

Developing cloud based secure data monitoring system

A major project report submitted in partial fulfillment of the requirement for the
award of degree of

Bachelor of Technology

in

Computer Science & Engineering / Information Technology

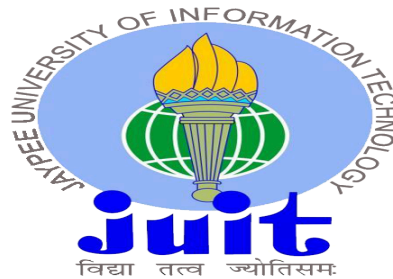
Submitted by

Surbhit Sharma(201521)

Akshat Sharma (201247)

Under the guidance & supervision of

Dr . Pradeep .K. Gupta



**Department of Computer Science & Engineering and
Information Technology Jaypee University of Information Technology,
Waknaghat, Solan - 173234 (India)**

CANDIDATE'S DECLARATION

We hereby declare that the work presented in this report entitled '**Developing cloud based secure data monitoring system**' in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering / Information Technology** submitted in the Department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology, Waknaghat is an authentic record of my own work carried out over a period from August 2023 to May 2024 under the supervision of **Dr Pradeep K. Gupta** (Professor, Department of Computer Science & Engineering and Information Technology).

The matter embodied in the report has not been submitted for the award of any other degree or diploma.

Student Name: Surbhit Sharma

Roll No.: 201521

Student Name: Akshat Sharma

Roll No.: 201247

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

Supervisor Name: Pradeep Kumar Gupta

Designation: Professor Department: Computer Science and engineering

Dated:

ACKNOWLEDGEMENT

Firstly, we express our heartiest thanks and gratefulness to almighty God for His divine blessing makes it possible for us to complete the project work successfully. We are really grateful and wish to be profoundly indebted to Supervisor Dr. Pradeep Kumar Gupta, Professor(SG), Department of CSE Jaypee University of Information Technology, Wakhnaghat. Deep Knowledge & keen interest of our supervisor in the field of “Machine/deep learning” to carry out this project. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior drafts and correcting them at all stages have made it possible to complete this project. We would like to express our heartiest gratitude to Dr.Pradeep Kumar Gupta Department of CSE, for his kind help to finish our project. We would also generously welcome each one of those individuals who have helped us straightforwardly or in a roundabout way in making this project a win. In this unique situation, we might want to thank the various staff individuals, both educating and non-instructing, which have developed their convenient help and facilitated our undertaking. Finally, we must acknowledge with due respect the constant support and patients of our parents

TABLE OF CONTENTS

CANDIDATE’S DECLARATION	I
ACKNOWLEDGEMENT	II
LIST OF ABBREVIATIONS	IV
LIST OF FIGURES	V
ABSTRACT	VI
1. Chapter-1 INTRODUCTION	1
1.1 Introduction	1
1.2 Problem Statement	2
1.3 Objectives	3
1.4 Methodology	5
1.5 Organization	7
2. Chapter-2 LITERATURE SURVEY	8
2.1 Literature Survey	8
3. Chapter-3 SYSTEM DEVELOPMENT	14
4. Chapter-4 PERFORMANCE ANALYSIS	37
5. Chapter-5 CONCLUSIONS	41
5.1 Conclusions	41
5.2 Applications Contributions	42
6. REFERENCES	43
7. APPENDICES	45

LIST OF ABBREVIATIONS

Relu	Rectified linear unit
RF	Random forest
CNN	Convolutional neural network
Numpy	Numerical Python
DHT	Distributed hash table
SK learn	scikit-learn
ML	Machine learning
API	Application programming interface
GUI	Graphic user interface
POR	Proof-of retrievability
PKI	Public key infrastructure

LIST OF FIGURES

Figure 1. DFD lvl 1

Figure 2. DFD lvl 2

Figure 3. Workflow of processed model

Figure 4. Code snippet of Support vector Machine

Figure 5. Code snippet of Random Forest

Figure 6. Code snippet of Imported libraries

Figure 7. CNN Architecture

Figure 8. Code snippet of combined dataset

Figure 9. Code snippet of cleaning and modified data

Figure 10. Code snippet of clean data

Figure 11. Code snippet of implementation of decision tree

Figure 12. Login page for user authentication

Figure 13. Dashboard of proposed model

Figure 14. Power BI dashboard snapshot

Figure 15. Exploring Employee Data with Slicer Search Bar: A Snapshot from Power BI Dashboard

Figure 16. Implementation and training of RF

Figure 17. Code snippet of training decision tree

Figure 18. Code snippet of SVM's accuracy

Figure 19. Heatmap of dataset

Figure 20. Snippet of MongoDB terminal

Figure 21. Snapshot of Machine Learning Model Accuracy Evaluation

ABSTRACT

In the contemporary world, when data and digital technologies are the foundation of everything, monitoring a company's financial expenditures is crucial. Imagine a system that analyzes and comprehends employee monthly spending by using the power of technology, particularly cybersecurity (which defends computer systems against threats) and machine learning (a form of artificial intelligence that learns from data).

These days, protecting information is one of the main issues. Numerous individuals attempt to breach computer systems in order to take advantage of crucial data, particularly financial data.

We have worked very hard to ensure our system is extremely secure in order to fight this.

We have developed an intelligent system by merging the capabilities of machine learning with cybersecurity. It looks at more than just the amount of money spent; it also ensures that personal financial information is protected from outside dangers. This creative method establishes a new benchmark for safe, effective, and environment-specific spending management for today's hectic business world.

CHAPTER -01 INTRODUCTION

1.1 Introduction

In the current digital age, when organizations heavily rely on data, it is imperative that these data be monitored through dependable, secure, and efficient methods. The technology in question is an intricate fusion of cybersecurity and machine learning, two fundamental technologies. Its main objective is to provide a thorough grasp of the monthly costs incurred by workers in an organization.

Strong security focus is the cornerstone of our system's design philosophy. The ever-growing complexity and constant change in cyber threats and data breaches puts organizations at serious risk when it comes to protecting their confidential financial information. Our strategy gives the deployment of strict security measures first priority in order to mitigate these threats.

Advanced cybersecurity protocols are integrated into our system to provide defense against possible threats [1]. Sensitive information is encoded using encryption techniques to render it unreadable for unauthorized parties or cyberattackers. Furthermore, carefully designed access controls guarantee that only individuals with permission can access particular data. Mechanisms for continuously detecting threats are in place to quickly spot any questionable activity or possible breaches. As essential components of our security structure, real-time reaction tactics and routine monitoring allow us to take prompt action in the event that any anomalies or threats are discovered.

Furthermore, adding machine learning capabilities to our system gives it a dynamic layer[1]. It strengthens our security measures and makes it easier to analyze enormous amounts of financial data. The training of machine learning algorithms to identify patterns and abnormalities in data helps identify any security risks or irregularities in employee spending early on.

Our system's ultimate goal is to provide a safe and effective framework for thoroughly tracking and evaluating employee spending. Our goal is to redefine data protection and analytics by combining the adaptability of machine learning with strong cybersecurity protocols[1]. Our goal is to meet the complex needs of a data-driven business environment while maintaining the privacy, accuracy, and security of financial data.

1.2 Problem Statement

In the contemporary corporate landscape, seamless connectivity through SIM cards equipped with unrestricted internet access has become a norm, enhancing communication and productivity. However, this practice has inadvertently led to a critical challenge. Despite the availability of connectivity, employees often misuse this resource, resulting in significant financial strain for organizations [1]. This habitual mishandling of data translates into unnecessary expenses, potentially affecting profitability, budget allocations, and growth prospects. These unwarranted data expenditures have ramifications beyond inconvenience, resonating through an organization's financial structure. Decreased profitability, compromised budgets, and hindered growth prospects are potential consequences, necessitating innovative solutions that align connectivity provisions with prudent financial management. In response, a pioneering Python-based model has been meticulously developed. This sophisticated analytical tool delves into the intricate landscape of data usage patterns, aiming to equip decisionmakers with comprehensive insights into organizational data consumption. By doing so, it empowers leaders to curtail extravagant internet usage, thereby relieving the strain on financial resources. Beyond its primary role in managing data expenses, this model offers solutions for various challenges inherent in contemporary data management. Notably, it upholds data privacy through stringent encryption mechanisms and prevents unauthorized access via robust authentication protocols. The proposed system's distinctive feature is its scalability, deftly navigating diverse organizational structures and seamlessly integrating with existing infrastructures. This adaptability, combined with the model's capability to handle diverse data formats and sources, positions it as an indispensable asset for organizations navigating the fluid digital landscape. However, this model transcends being a mere problem-solving tool. It aspires to redefine data monitoring, merging technological prowess with nuanced organizational understanding to become a transformative force. Its applicability spans various industries,

uniting them under streamlined data governance. Furthermore, the model's adaptability to evolving regulatory and industry standards underscores its forward-looking design. Amid perpetual change, it instills confidence in organizations to oversee their data assets with unwavering assurance. In essence, this model seeks to remedy the challenge of reckless data usage, enabling organizations to strike a harmonious balance between connectivity and financial prudence.

1.3 Objectives

- **Enhanced Security Measures:** By putting strict security measures in place, the main goal is to protect the system from possible cyberattacks and data breaches[1]. To protect sensitive financial data, this entails strong encryption methods, access restrictions, threat detection, and real-time monitoring.
- **Detailed Expense Analysis:** Create models and algorithms that make use of machine learning capabilities in order to analyze and interpret employee spending in great detail[1]. The goal is to give the business comprehensive insights into spending trends, anomalies, and patterns so that it can make smarter financial decisions.
- **Real-time Monitoring and Threat Detection:** Put in place systems that keep an eye on data activity all the time, and use machine learning techniques to quickly spot any anomalies or questionable activity. Enabling quick reactions to possible security threats or unauthorized access attempts is the aim.
- **Adaptive machine learning :** it involves teaching models to change and grow in response to fresh trends or new risks in labor costs. One of the most important goals is to gradually increase the system's accuracy and capacity to identify abnormalities.
- **User-Friendly Interface:** Create an easy-to-use interface that makes it possible for stakeholders to obtain and comprehend spending data. ensuring that authorized personnel may explore the system with ease and accessibility and obtain useful insights.

- **Scalability and Efficiency:** Create a system design that can manage massive data volumes with ease and still deliver excellent performance[1]. Make that the system is scalable to handle future expansion and rising data volumes without sacrificing system performance.
- **Regulations and Compliance:** Verify that the system complies with all applicable laws, industry standards, and data protection guidelines. This entails upholding data confidentiality, integrity, and compliance with privacy laws like HIPAA and GDPR.
- **Continual Adaptability and Improvement:** Set up procedures for ongoing system assessment, feedback incorporation, and enhancement[1]. To preserve system efficacy, allow for flexibility to change cybersecurity threats and developments in machine learning technology.
- **Cost-effectiveness:** Provide top-notch security and analytics capabilities while minimizing resources and expenses related to system setup and maintenance.

All of these goals work together to build a strong, safe, and effective system that analyzes employee spending and places data security and integrity at the forefront of the business's operating structure.

1.4 Methodology

Developing a process for combining machine learning and cybersecurity to analyze employee spending entails a number of crucial actions and strategies. The approach is outlined as follows:

Recognising Data Sources and Gathering

Find and compile pertinent sources of information about employee spending in the organization. This could include invoices, transaction data, spending reports, etc. Make sure the

information gathered complies with all applicable privacy and regulatory compliance requirements[2].

Preparing data:

Make sure the gathered data is accurate, consistent, and prepared for analysis by cleaning and preprocessing it. This covers formatting for machine learning compatibility, normalization, and managing missing values.

Development of Machine Learning Models:

Select appropriate machine learning algorithms according to the type and intricacy of the spending information. Teach machine learning models to evaluate and decipher trends, anomalies, and spending patterns[2]. Regression analysis for forecast analysis, anomaly detection for finding irregularities, and classification algorithms for classifying spending are a few examples of models.

Selection and Feature Engineering:

Extrapolate significant characteristics from the data that support cost analysis. In order to improve the accuracy and efficiency of the machine learning models, identify the most pertinent qualities using feature selection approaches[2].

Development of User Interfaces:

Provide a simple, easy-to-use interface so that stakeholders can safely access and understand spending analytics[2]. Ascertain that permissions and secure access controls are in place to restrict access to just authorized individuals. we have shown the data flow diagrams as shown Figure1, Figure 2.

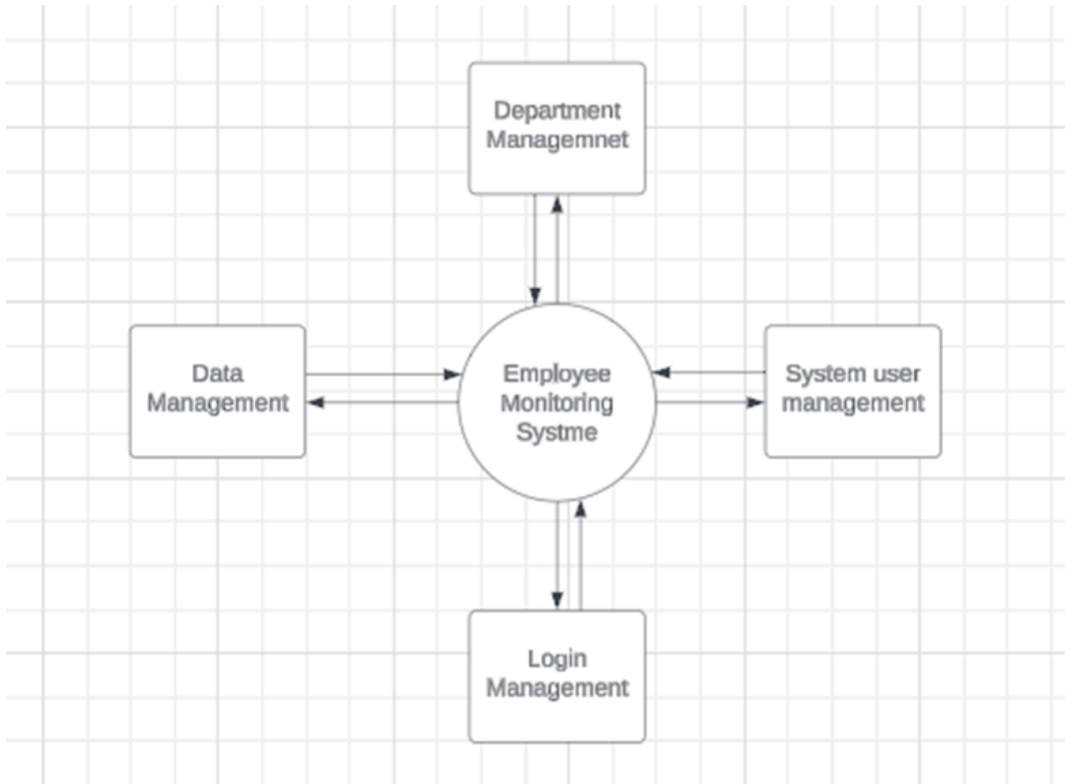


Figure 1. DFD lvl 1

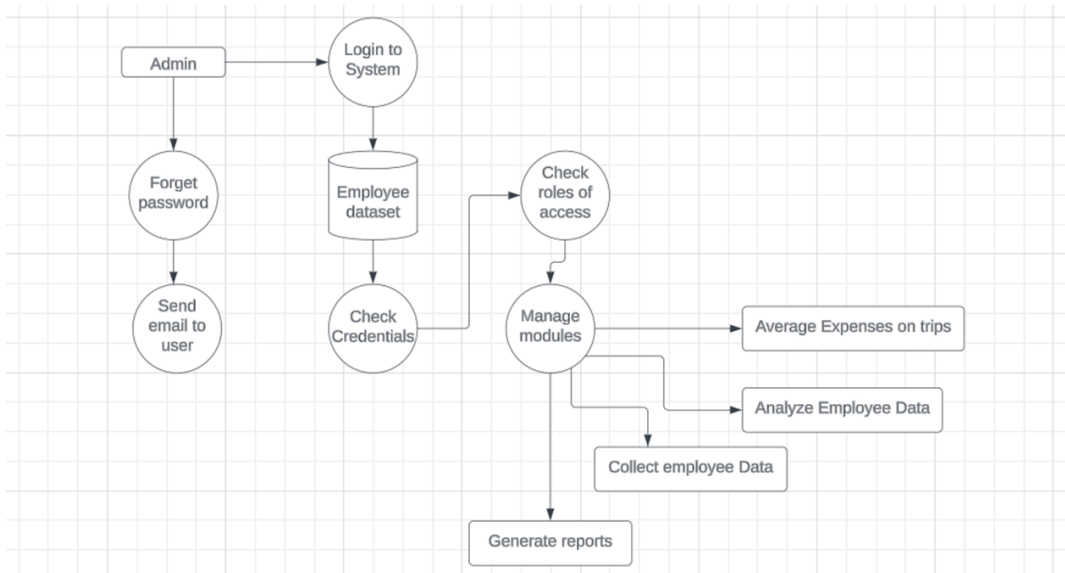


Figure 2. DFD lvl 2

1.5 Organization of the Report

Our report comprises of the chapters and has been divided as follows:

Chapter 2 gives an overview of the literature study performed.

Chapter 3 discusses the system development and workflow.

Chapter 4 shows the performance analysis..

Chapter 5 highlights the conclusion, future scope and application contribute

CHAPTER -02 LITERATURE REVIEW

Integrating several elements to effectively track and manage employee internet connections and expenses is necessary when building a data monitoring system based on automated human resources. The parts and their potential functions are broken out as follows:

Monitoring of Internet Connections:

Tracking Usage: Tracking Usage: Throughout business hours, the system keeps track of how much data, which websites employees visit, and how much time they spend online[3]. This helps with resource optimisation and educated decision-making by offering insightful information about both individual and group online behavior.

Security Measures: By actively spotting and reporting possible security lapses or unauthorized online activity, the system serves as a proactive security sentinel. This entails spotting erroneous data transfers, strange internet visits, and any attempt to breach established security procedures[3].

Performance analysis: To find bottlenecks and opportunities for improvement, the system examines metrics related to network performance, such as internet speeds and connectivity problems. This makes it possible to implement focused interventions to guarantee that every employee has the best possible internet performance, which boosts output and enhances user satisfaction.

Managing Expenses:

Expense tracking: Expense tracking: The system keeps a thorough record of all employee costs in a number of areas, such as travel, housing, food, and other company-provided allowances[3]. Making educated financial decisions is made possible by this centralized platform, which offers clear visibility into employee spending trends.

Verification of Recipient: By utilizing integrated receipt scanning technology, the system reduces the need for manual data entry and improves accuracy by automatically verifying and

classifying expenses[3]. This simplifies cost reporting procedures and lowers administrative overhead.

Budgeting and Reporting: With the system's ability to provide thorough reports and data visualizations, management is better equipped to track budget adherence, examine spending patterns, and make informed decisions. Informed resource allocation is encouraged, and total financial performance is maximized.

Alerts and Automation:

Automated Processes: By minimizing manual intervention and cutting down on processing time, the use of automated workflows optimizes the approval process for online and expenditure spending. This promotes cost savings, openness, and efficiency.

Alert systems: Predefined thresholds or anomalies, such as surpassing spending caps or atypical online behavior, set off proactive alerts[3]. This guarantees prompt reporting of any fraud, policy infractions, or important decision-making to pertinent staff.

Adherence to and Implementation of Policy: Continuous monitoring and enforcement of company policies and guidelines governing employee internet usage are key to maintaining a secure and compliant environment. This promotes responsible digital behavior and mitigates potential risks.

Regulatory Compliance: The system adheres to all relevant laws and regulations regarding cost reporting, data privacy, and online security[3]. This ensures legal and ethical business practices while safeguarding sensitive information.

Data insights and analytics:

Data visualization is the process of presenting data in clear, understandable visual formats, such as graphs or dashboards, to reveal trends in internet usage and spending patterns. Predictive analysis is the process of forecasting future costs or streamlining internet usage to save money[3].

Scalability and Integration:

Integration with HR Systems: Establishing connections with current HR databases and systems to guarantee correctness and expedite employee data.

Scalability: Building the system to support the expansion of the business by adding additional personnel and adding more functionalities as needed[3].

Table 1. Literature review

S. No.	Paper Title	Journal/Conference (Year)	Tools/Techniques/Dataset	Results	Limitations
1.	Research on data security technology based on cloud storage	13th Global Congress on Manufacturing and Management, GCOMM 2016	Symmetric encryption, Erasure codes, Proof-of retrievability (POR), , data label verification, replica strategy, height of authentication, attribute encryption, time encryption, and DHT network.	This paper describes a data secure storage scheme based on Tornado codes (DSBT) that is designed to address the problems of data availability, confidentiality, and recovery in cloud storage systems	<ul style="list-style-type: none"> • It is limited to the use of Tornado codes. • It does not address the issue of data recovery from malicious behavior of the cloud service provider.
2.	Addressing cloud computing security issues	Future Generation Computer Systems 3, March 2012.	Public Key Infrastructure (PKI), Single Sign-On (SSO), Lightweight Directory Access Protocol (LDAP)	It proposes a solution to the security challenges of cloud computing based on a Trusted Third Party (TTP) that uses cryptography, Public Key	<ul style="list-style-type: none"> • Reliance on a TTP • Use of PKI • Use of SSO and LDAP

				Infrastructure (PKI), Single Sign-On (SSO), and LDAP.	
3.	Exploring human resource management intelligence practices using machine learning models	Journal of High Technology Management Research, 2023	Natural Language Processing (NLP), Machine learning algorithms, Data mining, Statistical analysis	Machine learning (ML) can help HR professionals to see patterns or trends in HR data, which can then be used to inform better strategy. It is important to use ML in an honest and open manner, and to check that no discriminatory results are formed.	<ul style="list-style-type: none"> • It does not discuss the ethical implications of using machine learning for recruitment. • It does not provide any concrete recommendations for how organizations can successfully implement ML for recruitment.
4.	BAMHealthCloud: A biometric authentication and data management system for healthcare data in cloud	Journal of King Saud University, 2020	ALGOHealthSecurityCheck, MapReduce programming model. Different smartphones, tablets, phablets and PDAs are used for acquiring data.	The results obtained from the proposed framework are better than the other approaches. BAMHealthCloud achieves an EER of 0.12, sensitivity of 0.9.	<ul style="list-style-type: none"> • In different levels of biometric security there may be some threat. • It contains lots of mathematical operations.

5.	Secure Data Sharing and Analysis in Cloud-Based Energy Management Systems	ICST Institute for Computer Sciences, Social Informatics and Telecommunications Engineering, 2018	Cloud based local energy management Systems, Iot devices at the network edge, The LEMS algorithm and GUI.	In this paper we have given an overview of the cloud based energy management system using live examples of cloud based LEMS. The aim is to reduce the aggregated energy demand. Also, security is the main concern. So for the different attacks involved we have the proper measures.	<ul style="list-style-type: none"> • Data leakage, • Spoofing, • Disruption of services, • Energy Bleeding, • Hardware issues
6.	An Analysis of the Cloud Computing Security Problem	Journal of Cornell University, 2016	The cloud computing model depends on a deep stack of dependent layers of objects. IaaS Model, PaaS Model.	Cloud security must address inherited technology problems and multi-tenancy complexities	<ul style="list-style-type: none"> • PaaS Security Issues • SaaS security Issues • IaaS security

CHAPTER - 03 SYSTEM DEVELOPMENT

3.1 Applications

Cost Management and Optimization:

- **Expense Control:** Organizations can effectively control costs and optimize budgets by implementing effective practices for tracking and managing expenses [3]. This proactive strategy improves organizational sustainability by fostering a culture of financial accountability and responsibility.
- **Finding Cost-saving Opportunities:** Organizations can find hidden places where cost-saving solutions can be adopted by methodically analyzing spending patterns. Targeted cost reduction measures are made possible by this data-driven approach, which enhances financial efficiency and resource allocation[4].

Employee Productivity and Performance:

- **Internet Usage Insights:** Examining how employees use the internet can reveal important information about possible productivity snags and areas in need of further funding or instruction[4]. Organizations can improve employee performance and overall organizational productivity by optimizing internet policy and infrastructure by identifying websites, programmes, and online activities that take up a substantial amount of time or resources.
- **Performance Monitoring:** Ensuring a constant and ideal working environment for staff members requires constant monitoring of internet connectivity and speed. By detecting and resolving network problems early on, this proactive strategy reduces interruptions and facilitates effective online workflows.

Policy Adherence and Compliance:

- **Ensuring Policy Compliance:** Keeping an eye on internet usage makes it easier to enforce rules on appropriate resource use. This involves following security guidelines, such as those governing data protection and the complexity of passwords. It also

guarantees adherence to pertinent regulatory requirements, reducing legal risks and maintaining organizational integrity[4].

- **Expense Policy Enforcement:**Real-time expense verification against company policies safeguards against fraudulent claims and ensures adherence to established spending guidelines[4]. This promotes fiscal responsibility and prevents unnecessary financial burdens on the organization.

Security and Risk Mitigation:

- **Detecting Security Threats:** Continuous monitoring of internet activity serves as a crucial early warning system for potential security threats. This includes identifying anomalies indicative of unauthorized access attempts, malware infections, or data exfiltration efforts[4]. By proactively detecting such threats, the system facilitates timely interventions and mitigates potential damage.
- **Expense Fraud Detection:** Analyzing expense data for unusual patterns or inconsistencies can expose fraudulent activities like expense padding, duplicate claims, or unauthorized spending. By uncovering such irregularities, the system empowers organizations to take swift action to minimize financial losses and implement preventative measures against future fraud attempts.

Data-driven Decision Making:

- **Strategic Planning:** Data-driven insights gleaned from expense and internet usage patterns can empower strategic decision-making, optimize resource allocation, and inform future investment plans.
- **Predictive Analysis:**Leveraging historical data to forecast future expenses enables proactive budget planning and precise financial projections[4].

Enhanced Transparency and Accountability:

- **Clear Reporting:**Generating comprehensive reports and visualizations facilitates both internal and external accountability in expense and resource utilization. This transparency improves internal decision-making and provides readily accessible data for audits or stakeholder inquiries.

- **Fair Expense Policies:** Transparent expense management policies foster employee trust and satisfaction through fair and timely reimbursements[4]. This fosters a culture of accountability and encourages responsible spending habits.
- **Resource Availability:** Ensuring readily accessible and reliable internet resources directly impacts employee morale and engagement. This minimizes frustrations, improves productivity, and promotes a positive working environment.
- **Feedback Loop:** Insights gleaned from ongoing monitoring fuel continuous refinement of policies, training programs, and infrastructure. This iterative process optimizes efficiency, effectiveness, and cost-savings potential[4].

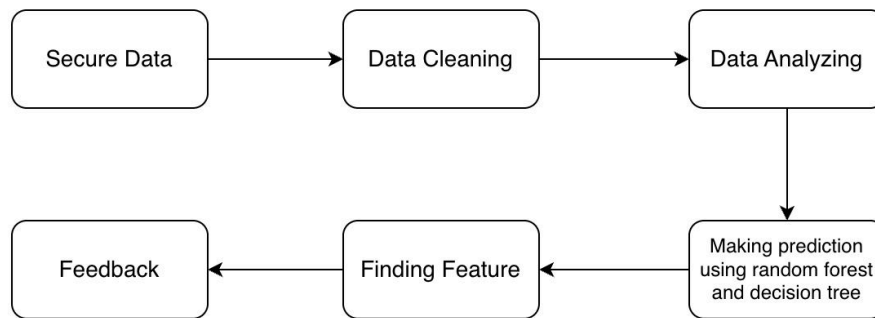


Figure 3. Workflow of proposed model

3.2 Proposed work

This proposed workflow Figure 3 provides a structured approach to implementing a comprehensive monitoring system while emphasizing continuous improvement and alignment with organizational goals and policies. Adjustments can be made based on specific organizational needs and the nature of the workforce.

Phase 1: Technology and System Implementation

Install software and monitoring tools to keep tabs on internet usage and handle expenses. For smooth data flow, integrate these solutions with the current HR and finance systems. Configuring Data Gathering and Analysis. Set up procedures to gather and examine records of expenses and internet usage. Provide systems for classifying and verifying receipts. Measures

for Security and Compliance Put security measures in place to guarantee data protection and adherence to laws (such as GDPR and HIPAA). Provide data privacy and usage policy training to employees.

Phase 2: Testing, Activation, and Monitoring of the System

For the first testing, turn on the monitoring system at a safe location. To make sure the system is working and the data is accurate, conduct trial runs. **Constant Observation and Evaluation** Keep an eye on real-time spending data and internet usage on a regular basis. Make use of dashboards and reports to examine compliance levels, trends, and anomalies. **Notifications and Alerts** To take prompt action, set up alerts for anomalous internet usage or spending trends. Create procedures for responding to alarms and resolving problems that are found.

Phase 3: Assessment and Enhancement

Performance Assessment Evaluate the monitoring system's performance in relation to predetermined goals. Get input on the effect and usefulness of the system from stakeholders and staff. **Enhancement and Streamlining** Make the required adjustments in light of the evaluation's comments. Optimize policies, procedures, or algorithms to improve efficiency and accuracy. **Instruction and Interaction** Hold recurring training sessions for staff members to help them grasp spending policies and make efficient use of resources. Share best practices, updates, and modifications to the system with all pertinent parties.

3.3 Module and implementation

SVM(Support vector machine)

Support vector machine, or SVM, is a potent machine learning technique that is employed for a variety of applications. This model of learning is supervised. SVMs are trained using labeled data, which is data in which each of the samples has a label or category assigned to it[13].

Among its fundamental ideas are:

- **Hyperplane:** A hyperplane is a decision boundary that divides the data points into the appropriate classes. It looks like a line in 2D or a plane in 3D[13].

- **Margin:** The distance, known as the support vectors, between the hyperplane and the nearest data points from each class is referred to as the margin. To achieve the optimal separation, SVMs try to maximize this margin[13].
- **Support Vectors:** These represent the margin and are the data points that are closest to the hyperplane. They are essential to the SVM's ability to make decisions.

Applications of Vector Machine Support

- It is mostly employed for categorization SVM is very good at classifying data points into distinct groups. It does this by identifying the decision boundary that best splits the data, which is frequently a hyperplane in higher dimensions.[13]
- It is suitable for regression as well: Although they are less popular, SVMs can also be used to foresee a continuous value in regression issues.

Application of SVM

- SVMs work well with high-dimensional data since they don't significantly reduce performance when dealing with large amounts of features[13].
- Excellent performance on complicated datasets: SVMs can perform well even when dealing with non-linearly separable data by employing kernel tricks, which are methods for transforming data into a higher dimension.
- Robust against outliers: Compared to certain other algorithms, SVMs are less vulnerable to the impact of outliers in the data.[13]
- It helps in text categorization which includes things like labeling emails as spam or not, topic-based news article classification, and more.
- Also used in Image classification which includes things like facial recognition and object recognition in photos.
- It helps in anomaly detection which is finding anomalous patterns in data that diverge from the norm.[14]

Our study describes the application of a Support Vector Machine (SVM) model for forecasting employee turnover. We employed the SVM Model. The model attempts to categorize workers according to their data as either likely to stay with the company or leave (turnover).

Data Preparation

The Python code builds models and manipulates data using tools like scikit-learn and pandas.

- **Data Loading:** The code expects that employee data is contained in a pandas DataFrame named data. During implementation, the actual data source must be used in place of this DataFrame.
- **One-Hot Encoding:** We used `pd.get_dummies`, categorical characteristics such as Education, City, Gender, and EverBenched are transformed into one-hot encoded features. Since SVM models generally perform better with numerical information, this translation is required.
- The feature separation divides the data into the target variable (Y) and its features (X). With the exception of the goal variable, which indicates whether or not the employee left the company (1) or (0), the characteristics represent all employee data.
- **Train-Test Split:** We used the `Train_test_split` to divide the data into training and testing sets. The model is trained on the training set (80%), and its performance on unseen data is assessed on the testing set (20%). It is set to a random state for reproducibility.
- **Preprocessing:** The Mean technique is used when handling missing values with `SimpleImputer`. This substitutes the related feature's mean value for any missing entries.
- **Feature Scaling:** To guarantee that every feature has a comparable range, features are scaled using `StandardScaler`. This can enhance the SVM model's performance by evenly allocating weight to each feature throughout training.

Model Training

The parameters which we used to define an SVM classifier are as follows:

- **Radial Basis Function (RBF) kernel:** The model can capture intricate correlations between features because it is designed for non-linear decision boundaries.
- **Parameter for Regularisation (C):** To manage the trade-off between training error and model complexity, C is set to 1.0.

- **Gamma:** In order to pick the kernel coefficient automatically based on the data, Gamma is set to scale.

The model is then trained on the preprocessed training data (Xtrain_scaled, Ytrain).

Model Evaluation

- **Accuracy:** Accuracy_score is used to compute the model's accuracy. This measure shows the overall percentage of accurate predictions the model made using the test set.
- **Classification Report:** The classification_report function offers a thorough analysis of the model's performance for every class (leaving or remaining employees). Typically, this report contains metrics such as:

Precision: The percentage of employees who were actually positive (left) out of those who were expected to be positive (predicted to quit).[14]

Recall: The percentage of accurately anticipated actual positives (i.e., employees who actually left).[14]

F1-Score: A single score derived from the harmonic mean of precision and recall.

Support: The total number of actual cases for every class (the sum of the departing and remaining employees).

We can learn more about the model's capacity to recognise departing employees and how well it keeps individuals who are most likely to stay in the workforce by examining these indicators mentioned above in Report.

```

from sklearn.model_selection import train_test_split
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report

# Load your dataset
# Assuming your dataset is stored in a DataFrame named 'data'
# Replace 'data' with the actual name of your DataFrame
# data = ...

# Convert categorical variables to one-hot encoded features
data_encoded = pd.get_dummies(data, columns=['Education', 'City', 'Gender', 'EverBenched'])

# Separate features (X) and target variable (Y)
X = data_encoded.drop(columns=['LeaveOrNot'])
Y = data_encoded['LeaveOrNot']

# Split the dataset into training and testing sets
Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, Y, test_size=0.2, random_state=2)

# Data preprocessing
imputer = SimpleImputer(strategy='mean')
scaler = StandardScaler()

Xtrain_imputed = imputer.fit_transform(Xtrain)
Xtrain_scaled = scaler.fit_transform(Xtrain_imputed)
Xtest_imputed = imputer.transform(Xtest)
Xtest_scaled = scaler.transform(Xtest_imputed)

# Train the SVM classifier
SVM = SVC(kernel='rbf', C=1.0, gamma='scale')
SVM.fit(Xtrain_scaled, Ytrain)

# Make predictions on the test set
predicted_values = SVM.predict(Xtest_scaled)

# Calculate accuracy
accuracy = accuracy_score(Ytest, predicted_values)
print("SVM's Accuracy is:", accuracy*100)

# Print classification report
print(classification_report(Ytest, predicted_values))

```

Figure 4. Code snippet of Support Vector Machine

Random forest

The Random Forest model Figure 4 is an ensemble learning technique that builds several decision trees during training and produces a class that represents the mean prediction (regression) or mode of the classes (classification) of the individual trees[4].

Trees of Decisions:

Basic Building Blocks: The cornerstones of a Random Forest are decision trees. In the Random Forest, every tree is constructed separately and functions as a different classifier or regressor.

Node splitting: Decision trees divide based on features at each node to form branches, which ultimately lead to leaf nodes, which are the prediction-making nodes[4].

Decision Criteria: The characteristic that offers the best data separation is used to produce the divides. For classification, the criteria that are frequently used are Gini impurity, and for regression, mean squared error.

Ensemble Random Forest:

Bagging (Bootstrap Aggregating): Random Forest employs a method known as bootstrap aggregating, or bagging. By sampling with replacement, or bootstrapping, it generates numerous subsets of the training data, each of which is used to train a different decision tree.

Random Feature Selection: At each split in the process of building a single tree, a random subset of features is taken into account. By adding variation and decorrelating the trees, this randomness lessens overfitting and enhances generalization[5].

Voting or Averaging: In classification problems, the Random Forest model's Figure 5 final prediction is decided by a majority vote among each individual tree's predictions. It's usually the mean prediction of all the trees for regression tasks.

Benefits of Random Forest:

Reduction of Overfitting: By mixing numerous decision trees and employing random feature subsets, Random Forest has the tendency to reduce overfitting in comparison to individual decision trees[5]. Without requiring a lot of hyperparameter tweaking, it frequently offers good accuracy and generalization on unseen data. It takes Care of Big Datasets It performs well with missing values and can manage big, highly dimensional datasets.

```
from sklearn.ensemble import RandomForestClassifier

RF = RandomForestClassifier(n_estimators=20, random_state=0)
RF.fit(Xtrain,Ytrain)

predicted_values = RF.predict(Xtest)

x = metrics.accuracy_score(Ytest, predicted_values)
acc.append(x)
model.append('RF')
print("RF's Accuracy is: ", x)

print(classification_report(Ytest,predicted_values))
```

Figure 5. Code snippet of Random Forest

```
# Importing libraries

from __future__ import print_function
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import classification_report
```

Figure 6. Code Snippet of Imported Libraries

CNN (Convolutional neural network)

Convolutional neural networks, or CNNs for short, are a particular kind of artificial neural network that function well for tasks involving image recognition and classification. CNNs are built with the ability to automatically and adaptably identify feature spatial hierarchies from input data[6].

Important CNN Components:

1. Convolutional Layer

(a)Feature extraction: The convolutional layer is the fundamental unit of a CNN as shown in Figure 6[6]. It extracts different features from the input data by applying a collection of learnable filters, or kernels.

(b)Feature Maps: These filters conduct element-wise multiplication and summation as they slide (convolve) over the input image, creating feature maps that illustrate various learnt features.

(c)Layers of Pooling: Following convolutions, each feature map's spatial dimensions are decreased while crucial information is preserved using pooling layers, such as max-pooling.

2. Activation Function:

(a)Non-linearity: To introduce non-linearity and uncover intricate patterns in the data, activation functions such as ReLU (Rectified Linear Unit) are typically applied after convolution[6].

3. Completely Connected Layers: Convolutional layers extract features, which are then used by fully connected layers at the end of the network to carry out classification or regression tasks. Before being fed into fully connected layers, the output from convolutional and pooling layers is flattened. Backpropagation is the training method used for CNN[6]. The network is optimized to minimize the error between the expected and actual outputs by calculating gradients as they flow through it to update the weights and biases.

CNN's benefits

Automatically recognising the spatial hierarchies of features, CNNs are capable of identifying intricate patterns, textures, and edges. Parameter Sharing: Convolutional layers with shared weights have fewer parameters than fully connected networks, which increases CNN's computational efficiency. Translation Invariance: CNNs are resilient to changes in the location and orientation of observed features because of weight sharing and local receptive fields.

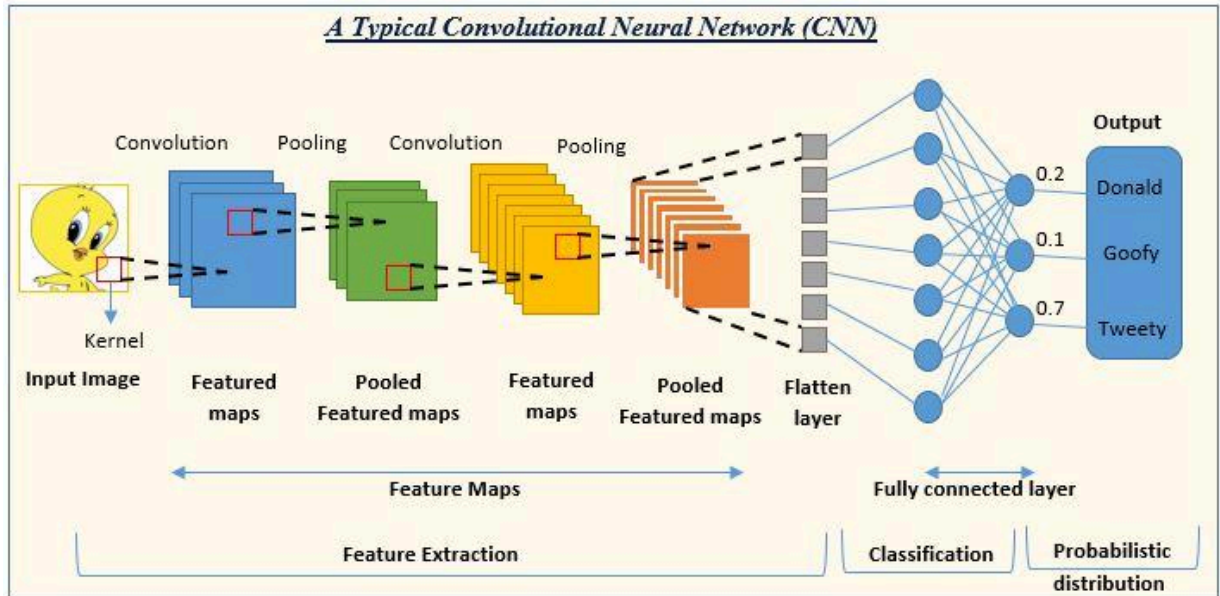


Figure 7. CNN Architecture[6]

3.4 Workflow of proposed model

- **Data creation**

We have created our datasets by using an online github repository and other online platforms like kaggle etc.

1. Repositories for Data Identification and Discovery:

We will explore relevant government data portals and academic repositories to identify publicly accessible datasets pertaining to the desired research area. This may include resources like the UCI Machine Learning Repository, institutional repositories, and subject-specific data archives. We used Kaggle's search feature to locate datasets associated with particular subjects, sectors, or machine learning contests[7]. This platform offers a vast collection of datasets with varying degrees of complexity and structure, allowing for targeted exploration.

2. Gathering and Combining Data:

(a)Data Collection: We have collected our data with the help of above resources and connected them to make one.

(b)Data Aggregation: We have obtained information from many sources and combined it into a structured manner that may be used for activities involving analysis or machine

learning[7]. To guarantee consistency, this entails cleaning, combining, and formatting data. The final data received after entailment has been shown in Figure 7.

```
merged_df = pd.merge(data_1, data_2, on=common_key)
✓ 0.0s Python
```

```
merged_df.head()
✓ 0.0s Python
```

	EEID	January	February	March	April	May	June	July	August	September	...	Bonus %	Country	City	Exit Date	week_1	week_2	week_3	week_4	Expenses on trips	Sum_Columns
0	E02387	0	0	0	0	0	0	0	0	0	...	15%	United States	Seattle	10/16/2021	4	1	6	8	66773	19
1	E04332	0	0	0	0	0	0	0	0	0	...	0%	United States	Miami	5/20/2021	9	2	1	7	70137	19
2	E04332	0	0	0	0	0	0	0	0	0	...	0%	United States	Miami	5/20/2021	9	2	1	7	70137	19
3	E03496	0	0	0	0	0	0	0	0	0	...	0%	United States	Austin	3/9/2020	2	2	2	4	67826	10
4	E04732	0	0	0	0	0	0	0	0	0	...	0%	United States	Chicago	4/22/2006	8	9	4	2	78326	23

5 rows x 33 columns

Figure 8. Code snippet of combined dataset

- **Data Cleaning & processing:** Our next step is cleaning and analyzing data Figure 8 using pandas and numpy
- **NumPy:** The core package for numerical computations in Python is called NumPy (Numerical Python). Multidimensional arrays and mathematical operations on these arrays are supported.
- **NumPy Data Cleaning Tasks:**
 - Managing Value Missing:** Tools for handling missing or NaN (Not a Number) values in arrays are provided by NumPy. Arrays can be checked for NaN values using `np.isnan()`.
 - Sorting & Filtering:** we have sorted our data on the basis of Expenses on trips and excessive usage of the internet in descending order[7].
 - Operations in Mathematics:** NumPy allows one to perform element-by-element mathematical operations on arrays, including division, multiplication, addition, and subtraction. The mean, median, and standard deviation can be computed using the `np.mean()`, `np.median()`, and `np.std()` functions to calculate the the sum of internet usage of months we have used the above expressions to finalize our dataset

- **Pandas:** Based on NumPy, Pandas is a robust library that provides tools for manipulating and analyzing data as well as data structures (such Series and DataFrame).

(a)Managing Missing Data:

Methods like `data.isnull()` and `data.dropna()` are provided by Pandas to detect and manage missing data in DataFrames[7]. We have used these expressions to calculate the null values and drop them because many ml models does not work on null values in the other hand we can all use `fillna()` to fill our null values `Fillna()` is a techniques like forward-fill or backward-fill to replace missing values with a supplied value.

(b)Data Selection and Filtering: While selecting our columns to sum all the months usage we use `.loc` and `.iloc` to select specific columns.

(c)Combining and Joining DataFrames: We have used `pd.concat()` to combine multiple DataFrames to make one dataset for analysis and `pd.merge()` to merge two datasets as we can see in Figure 8.

```
common_key = 'i»EEID'.strip('i»')

data_1.rename(columns={'i»EEID': 'EEID'}, inplace=True)

columns_to_drop = ['week_1', 'week_2', 'week_3', 'week_4', 'February', 'November']
data.columns = data.columns.str.strip()
data = data.drop(columns=columns_to_drop)
```

Figure 9. Code snippet of cleaning & Modifying data

```

data['March'] = np.random.randint(8, 40, len(data))
data['April'] = np.random.randint(10, 35, len(data))
data['May'] = np.random.randint(11, 36, len(data))
data['June'] = np.random.randint(12, 40, len(data))
data['July'] = np.random.randint(8, 37, len(data))
data['August'] = np.random.randint(10, 38, len(data))
data['September'] = np.random.randint(11, 33, len(data))
data['October'] = np.random.randint(9, 38, len(data))

data['December'] = np.random.randint(6, 39, len(data))

```

```

data['week_1'] = np.random.randint(0, 10, len(data))
data['week_2'] = np.random.randint(0, 10, len(data))
data['week_3'] = np.random.randint(0, 10, len(data))
data['week_4'] = np.random.randint(0, 10, len(data))
data['Expenses on trips'] = np.random.randint(10000, 100000, len(data))

```

```

data.head()

```

Figure 10. Code snippets of Cleaned Data

Implementation of decision tree

The provided code demonstrates the implementation of a Decision Tree classifier using the scikit-learn library in Python for predictive analysis. Here's a breakdown of the steps:

(a)Data Preprocessing: The missing values in the training dataset Xtrain are handled using the SimpleImputer from scikit-learn. The strategy employed for imputation is to fill missing values with the mean of the respective feature. Xtrain_imputed stores the transformed training dataset after imputation[8].

(b)Imputation on Test Data: The same imputer that was fitted on the training set is used to transform the test dataset Xtest (Xtest_imputed). This ensures consistency in data handling between the training and test sets.

(c)Decision Tree Classifier Initialization: A Decision Tree classifier is instantiated with specific parameters **criterion = "entropy"**. This criterion measures the quality of a split in the decision tree. Here, it utilizes information gain based on entropy. **random_state = 2**: Setting a seed for randomization ensures reproducibility of results. **max_depth=5**: Restricts the maximum depth of the decision tree to 5 levels, controlling overfitting.

(d)Training the Decision Tree Classifier: The Decision Tree classifier (DecisionTree) is trained using the imputed training dataset Xtrain_imputed along with corresponding target values Ytrain.

(e)Prediction and Evaluation: Predictions are made on the imputed test dataset Xtest_imputed using the trained Decision Tree classifier (DecisionTree.predict()). The accuracy_score function from scikit-learn computes the accuracy by comparing predicted values with the actual target values YTest. The calculated accuracy score is stored in acc and the model name ('Decision Tree') in the model list for potential comparative analysis.

(f)Output: The accuracy score of the Decision Tree classifier on the test set is displayed to provide an assessment of its predictive performance[9].

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.impute import SimpleImputer

# Assuming Xtrain has missing values and you've applied imputation
imputer = SimpleImputer(strategy='mean')
Xtrain_imputed = imputer.fit_transform(Xtrain)

# Transform Xtest using the same imputer
Xtest_imputed = imputer.transform(Xtest)

# Now, create and train the DecisionTreeClassifier
DecisionTree = DecisionTreeClassifier(criterion="entropy", random_state=2, max_depth=5)
DecisionTree.fit(Xtrain_imputed, Ytrain)

# Predict using the trained classifier and the transformed test set
predicted_values = DecisionTree.predict(Xtest_imputed)

✓ 0.0s

from sklearn.metrics import accuracy_score

# Predict using the trained classifier and the transformed test set
predicted_values = DecisionTree.predict(Xtest_imputed)

# Calculate accuracy
x = accuracy_score(Ytest, predicted_values)
acc.append(x)
model.append('Decision Tree')
print("DecisionTrees's Accuracy is: ", x*500)

✓ 0.0s

DecisionTrees's Accuracy is: 90.9090909090909
```

Figure 11. Code snippet of Implementation of Decision tree

The Frontend work : It comprises of creating a dashboard (Figure 13) using HTML CSS and Java script which displayed our data that has been given by our backend as well as we have been created our login page which is connected to our dashboard which gives the data of user signing in which will be displayed in database which has been created in mongo database.

WORKFLOW

- **Frontend User Authentication Flow (Client-Side)** Upon using an online application, users are required to enter their login credentials, which typically consist of a username and password. Users enter their information on a login screen presented by the frontend, which is constructed with HTML, CSS, and JavaScript. JavaScript controls how the form is submitted, does preliminary validation to make sure that all required fields are filled out, and can carry out simple tasks like confirming the length of the password or the structure of the email.

- **Backend (Side-Server)** JavaScript sends an authentication request to the backend API upon form submission. This request is received by the backend server, which is often constructed using a server-side language like Node.js. It verifies the credentials entered against a database Mongo db where user information is safely kept[9]. The backend indicates successful authentication to the frontend by generating a session token or sending a confirmation message if the credentials match those in the database.
- **Dashboard Access:** The frontend gives the user access to the dashboard Figure 13 as soon as the backend confirms that the authentication process was successful. Typically, this access is managed by confirming that the session token that was obtained during the authentication procedure is present and valid. After successful authentication, the user can access the HTML, CSS, and JavaScript dashboard interface[10].
- **MongoDB Data Retrieval Figure :** JavaScript code on the dashboard makes queries to backend APIs to obtain certain data needed for display. The backend server receives these requests and uses its interface with the MongoDB database to retrieve the required data.
- **Frontend technologies:** The dashboard's Figure 11 hierarchical framework is provided using HTML. It outlines the composition and organization of the webpage, showcasing components such as headers, forms, content sections, and navigation menus. Users can interact with the dashboard and input login credentials. Cascading Style Sheets, or CSS, fashion the HTML elements to improve the dashboard's visual appeal and user experience. It specifies the interface's general look, layout, fonts, and colors. The dashboard can be made aesthetically pleasing with CSS, guaranteeing readability and consistency on many screens and devices.
- **JavaScript:** JavaScript gives the dashboard more functionality and interactivity. It controls form submissions, handles dynamic content updates without necessitating a page reload, and carries out client-side validation to guarantee accurate user credential

input[10]. JavaScript enables smooth data retrieval and manipulation without interfering with the user's experience by facilitating communication with the backend through API queries.

- **APIs to access the backend system (MongoDB database):** Backend Server: The backend server serves as a bridge between the MongoDB database and the frontend and is constructed in languages such as Node.js , Express js, and others. It processes requests from the frontend and communicates with the database to get or change data.
- **MongoDB Database:** MongoDB Figure 13 is a NoSQL database that uses an adaptable JSON-like structure to store data. User passwords and other application data are safely stored by it[10]. Using APIs (Application Programming Interfaces), the backend server and MongoDB interact to carry out CRUD (Create, Read, Update, Delete) actions on the database in response to frontend queries.
- **APIs:** The actions that can be carried out on the database are defined by the APIs (HTTP endpoints) that are exposed by the backend server. By cross-referencing user credentials with the data that has been stored, these APIs manage authentication. In order to guarantee that only authorized users can access particular information, they also control data retrieval based on user access credentials.
- **Data Display and User Identification:** JavaScript controls form submission when a user enters their credentials in the frontend login screen Figure 12 .Through API queries, the backend server receives these credentials[10].The user's identity is validated by the backend by cross-referencing the credentials with the data kept in MongoDB.
- **Privileges of Access:** The backend recognises the user and ascertains their access privileges following successful authentication.The backend builds answers to API queries based on the user's permissions, granting access to just authorized information and features.

- Data Showcase:** The dashboard interface becomes accessible to the frontend upon getting the authenticated status from the backend. JavaScript initiates queries for particular data from backend APIs. After establishing a connection with MongoDB, the backend returns the required data to the frontend Figure 11. JavaScript makes sure that users only see data that they are authorized to view by dynamically updating the HTML components of the dashboard with the data that is received[10].

Figure 12. Login page for user authentication

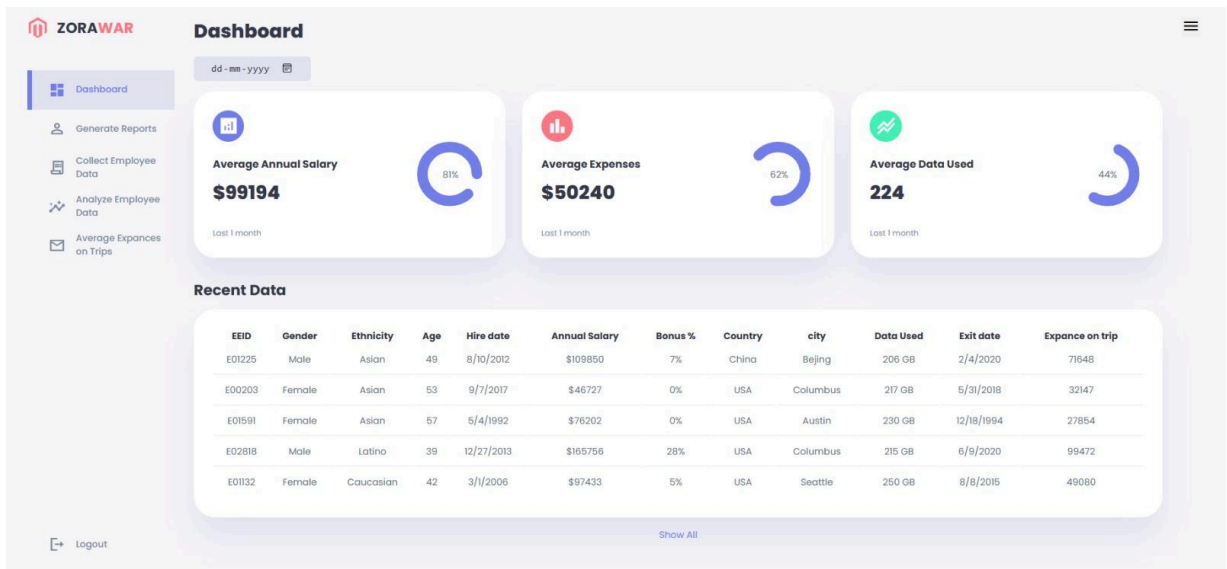


Figure 13. Dashboard of proposed model

Power BI

The generated Power BI dashboard (Figure.14) offers us a thorough summary regarding employee spending, travel information, internet usage, and employee locations in Bangalore, Delhi, and Pune. We organized the dashboard into five tiles, each of which uses a different representation to highlight a key insight. [15]

1. Employee Costs in Different Cities: The map visualization offers a detailed perspective of the spending hotspots in Bangalore, Delhi, and Pune in addition to displaying the distribution of employee expenses throughout these three cities. [15]Managers can more effectively distribute resources and discover high-spending areas thanks to their geographic awareness. For instance, if costs are excessively high in one city, it could be necessary to look into the underlying causes, which could include staff conduct, operational activities, or customer engagements. Furthermore, correlations between spending habits and outside variables can be investigated by superimposing demographic or corporate data into the map. This enables more well-informed decision-making and strategic planning.

2. Travel Charges and Staff Identity: The card visualization gives a brief summary of spending patterns by displaying the entire amount spent on staff travel expenses. Managers may immediately identify high spenders or outliers by seeing employee IDs next to trip charges. This makes it possible to implement targeted cost-control measures or policy changes. Additionally, adding employee IDs enables more in-depth analysis, including figuring out which employees travel a lot or connecting trip costs to particular tasks or divisions[15]. With the help of this knowledge, businesses can guarantee that spending regulations are followed, negotiate advantageous contracts with travel suppliers, and optimize their travel budgets.

3. Internet Data Usage Every Quarter: Each employee's quarterly internet data usage is displayed in a table representation that provides a thorough analysis of usage trends over time. Managers possess the ability to scrutinize patterns, detect deviations, and pinpoint prospects for enhancement or expense mitigation. For instance, it can be a sign of workflow inefficiencies or the need for more training on data management procedures if specific staff members often use more data than is permitted[15]. Managers can customize interventions or incentives to encourage safe usage while eliminating unnecessary costs by breaking down internet usage by employee IDs. Furthermore, past internet usage data can guide future infrastructure and capacity planning decisions, guaranteeing that resources are distributed efficiently to support corporate operations.

4. The ability to search employees: The ability to target data retrieval based on particular employee characteristics is one way that the slicer capability improves user experience. It is simple for users to look for specific employees by name, ID, or other pertinent criteria, which makes data exploration and analysis more efficient[15]. The dashboard(Figure 15.)

dynamically adjusts to show personalized details, such as trip expenses, internet usage, and allocated city, once an employee is selected. With the help of this interactive tool, users may examine employee records in greater detail, spot trends, and come to informed judgments[15]. Furthermore, managers may attend to individual needs, keep an eye on performance, and make sure that organizational regulations are being followed thanks to the ability to filter data based on personal traits.

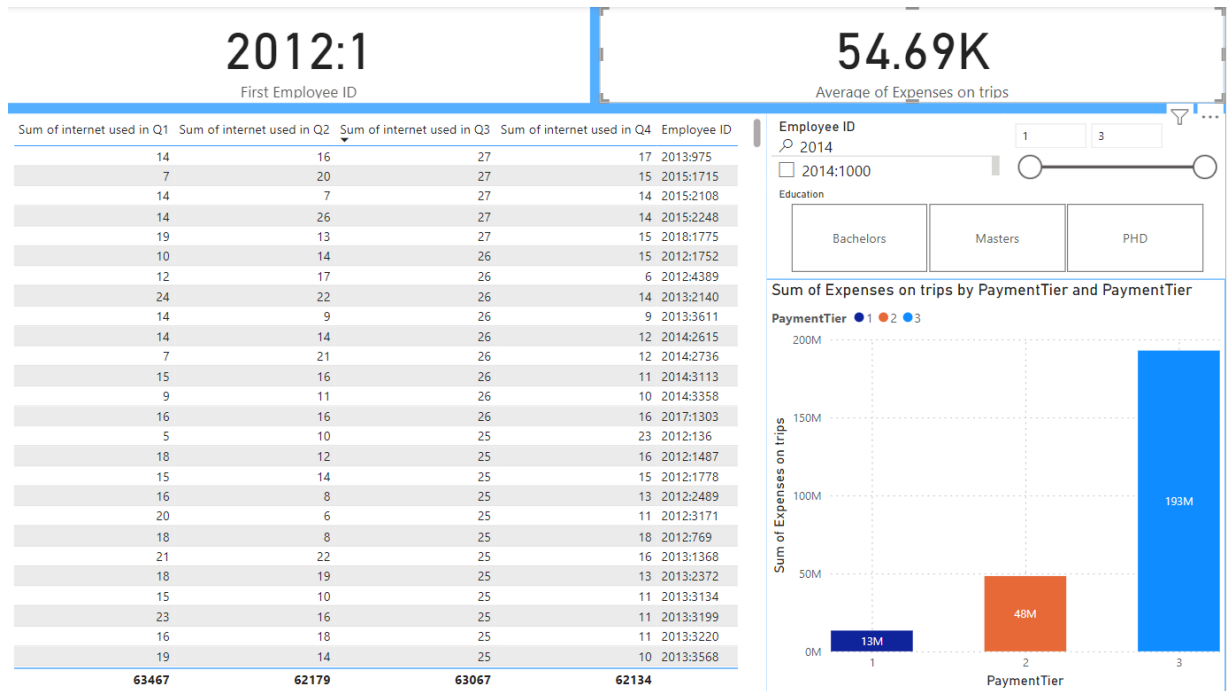


Figure 14. Power BI Dashboard Snapshot

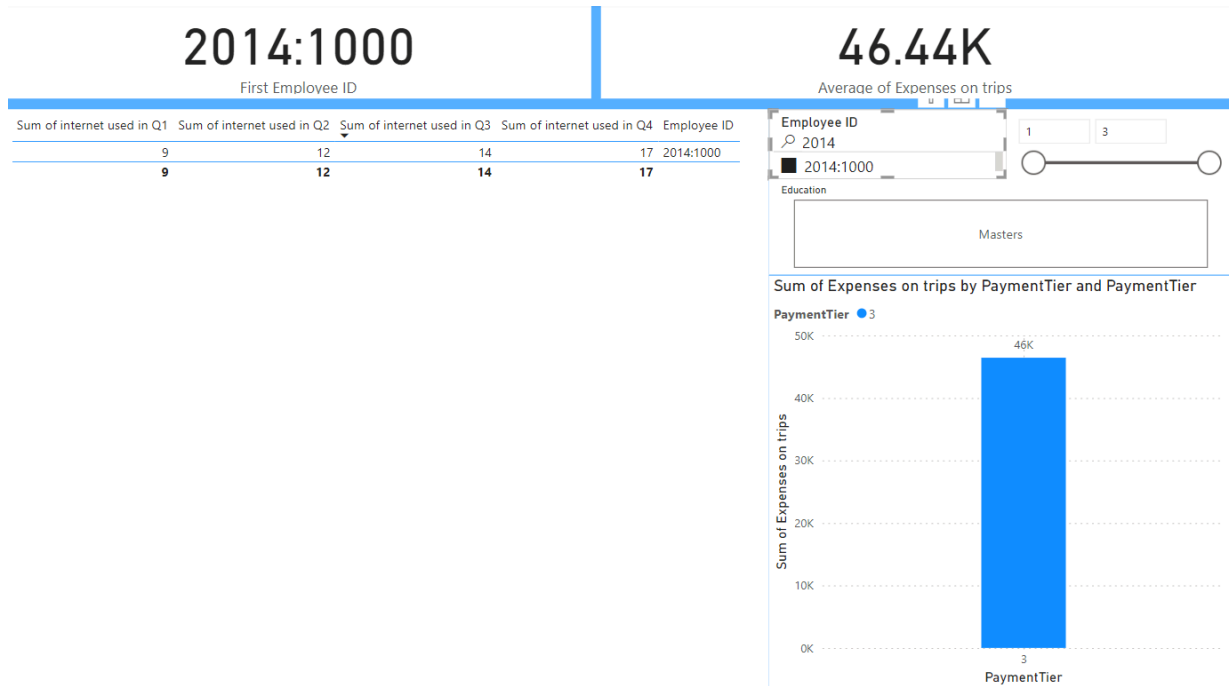


Figure 15. Exploring Employee Data with Slicer Search Bar: A Snapshot from Power BI Dashboard

CHAPTER - 04 PERFORMANCE ANALYSIS

The dataset used for testing was taken from the original dataset where it was divided between training and testing data. The evaluation metrics we use is Accuracy. The Random forest model and Decision Tree is working successfully as depicted in Figure 12 and Figure 13.



```
Random Forest

from sklearn.ensemble import RandomForestClassifier

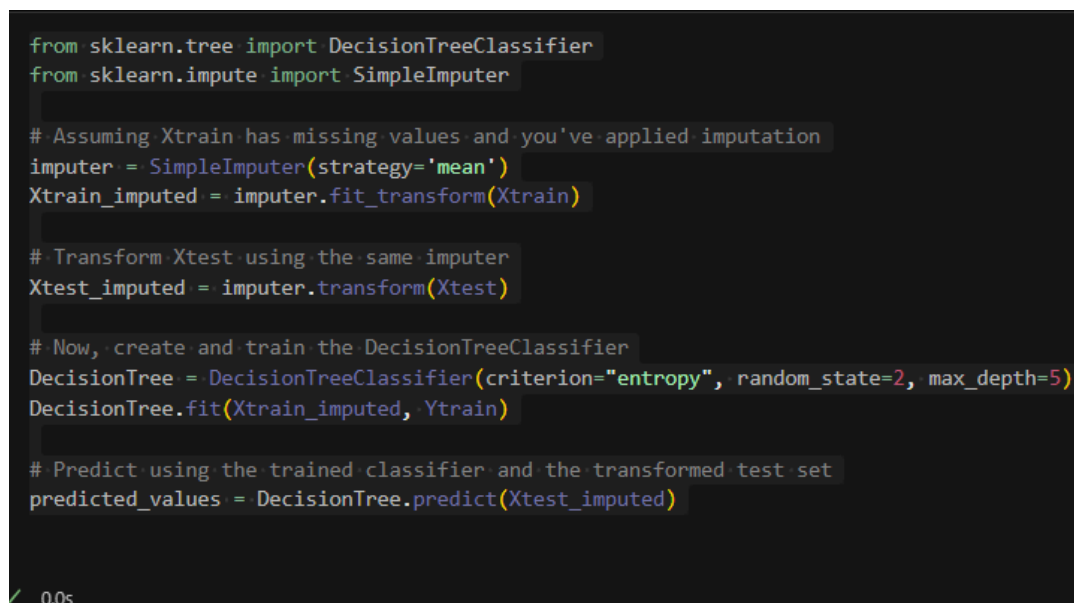
RF = RandomForestClassifier(n_estimators=20, random_state=0)
RF.fit(Xtrain,Ytrain)

predicted_values = RF.predict(Xtest)

x = metrics.accuracy_score(Ytest, predicted_values)
acc.append(x)
model.append('RF')
print("RF's Accuracy is: ", x)

print(classification_report(Ytest,predicted_values))
```

Figure 16. Implementation and training of RF



```
from sklearn.tree import DecisionTreeClassifier
from sklearn.impute import SimpleImputer

# Assuming Xtrain has missing values and you've applied imputation
imputer = SimpleImputer(strategy='mean')
Xtrain_imputed = imputer.fit_transform(Xtrain)

# Transform Xtest using the same imputer
Xtest_imputed = imputer.transform(Xtest)

# Now, create and train the DecisionTreeClassifier
DecisionTree = DecisionTreeClassifier(criterion="entropy", random_state=2, max_depth=5)
DecisionTree.fit(Xtrain_imputed, Ytrain)

# Predict using the trained classifier and the transformed test set
predicted_values = DecisionTree.predict(Xtest_imputed)
```

Figure 17. Code Snippet of training Decision Tree

```

SVM's Accuracy is: 74.86573576799141
      precision    recall  f1-score   support

     0       0.74       0.94       0.83       597
     1       0.79       0.41       0.54       334

 accuracy                   0.75       931
 macro avg       0.76       0.67       0.68       931
 weighted avg    0.76       0.75       0.72       931

```

Figure 18. Code snippet of SVM's accuracy

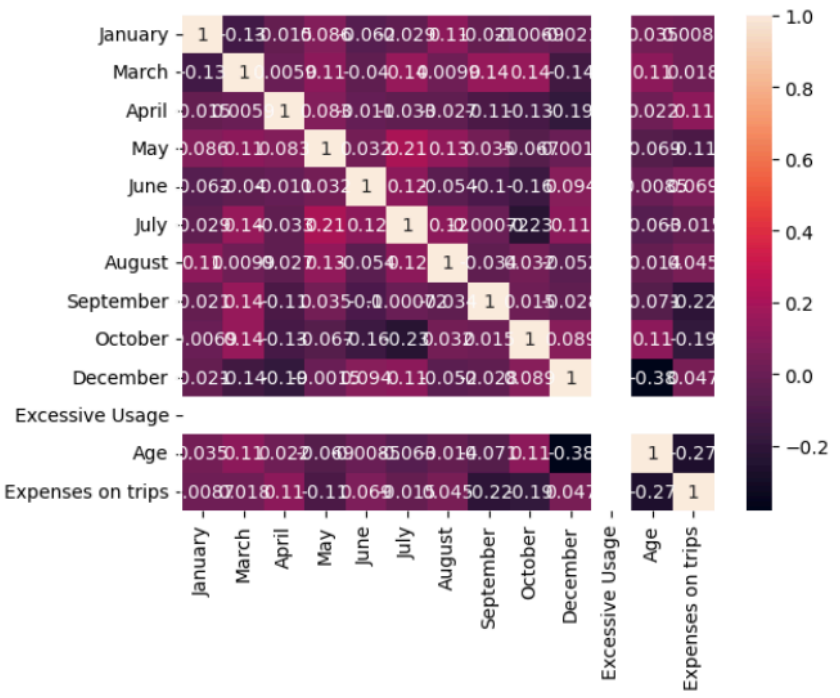


Figure 19. Heatmap of dataset

Reverse Connectivity: The data on the dashboard is provided by backend services. A MongoDB database or similar backend system may hold this data. Application Programming Interfaces, or backend APIs, are made to handle frontend (dashboard) requests. These APIs send data back to the frontend for display after retrieving it from the database.

Integration of Login Pages: Before being able to access the dashboard, users can authenticate themselves on the login page. The login page's user credentials, such as the password and

username, are compared to those kept in the backend (MongoDB database). Users get access to the dashboard after their authentication is accepted.

Database Interaction and Data Display: The purpose of the dashboard is to show data that has been pulled from the backend. JavaScript is used to send requests to the backend APIs (e.g., XMLHttpRequest or Fetch API) in order to retrieve the necessary data from the MongoDB database. The HTML structure and CSS styling of the retrieved data are then used to dynamically render and display it in the dashboard's user interface. JavaScript functions may be triggered by user activities (e.g., clicking buttons or selecting options) in the dashboard to retrieve updated or additional data from the backend.

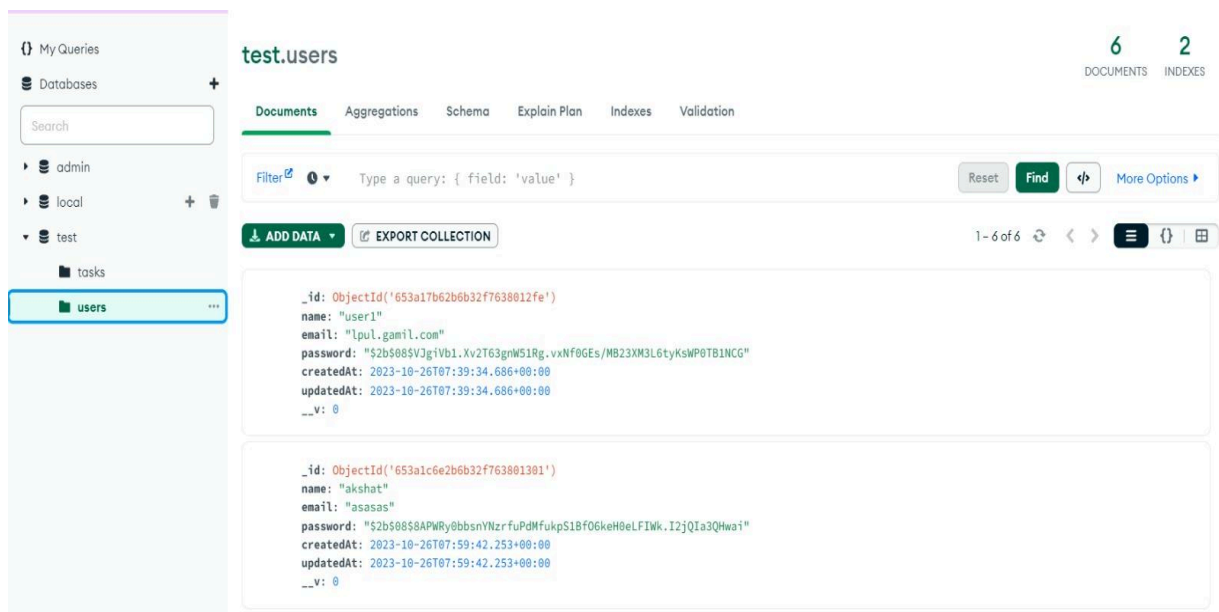


Figure 20. Snippet of Mongoddb Terminal

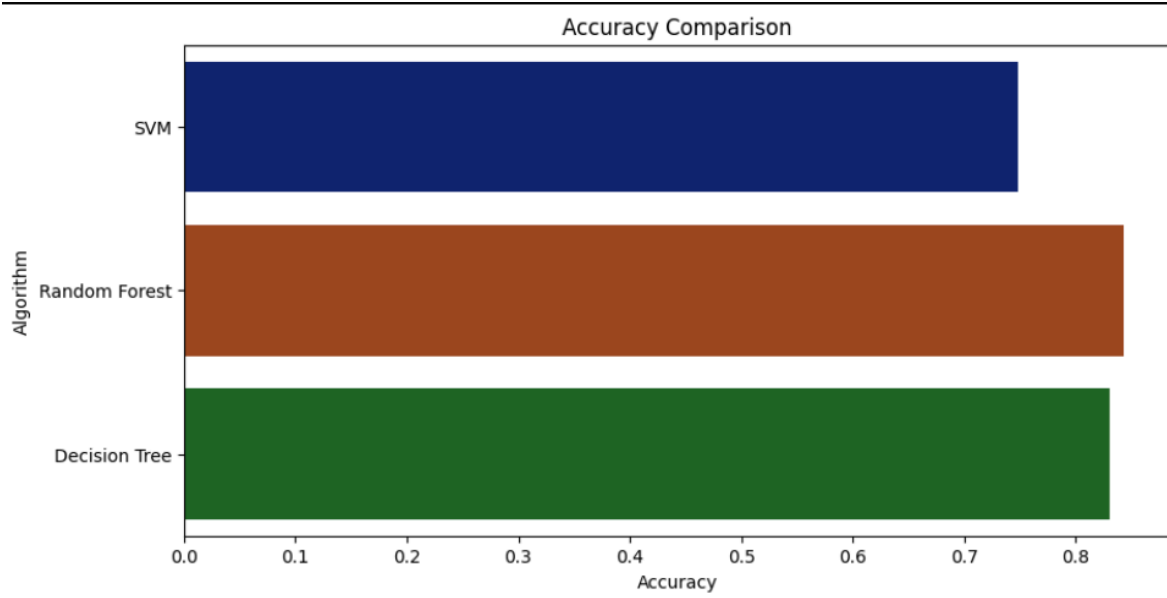


Figure 21. Snapshot of Machine Learning Model Accuracy Evaluation

CHAPTER - 05 CONCLUSION

The development and implementation of an employee internet connection and expense monitoring system represent a crucial stride toward optimizing resource utilization, enhancing security measures, and fostering a culture of compliance within an organization. By systematically tracking internet usage and managing expenses, this system offers multifaceted benefits across various operational domains. It empowers companies to exercise better cost management, gain insights into employee productivity, and fortify security measures against potential risks. Moreover, such a system bolsters transparency and accountability, promoting fair expense policies and adherence to established guidelines. It enables data-driven decision-making by providing actionable insights derived from comprehensive analytics and reporting. However, successful implementation necessitates careful planning, robust technological infrastructure, and a clear communication strategy. Striking a balance between monitoring and respecting employee privacy while ensuring compliance with regulatory standards is pivotal to fostering a supportive work environment. Continuous evaluation, feedback integration, and iterative improvements are vital components of the system's lifecycle, ensuring it remains aligned with evolving organizational needs and technological advancements.

Ultimately, an employee internet connection and expense monitoring system serve as a catalyst for organizational efficiency, security, and growth while prioritizing the well-being and productivity of the workforce. Its successful integration into the operational framework can significantly contribute to achieving strategic objectives and maintaining a competitive edge in the ever-evolving business landscape.

5.2 Application Contribution

Efficiency of Operations: Contribution: Reduces inefficiencies, streamlines procedures, and maximizes the use of available resources.

Prudence in finances: Contribution: Maintains responsible financial management, minimizes wasteful spending, and controls costs. Employee Contentment and Engagement Contribution: Increases employee satisfaction by fostering trust through equitable expense management. Risk Reduction: Contribution: Protects the interests of the organization by identifying and mitigating potential hazards. Growth and Strategic Planning:Contribution: Promotes forecasting, data-driven initiatives, and organizational expansion. Observance and Moral Conduct Contribution: Promotes a culture of integrity and accountability by making sure rules are followed. Constant Enhancement Contribution: Offers feedback loops to support continuous improvements, guaranteeing flexibility and development.

REFERENCES

- [1] Huang Ruwei, Gui Lin, Yu Si, Zhuang Wei. Cloud environments in support of the privacy protection can be calculated using encryption method [J]. Journal of the computer. 2011 (12)
- [2] Mao Jian, Li Kun, Xu Xiandong. Privacy protection scheme in cloud computing environment [J]. Journal of Tsinghua University (NATURAL SCIENCE EDITION). 2011 (10).
- [3] Lv Zhiquan, Aman Chang, Feng Dengguo. Cloud storage access control scheme [J]. computer science and exploration.2011 (09)
- [4] K. Stanoevska-Slabeva, T. Wozniak Grid and Cloud Computing-A Business Perspective on Technology and Applications Springer-Verlag, Berlin,Heidelberg (2010)
- [5] G. Reese Cloud Application Architectures: Building Applications and Infrastructure in the Cloud, Theory in Practice, O'Reilly Media (2009)
- [6] S. Shah, "Convolutional Neural Network: An Overview," Analytics Vidhya, Jan. 27, 2022. <https://www.analyticsvidhya.com/blog/2022/01/convolutional-neural-network-an-overview/>
- [7] D. Lakkas Establishing and managing trust within the public key infrastructure Computer Communications, 26 (16) (2003)
- [8] R. Sherman Distributed systems security Computers & Security, 11 (1)
- [9] J. Lee The era of Omni-learning: Frameworks and practices of the expanded human resource development Organizational Dynamics (2023)(1992)
- [10] M.Y.Abdel Sadek et al.Matching-based resource allocation for critical MTC in massive MIMO LTE networks IEEE Access (2019)
- [11] DoD Computer Security Center ,Trusted computer system evaluation criteria,DoD 5200.28-STD, 1985.
- [12]<https://www.analyticsvidhya.com/blog/2021/06/support-vector-machine-better-understanding/>
- [13] P. S. [Premanand S.], "Support Vector Machine: Beginners Guide - Analytics Vidhya [SNIPPET]", Analytics Vidhya, Jun. 2021, [Online].
- [14] A. [alokesh985], "Introduction to Support Vector Machines (SVM)", GeeksforGeeks, [Online].
- [15] L. Liu and J.-F. Chamberland, "On the effective capacities of multiple antenna Gaussian channels," in Proc. IEEE Int. Symp. Inf. Theory, Jul. 2008, pp. 2583–2587.

- [16] J. Tang and X. Zhang, "Cross-layer-model based adaptive resource allocation for statistical QoS guarantees in mobile wireless networks," IEEE
- [17] *Trans. Wireless Commun.*, vol. 7, no. 6, pp. 2318–2328, Jun. 2008.
- [18] F. Kelly, S. Zachary, and I. Ziedins, *Stochastic Networks: Theory and Applications*. London, U.K.: Oxford Univ. Press, 1996.
- [19] J. Tang and X. Zhang, "Quality-of-service driven power and rate adaptation over wireless links," *IEEE Trans. Wireless Commun.*, vol. 6, no. 8, pp. 3058–3068, Aug. 2007.
- [20] C. H. Papadimitriou and K. Steiglitz, *Combinatorial Optimization: Algorithms and Complexity*. Chelmsford, MA, USA: Courier, 1998.
- [21] E. Bodine-Baron, C. Lee, A. Chong, B. Hassibi, and A. Wierman, "Peer effects and stability in matching markets," in *Proc. Int. Symp. Algorithmic Game Theory*. Berlin, Germany: Springer, 2011, pp. 117–129.
- [22] L. Liu and J.-F. Chamberland, "On the effective capacities of multiple antenna Gaussian channels," in *Proc. IEEE Int. Symp. Inf. Theory*, Jul. 2008, pp. 2583–2587.
- [23] J. Tang and X. Zhang, "Cross-layer-model based adaptive resource allocation for statistical QoS guarantees in mobile wireless networks," IEEE
- [24] *Trans. Wireless Commun.*, vol. 7, no. 6, pp. 2318–2328, Jun. 2008.
- [25] F. Kelly, S. Zachary, and I. Ziedins, *Stochastic Networks: Theory and Applications*. London, U.K.: Oxford Univ. Press, 1996.
- [26] J. Tang and X. Zhang, "Quality-of-service driven power and rate adaptation over wireless links," *IEEE Trans. Wireless Commun.*, vol. 6, no. 8, pp. 3058–3068, Aug. 2007.
- [27] C. H. Papadimitriou and K. Steiglitz, *Combinatorial Optimization: Algorithms and Complexity*. Chelmsford, MA, USA: Courier, 1998.
- [28] E. Bodine-Baron, C. Lee, A. Chong, B. Hassibi, and A. Wierman, "Peer effects and stability in matching markets," in *Proc. Int. Symp. Algorithmic Game Theory*. Berlin, Germany: Springer, 2011, pp. 117–129.

APPENDICES

Random forest:

Among the supervised learning methods is the well-known machine learning algorithm Random Forest. It can be applied to ML issues involving both classification and regression. Its foundation is the idea of ensemble learning, which is the process of merging several classifiers to solve a challenging issue and enhance the model's functionality.

According to its name, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Rather than depending on a single decision tree, the random forest forecasts the outcome based on the majority vote of projections from each tree.

CNN:

Tens or even hundreds of layers can be found in a convolutional neural network, each of which is trained to recognise a unique characteristic of an image. Every training image is subjected to various resolutions of filters, and the result of every convolved image serves as the input for the subsequent layer. The filters can begin with relatively basic criteria, such edges and brightness, and progress in sophistication to include features that specifically identify the object.

SK learn:

A machine learning package for the Python programming language, scikit-learn (formerly known as scikits.learn and also named sk learn) is available as free software.[3] With support-vector machines, random forests, gradient boosting, k-means, DBSCAN, and other classification, regression, and clustering techniques, it is compatible with the NumPy and SciPy scientific and numerical libraries for Python. Scikit-learn is a financially supported project by NumFOCUS.

Project report May 6

ORIGINALITY REPORT

19%

SIMILARITY INDEX

17%

INTERNET SOURCES

10%

PUBLICATIONS

9%

STUDENT PAPERS

PRIMARY SOURCES

1	ir.juit.ac.in:8080 Internet Source	3%
2	Mohammed Y. Abdelsadek, Yasser Gadallah, Mohamed H. Ahmed. "Matching-Based Resource Allocation for Critical MTC in Massive MIMO LTE Networks", IEEE Access, 2019 Publication	3%
3	www.ir.juit.ac.in:8080 Internet Source	2%
4	www.irejournals.com Internet Source	1%
5	cyberleninka.org Internet Source	1%
6	Sai Rama Krishna Indarapu, Swathy Vodithala, Naveen Kumar, Siripuri Kiran, Soora Narasimha Reddy, Kumar Dorthi. "Exploring human resource management intelligence practices using machine learning models",	1%

JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT
PLAGIARISM VERIFICATION REPORT

Date: 10 May 2024
 Type of Document (Tick): **PhD Thesis** **M.Tech/M.Sc. Dissertation** **B.Tech./B.Sc./BBA/Other**
 Name: SURBHIT, AKSHAT SHARMA Department: CSE-IT Enrolment No. 201521, 201247
 Contact No. 9794322570 E-mail: Surbhisharma16@gmail.com
 Name of the Supervisor: Dr Pradip K. Gupta
 Title of the Thesis/Dissertation/Project Report/Paper (In Capital letters): DEVELOPING CLOUD BASED SECURE DATA MONITORING SYSTEM

UNDERTAKING

I undertake that I am aware of the plagiarism related norms/ regulations, if I found guilty of any plagiarism and copyright violations in the above thesis/report even after award of degree, the University reserves the rights to withdraw/ revoke my degree/report. Kindly allow me to avail Plagiarism verification report for the document mentioned above.

- Total No. of Pages =
- Total No. of Preliminary pages =
- Total No. of pages accommodate bibliography/references =

[Signature]
 (Signature of Student)

FOR DEPARTMENT USE

We have checked the thesis/report as per norms and found **Similarity Index** at 19% (%). Therefore, we are forwarding the complete thesis/report for final plagiarism check. The plagiarism verification report may be handed over to the candidate.

[Signature]
 (Signature of Guide/Supervisor)

[Signature]
 Signature of HOD

FOR LRC USE

The above document was scanned for plagiarism check. The outcome of the same is reported below:

Copy Received on	Excluded	Similarity Index (%)	Abstract & Chapters Details	
	<ul style="list-style-type: none"> • All Preliminary Pages • Bibliography/Images/Quotes • 14 Words String 		Word Counts	
Report Generated on			Character Counts	
		Submission ID	Page counts	
			File Size	

Checked by
 Name & Signature

Librarian

Please send your complete Thesis/Report in (PDF) & DOC (Word File) through your Supervisor/Guide at plagcheck.juit@gmail.com

