="District_Loss" step="any" required><br><br>

  <label for="History">History:</label>
  <input type="number" name="History" step="any" required><br><br>

  <label for="Score">Score:</label>
  <input type="number" name="Score" step="any" required><br><br>

  <label for="Inherent_Risk">Inherent Risk:</label>
  <input type="number" name="Inherent_Risk" step="any" required><br><br>

  <label for="Audit_Risk">Audit Risk:</label>
  <input type="number" name="Audit_Risk" step="any" required><br><br>
  <input type="submit" value="Submit">
</form>
</div>
</body>
</html>
```

## **CHAPTER-7**

### **RESULTS**

## Model Performance Testing:

Project team shall fill the following information in the model performance testing template.

S.No.	Parameter	Values	Screenshot
1.	Metrics	<p><b>Knn Model:</b> MAE - , MSE - , RMSE - , R2 score –</p> <p><b>Classification Model:</b> Confusion Matrix - , Accuracy Score- &amp; Classification Report -</p>	<pre>[139] mae = mean_absolute_error(y_test, knn_test_pred)  # Calculate MSE mse = mean_squared_error(y_test, knn_test_pred)  # Calculate RMSE rmse = mean_squared_error(y_test, knn_test_pred, squared=False)  # Calculate R2 score r2 = r2_score(y_test, knn_test_pred)  # Print MAE, MSE, RMSE, R2 score print("MAE:", mae) print("MSE:", mse) print("RMSE:", rmse) print("R2 score:", r2)</pre> <p>MAE: 0.04721030042918455 MSE: 0.04721030042918455 RMSE: 0.21727931431497235 R2 score: 0.8008547008547009</p> <pre># Print Confusion Matrix print("Confusion Matrix:") print(confusion_matrix(y_test, knn_test_pred))  # Print Accuracy Score accuracy = accuracy_score(y_test, knn_test_pred) print("Accuracy:", accuracy)  # Print Classification Report print("Classification Report:") print(classification_report(y_test, knn_test_pred))</pre> <pre>Confusion Matrix: [[140  3]  [ 8 82]] Accuracy: 0.9527896995708155 Classification Report:               precision    recall  f1-score   support       0       0.95       0.98       0.96         143      1       0.96       0.91       0.94          90     accuracy          0.95         0.95         0.95         233   macro avg       0.96       0.95       0.95         233  weighted avg       0.95       0.95       0.95         233</pre>

2.	Tune the Model	Hyperparameter Tuning	<pre> from sklearn.neighbors import KNeighborsClassifier from sklearn.impute import SimpleImputer from sklearn.model_selection import GridSearchCV  # Assuming your DataFrame is named 'df' # Extract the features (x_train) and the target variable (y_train) x_train = df.drop('Risk', axis=1) y_train = df['Risk'] # Select only the desired 14 features from x_train selected_features = ['Sector_score', 'LOCATION_ID', 'PARA_A', 'Score_A', 'Risk_A',                     'PARA_B', 'Score_B', 'Risk_B', 'TOTAL', 'numbers',                     'Money_Value', 'Score_MV', 'District_Loss', 'History'] x_train_selected = x_train[selected_features] # Handle missing values by replacing them with the mean of each column imputer = SimpleImputer(strategy='mean') x_train_selected_imputed = imputer.fit_transform(x_train_selected) # Create the KNeighborsClassifier knn = KNeighborsClassifier() # Define the hyperparameter grid for tuning param_grid = {'n_neighbors': [3, 5, 7], 'weights': ['uniform', 'distance']} # Perform GridSearchCV with 5-fold cross-validation grid_search = GridSearchCV(knn, param_grid, cv=5) grid_search.fit(x_train_selected_imputed, y_train) # Print the best hyperparameters found print("Best Hyperparameters:", grid_search.best_params_) # Print the best model's score print("Best Model Score:", grid_search.best_score_) </pre> <p>Best Hyperparameters: {'n_neighbors': 3, 'weights': 'distance'} Best Model Score: 0.9213647642679901</p> <pre> [145] param_grid = {'n_neighbors': [3, 5, 7, 9],                   'weights': ['uniform', 'distance'],                   'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute']}  [146] from sklearn.model_selection import GridSearchCV  [147] knn = GridSearchCV(knn, param_grid, cv=5, n_jobs=-1)  [148] knn.fit(x_train_selected_imputed, y_train)  GridSearchCV   estimator: KNeighborsClassifier     KNeighborsClassifier  [149] # Print the best hyperparameters and corresponding mean cross-validated score print("Best hyperparameters: ", knn.best_params_)  Best hyperparameters: {'algorithm': 'auto', 'n_neighbors': 9, 'weights': 'distance'}  [150] print("Best mean cross-validated score: {:.2f}".format(knn.best_score_))  Best mean cross-validated score: 0.92 </pre>

**CHAPTER-8**  
**ADVANTAGES & DISADVANTAES**

## Advantages

1. **Improved Fraud Detection:** The project utilizes machine learning algorithms to analyze audit data and identify potential fraudulent activities. This can enhance the accuracy and efficiency of fraud detection compared to traditional manual methods.
2. **Time and Cost Efficiency:** By automating the fraud detection process, AuditAI can save significant time and resources for auditors. It can quickly process large volumes of data, identify suspicious patterns, and prioritize high-risk areas, leading to more efficient audits.
3. **Early Detection:** The project aims to detect fraud at an early stage, allowing organizations to take timely actions to mitigate the impact of fraudulent activities. Early detection can potentially minimize financial losses and reputational damage.
4. **Scalability:** Machine learning models used in the project can be scaled up to handle large datasets and adapt to changing business environments. This scalability ensures the system's effectiveness in detecting fraud as the volume and complexity of data increase.
5. **Continuous Learning:** The system can continuously learn from new data and update its algorithms to adapt to evolving fraud patterns. This enables the system to improve over time and stay up to date with emerging fraud techniques.

## Disadvantages

1. **False Positives:** Like any machine learning system, AuditAI may generate false positive results, flagging legitimate transactions as fraudulent. This can result in unnecessary investigations and potential disruption to normal business operations.

2. **Data Limitations:** The accuracy of the fraud detection system heavily relies on the quality and availability of data. If the input data is incomplete, inaccurate, or biased, it can negatively impact the system's performance and effectiveness.
3. **Complexity and Technical Expertise:** Implementing and maintaining the AuditAI system requires technical expertise in machine learning, data analytics, and system integration. Organizations may need to invest in skilled resources or collaborate with external experts to ensure successful implementation.
4. **Regulatory Compliance:** The project needs to comply with legal and regulatory requirements regarding data privacy and security. Ensuring compliance with industry standards and regulations can be a challenge and may require ongoing monitoring and updates to the system.
5. **System Integration:** Integrating AuditAI with existing audit processes, systems, and workflows can be complex and time-consuming. It may require significant changes to the organization's infrastructure and processes to ensure seamless integration and effective utilization of the system.
6. **Ethical Considerations:** The project should address ethical concerns related to privacy, data handling, and potential biases in the algorithms used for fraud detection. Ensuring fairness, transparency, and accountability in the system's operation is crucial.
7. **Human Expertise and Judgment:** While AuditAI can automate certain aspects of fraud detection, it should not replace human expertise and judgment. Human auditors are still essential for interpreting results, conducting in-depth investigations, and making informed decisions based on the detected anomalies..



**CHAPTER-9**  
**CONCLUSION**

## CONCLUSION

The project "AuditAI: A Machine Learning for Detecting Fraud in Audit Data" offers several advantages in enhancing fraud detection in the audit process. By leveraging machine learning algorithms, the project improves the efficiency, accuracy, and early detection of fraudulent activities. It can save time and resources for auditors, enable scalable analysis of large datasets, and continuously learn from new data.

However, it is important to consider the potential disadvantages and challenges associated with the project. False positives, data limitations, complexity in implementation and integration, regulatory compliance, ethical considerations, and the need for human expertise and judgment are factors that should be carefully addressed.

Despite these challenges, AuditAI presents a promising solution to detect and prevent fraud in audit data. With proper planning, implementation, and ongoing monitoring, the project can contribute to more effective and efficient audit processes, mitigating financial losses and reputational damage caused by fraudulent activities.

It is recommended to conduct thorough testing, validation, and feedback from stakeholders to ensure the project's effectiveness, reliability, and compliance with industry standards and regulations. Regular updates and improvements should be made to address emerging fraud patterns and changing business environments.

Overall, the AuditAI project has the potential to revolutionize fraud detection in audits, providing auditors and organizations with valuable insights and tools to safeguard their financial integrity and reputation.

**CHAPTER 10**  
**FUTURE SCOPE**

## FUTURE SCOPE

1. **Enhanced Accuracy:** Continuously improving the accuracy of fraud detection algorithms is essential. Researchers can explore advanced machine learning techniques, such as deep learning, to achieve higher precision and recall rates in identifying fraudulent patterns and anomalies.
2. **Real-Time Monitoring:** Integrating AuditAI into real-time monitoring systems can provide timely alerts and notifications when suspicious activities or anomalies are detected. This can help auditors and investigators respond quickly to potential fraud cases.
3. **Advanced Data Analytics:** Incorporating advanced data analytics techniques, such as network analysis, social network analysis, and graph algorithms, can enable the identification of complex fraud networks and uncover hidden relationships among individuals or entities involved in fraudulent activities.
4. **Automated Risk Assessment:** Developing AI algorithms that can automatically assess the risk level of different audit areas or transactions can optimize resource allocation and focus auditing efforts on high-risk areas, increasing efficiency and effectiveness.
5. **Integration with Other Systems:** Integrating AuditAI with existing audit management systems, data visualization tools, and reporting platforms can streamline the fraud detection process and facilitate seamless information exchange among auditors, investigators, and management.
6. **Continuous Learning:** Implementing a feedback loop mechanism that allows the system to learn from detected fraud cases, adjust algorithms, and improve over time can enhance the system's effectiveness in detecting evolving fraud patterns and adapting to new types of fraud.

7. **Regulatory Compliance:** Adapting AuditAI to comply with changing regulatory requirements and standards in the auditing and financial sectors is crucial. Staying up to date with regulations and incorporating necessary features into the system can ensure its relevancy and reliability.
8. **Privacy and Security:** Ensuring the privacy and security of sensitive audit data is of utmost importance. Implementing robust data encryption, access controls, and anonymization techniques can protect confidential information while allowing effective fraud detection.
9. **Integration with External Data Sources:** Incorporating external data sources, such as public records, news feeds, and social media, can enrich the fraud detection process by providing additional context and insights into potential fraudulent activities.
10. **Scalability and Performance:** Optimizing the performance and scalability of AuditAI to handle large volumes of data and accommodate growing user demands is essential. Employing distributed computing, cloud-based infrastructure, and parallel processing techniques can enhance system efficiency.

**CHAPTER 11**  
**APPENDIX**

## SourceCode

### **Build PYTHON FLASK Code:**

```
from flask import Flask, request, render_template
import pickle

app = Flask(__name__)
model = pickle.load(open('knn.pkl', 'rb'))

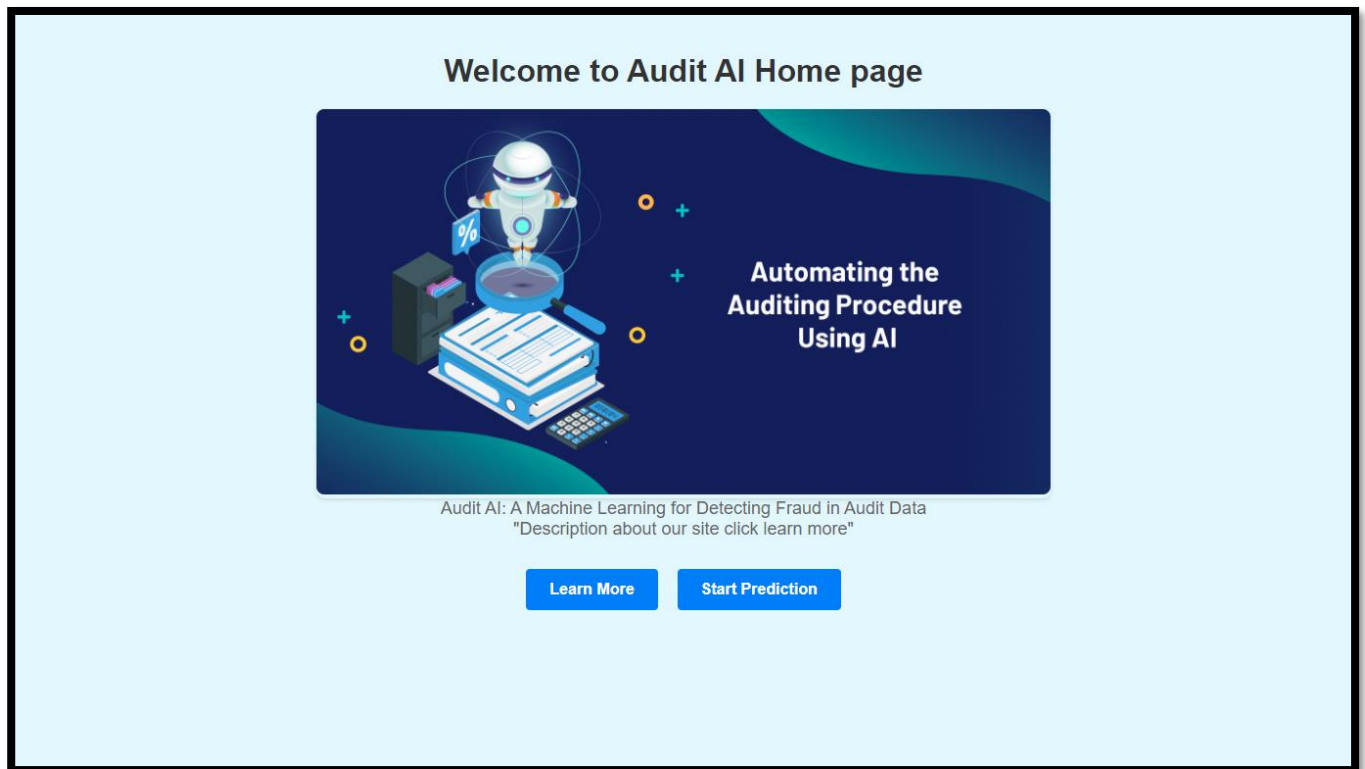
@app.route('/')
def home():
    return render_template('index.html')

@app.route('/project')
def project():
    return render_template('project.html')

@app.route('/riskpred', methods=['GET', 'POST'])
def predict_risk():
    if request.method == "POST":
        pred = [
            float(request.form['Sector_score']),
            float(request.form['PARAM_A']),
            float(request.form['Risk A']),
            float(request.form['PARAM_B']),
            float(request.form['Risk B']),
            float(request.form['TOTAL']),
            float(request.form['numbers']),
            float(request.form['Money_Value']),
            float(request.form['Score_MV']),
            float(request.form['District_Loss']),
            float(request.form['History']),
            float(request.form['Score']),
            float(request.form['Inherent Risk']),
            float(request.form['Audit_Risk'])
        ]
        output = model.predict([pred])[0]
        return render_template('result.html', predict="The Predicted Risk value is: " + str(output))
    else:
        return render_template('main.html')

if __name__ == '__main__':
    app.run(debug=True)
```

AUDIT AI home page:



Getting parameters from users

### Detecting Fraud in Audit Data

Sector Score:

PARA A:

Risk A:

PARA B:

Risk B:

TOTAL:

Numbers:

Money Value:



234

Money Value:

235

Score MV:

2345

District Loss:

234534

History:

324

Score:

2345

Inherent Risk:

2345

Audit Risk:

234

Submit

The result page

Result

The Predicted Risk value is: 1

**GITHUB LINK:**

<https://github.com/naanmudhalvan-SI/PBL-NT-GP--5717-1680798986>

**YOUTUBE LINK:**

<https://youtu.be/znhTCov1Tdc>