

Resilient Rivers- Empowering Flood-Free Futures

Submitted in partial fulfillment of the requirements of the
degree

BACHELOR OF ENGINEERING IN COMPUTER ENGINEERING

By

1. Aryan Girish Raje (D12A-51)
2. Arya Girish Raje (D12A-50)
3. Ishita Sudhir Marathe (D12A-41)
4. Prasad Kishor Lahane (D12A-36)

Name of the Mentor

Prof. Dr.Mrs.Gresha Bhatia



Vivekanand Education Society's Institute of Technology,

An Autonomous Institute affiliated to University of Mumbai

HAMC, Collector's Colony, Chembur,

Mumbai-400074

University of Mumbai (AY 2023-24)

CERTIFICATE

This is to certify that the Mini Project entitled “Resilient Rivers- Empowering Flood-Free Futures” is a bonafide work of Aryan Raje (D12A-51), Arya Raje(D12A-50), Ishita Marathe(D12A-44), Prasad Lahane(D12A-36) submitted to the University of Mumbai in partial fulfillment of the requirement for the award of the degree of “**Bachelor of Engineering**” in “**Computer Engineering**” .

(Prof. _____)

Mentor

(Prof. _____)

Head of Department

(Prof. _____)

Principal

Mini Project Approval

This Mini Project entitled “**Resilient Rivers- Empowering Flood-Free Futures**” by **Aryan Raje (D12A-51), Arya Raje(D12A-50), Ishita Marathe(D12A-44), Prasad Lahane(D12A-36)** is approved for the degree of **Bachelor of Engineering in Computer Engineering**.

Examiners

1.....
(Internal Examiner Name & Sign)

2.....
(External Examiner name & Sign)

Date: 13.04.2024

Place: Mumbai, India

Declaration

We declare that this written submission represents our ideas in our own words and where others' ideas or words have been included, we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

(Signature)

(Aryan Raje (D12A-51))

(Signature)

(Arya Raje(D12A-50))

(Signature)

(Ishita Marathe(D12A-44))

(Signature)

(Prasad Lahane(D12A-36))

Date: 13/04/2024

Contents

Abstract	ii
Acknowledgments	iii
List of Figures	iv
1 Introduction	5
1.1 Introduction	
1.2 Motivation	
1.3 Problem Statement & Objectives	
1.4 Organization of the Report	
2 Literature Survey	9
2.1 Survey of Existing System	
2.2 Limitation Existing system or Research gap	
2.3 Mini Project Contribution	
3 Proposed System	14
3.1 Introduction	
3.2 Architectural Framework / Conceptual Design	
3.3 Algorithm and Process Design	
3.4 Methodology Applied	
3.5 Hardware & Software Specifications	
3.6 Experiment and Results for Validation and Verification	
3.7 Result Analysis and Discussion	
3.8 Conclusion and Future work.	
References	32
Annexure	33

Abstract

This report conducts a comparative study on various models predicting overflow probability in different schemes within the Narmada Basin, a crucial water resource in central India. By employing statistical, hydrological, and machine learning methods and considering factors like rainfall patterns and basin topography, the study assesses model performance. Through literature review and empirical data validation, it aims to identify strengths and weaknesses in capturing overflow dynamics of water released from the various schemes of Narmada river Basin. The insights gained contribute to enhancing flood forecasting and management strategies, supporting sustainable water resource utilization in the Narmada region.

Acknowledgements

Perseverance, Inspiration & Motivation have always played a key role in the success of any venture. At this level of understanding it is difficult to understand the wide spectrum of knowledge without proper guidance and advice, Hence we take this time to express our sincere gratitude to our respected Project Guide Dr. Mrs. Gresha Bhatia who as a guide evolved an interest in us to work and select an entirely new idea for project work. She has been keenly cooperative and helpful to us in sorting out all the difficulties. We would also like to thank our Principal Mrs. J. M. Nair, for this golden opportunity. My deep sense of gratitude to Vivekanand Education Society's Institute of Technology for their timely advice and encouragement in our project development. I would also thank my Institution and my faculty members without whom this project would have been a distant reality.

List of Figures

Sr.No	Name	Page No
1	Figure 1.1.Flood hazard factors in Narmada River Basin	5
2	Figure 3.1.Architectural Framework	15
3	Figure.3.2.Methodology	20
4	Figure 3.3 Models Used	22
5	Figure 3.4.1 SVM Working	31
6	Figure 3.4.2 Classification Report of SVM	31
7	Figure 3.4.3 Classification Report of Logistic Regression	32
8	Figure 3.4.4 Classification Report of Random Forest Regression	33
9	Figure 3.4.5 Classification Report of Decision Tree	31
10	Figure 3.4.6 Classification Report of Naive Bayes	31
11	Figure 3.4.7 Working of LSTM	33
12	Figure 3.4.8 Classification Report of LSTM	34
13	Figure 3.4.9 Working of Recurrent Neural Network	35
14	Figure 3.4.10 Classification report of Recurrent Neural Network	39
15	Figure 3.5 Comparison of different machine learning algorithms	42
16	Figure 3.7.1 Nodal Officer Login	43
17	Figure 3.7.2 Predictor page	44
18	Figure 3.7.3 SMS sent via API	45

1. Introduction

1.1 Introduction

The release of approximately 18 lakh cusecs of water from the Sardar Sarovar dam on September 17, 2023, resulted in floods in Narmada, Bharuch, and parts of Vadodara district for two days. SSNNL reported an inflow of nearly 22 lakh cusecs, managed through diverting some water into canals for filling water bodies in other areas. The incident, attributed to state government mismanagement, highlights the need for a predictive monitoring system to prevent sudden water overflow. Our aim is to develop a system which can avoid such mishaps from happening again

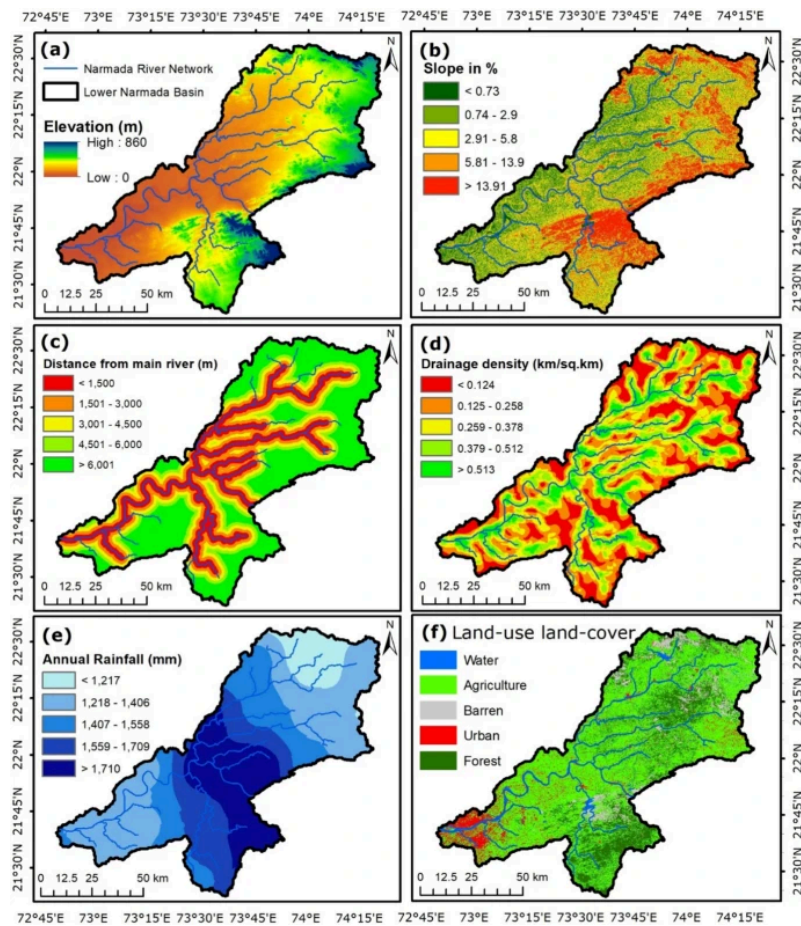


Figure 1.1. Flood hazard factors in Narmada River Basin

Numerous mathematical and statistical models have been developed for prediction of flood and various comparison models and surveys have been made using data-driven models like Machine Learning over them. Various machine learning algorithms are available, each tailored to different output requirements and datasets. For flood prediction, supervised learning models are deemed most effective [3].

1.2 Motivation

Motivated by the pressing need to enhance the accuracy and timeliness of flood forecasting, and driven by the imperative to expedite relief efforts during emergencies, this study endeavors to revolutionize the field of disaster preparedness and response. The aim is to develop a sophisticated flood prediction model that integrates cutting-edge machine learning techniques, historical meteorological data, and topographical insights to enable more precise and timely flood predictions. Additionally, the development of a user-friendly disaster relief dashboard serves to address the challenges associated with efficient coordination and resource allocation during critical flood situations.

By addressing these crucial gaps in current flood prediction and disaster management systems, this research aspires to empower authorities, relief organizations, and communities with the necessary tools and information to proactively plan and respond to flood disasters. The ultimate goal is to minimize the impact of floods on vulnerable communities, safeguard lives and livelihoods, and foster a more resilient and adaptive approach to mitigating the consequences of such natural calamities.

1.3 Problem Statement and Objectives

1.3.1 Problem Statement:

1. The release of approximately 18 lakh cusecs of water from the Sardar Sarovar dam on September 22, 2023, resulted in floods in Narmada, Bharuch, and parts of Vadodara district for two days.
2. The incident, attributed to state government mismanagement, highlights the need for a predictive monitoring system to prevent sudden water overflow.
3. Our aim is to develop a system which will monitor the level of the reservoir and predict efficiently the overflow safe conditions of the narmada basin

1.3.2 Objectives:

1. Develop an advanced flood prediction model that integrates historical meteorological data, topographical analysis, and machine learning techniques to improve the accuracy and lead time of flood forecasts.
2. Incorporate various data sources, such as rainfall patterns, river discharge levels, and terrain characteristics, into the predictive model to conduct a comprehensive assessment of flood risks.
3. Create and implement a user-friendly disaster relief dashboard that consolidates data on flood-prone regions, essential infrastructure, and potential evacuation areas to facilitate prompt decision-making during flood emergencies.
4. Enable smooth collaboration among disaster response teams, government entities, and humanitarian organizations through the disaster relief dashboard, ensuring the efficient allocation of resources for timely relief and recovery efforts.
5. Empower authorities and communities with the necessary resources and information to execute proactive measures and establish robust disaster management strategies, minimizing the adverse impact of floods on lives and livelihoods.

1.4 Organization of the Report

Reports are written to present and discuss research findings. They provide the reader with the rationale for the research, a description of the method used to conduct the research, the findings, results, a logical discussion, and conclusions/recommendations.

The following report consists of Introduction to my project, abstract and acknowledgments to all the people who have helped me. Moreover it contains a detailed literature review, surveying the existing research and finding additions to our project based on that knowledge. Finally, we have attached the results and outputs of the work done so far.

It majorly consists of the machine learning models that we have selected. The preprocessing of the data that has led to clean data and application of the models. In the report we have attached the precision and accuracy of the models chosen. Along with how these models work. The implementation of the model through deployment of a website is also attached.

2. Literature Survey

2.1 Survey of Existing System

2.1.1.Introduction

Flood prediction is a critical area of research aimed at mitigating the devastating impacts of flooding events. Various machine learning and statistical algorithms have been explored to develop accurate flood prediction models. In this literature survey, we review six research papers that investigate different methodologies for flood prediction using machine learning and statistical approaches. Each paper evaluates different methodologies, including Support Vector Machines (SVM), Naive Bayes, Decision Tree, Logistic Regression, Random Forest, Recurrent Neural Networks (RNN), and numerical simulations. The papers highlight the importance of accurate flood prediction for effective disaster management and infrastructure planning. Despite significant advancements, challenges such as data availability, model complexity, and generalization to diverse geographic regions persist in flood prediction research.

2.2.1.Paper 1 : Performance Evaluation of Different Machine Learning Based Algorithms for Flood Prediction and Model for Real Time Flood Prediction

Abstract:The authors evaluate various machine learning algorithms for flood prediction and find that Recurrent Neural Networks (RNN) perform exceptionally well due to their error correction capabilities. Their study underscores the importance of selecting appropriate algorithms for accurate flood prediction, with RNN demonstrating superior performance in terms of accuracy and precision.

Inference:RNN demonstrates superior accuracy and precision in flood prediction compared to other algorithms tested. Its ability to incorporate error correction through backpropagation contributes to its exceptional performance in predicting floods.

2.2.2.Paper 2: Urban Flash Flood Forecast Using Support Vector Machine and Numerical Simulation

Abstract: This paper proposes a flood forecast model combining Support Vector Machine (SVM) with numerical simulation to predict urban flash floods. By integrating SVM with numerical simulation, the model achieves accurate urban flash flood forecasting, facilitating timely disaster response and mitigation efforts.

Inference: SVM integrated with numerical simulation enables accurate urban flash flood forecasting, but faces challenges related to data quality and computational resources. Despite these challenges, the model holds promise for enhancing urban flood prediction capabilities.

2.2.3. Paper 3: A Review on Flood Prediction Algorithms and A Deep Neural Network Model for Estimation of Flood Occurrence

Abstract: The paper presents a review of flood prediction algorithms and outlines a methodology for flood prediction using Deep Neural Networks (DNN). The study highlights the potential of DNNs in improving flood prediction accuracy and discusses the challenges associated with their implementation.

Inference: Deep Neural Networks (DNN) offer improved flood prediction accuracy but require substantial computational resources and face challenges in generalization and model interpretability. Despite these challenges, DNNs represent a promising approach for enhancing flood prediction capabilities.

2.3.4. Paper 4: Scenario-Based Real-Time Flood Prediction with Logistic Regression

Abstract: Logistic regression is utilized for real-time flood prediction by establishing relationships between independent variables and flood occurrences. The study demonstrates the effectiveness of logistic regression in providing accurate binary predictions for flood occurrences in real-time scenarios.

Inference: Logistic regression provides accurate binary predictions for flood occurrences but is sensitive to data assumptions and may not capture complex relationships effectively. Despite these limitations, logistic regression remains a valuable tool for real-time flood prediction applications.

2.3.5. Paper 5: Water Level Prediction Model Applying a Long Short-Term Memory (LSTM)–Gated Recurrent Unit (GRU) Method for Flood Prediction

Abstract: The paper proposes a water level prediction model using Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) architectures. The study highlights the potential of LSTM and GRU models in improving flood prediction accuracy and discusses the challenges associated with their implementation.

Inference:LSTM and GRU models offer effective flood prediction capabilities but require extensive hyperparameter tuning and computational resources. Despite these challenges, LSTM and GRU architectures represent promising approaches for enhancing flood prediction accuracy.

2.2.6.Paper 6: Using Multi-Factor Analysis to Predict Urban Flood Depth Based on Naive Bayes

Abstract:Naive Bayes is employed for predicting urban flood depth based on multi-factor analysis. The study demonstrates the effectiveness of Naive Bayes in providing accurate predictions for urban flood depth, highlighting its potential applications in flood risk assessment and management.

Inference:Naive Bayes provides accurate predictions for urban flood depth but assumes feature independence and may require adjustments for imbalanced datasets. Despite these limitations, Naive Bayes remains a valuable tool for predicting urban flood depth and assessing flood risk in urban areas.

2.3 Limitation Existing system or Research gap

1. **Data Quality and Completeness:** Historical data quality and accuracy are potential concerns, impacting the reliability of flood predictions (Reference Paper 1).
2. **Generalizability:** Models developed for specific regions might lack generalizability to diverse or changing environments, raising questions about their applicability beyond the studied area (Reference Paper 1).
3. **Computational Demands:** Developing urban flood prediction models demands significant computational resources and expertise, making them inaccessible for regions with limited capabilities (Reference Paper 2).
4. **Complexity and Interpretability:** Complex models lack interpretability, posing challenges for stakeholders in understanding the underlying processes, particularly in urban flood contexts (Reference Paper 2).
5. **Ethical Concerns:** Issues related to data privacy, biases, and model reliability raise ethical concerns in the application of flood prediction systems, emphasizing the need for ethical considerations (Reference Paper 2).
6. **Real-time Deployment:** Real-time deployment of flood prediction models demands continuous updates and substantial resources, which might not be feasible for many regions (Reference Paper 2).
7. **Handling Non-linearity:** Many existing flood prediction models struggle with capturing non-linear relationships present in urban environments, affecting prediction accuracy (Reference Paper 3).
8. **Spatial-Temporal Resolution:** Data resolution limitations might miss localized variations, impacting the ability to capture specific patterns in flood occurrences, indicating challenges in capturing detailed spatial-temporal information (Reference Paper 3).
9. **Resource Intensity:** The resource-intensive nature of flood prediction model development and deployment poses challenges, especially for regions with limited resources (Reference Paper 3).
10. **Interdisciplinary Understanding:** Integrating knowledge from various disciplines like hydrology, climatology, and data science is crucial for effective flood prediction but poses challenges in interdisciplinary understanding and collaboration (Reference Paper 4).
11. **Limited Causality Understanding:** Existing flood prediction models may lack a deep understanding of causal relationships, limiting their predictive power in complex scenarios (Reference Paper 5).

2.3 Mini Project Contribution

1. **Improved Data Quality and Quantity:** Our system understands that having a diverse range of parameters enriches our dataset. This diversity can lead to more accurate predictions and a deeper understanding of the factors influencing dam overflow and flooding.
2. **Enhanced Feature Selection:** We're experimenting with different feature combinations and machine learning algorithms to identify the most relevant features for flood prediction. By doing so, we can significantly enhance the accuracy of our predictive models.
3. **Spatial and Temporal Considerations:** The inclusion of parameters like scheme ID, scheme name, and gate positions allows our system to account for spatial variability. Additionally, parameters such as cumulative rainfall provide temporal data, helping us consider changing weather patterns over time.
4. **Real-time Data Integration:** Parameters like present water level, inflow, outflow, and cumulative rainfall offer real-time insights. Integrating this data into our system ensures that our predictions are always based on the most current information, enabling timely and accurate flood forecasts.
5. **User Interface and Accessibility:** We're focusing on designing a user-friendly interface that visualizes real-time data and flood predictions. A clear and accessible interface ensures that dam operators and relevant authorities can quickly comprehend the information, facilitating faster decision-making during flood events.
6. **Continuous Monitoring and Improvement:** Our system is committed to continuously monitoring its performance.

3. Proposed System

3.1 Introduction

Machine Learning Algorithms utilize an automatic inductive approach to “learn” from data and visualize trends which can be used to develop a prediction model to pre-process the new datasets to develop predictions depending on the accuracy of the prediction model. The employment of Machine Learning is crucial as it allows for vast amounts of data to be processed which can be fed onto the Machine Learning algorithms and then trained using supervised or unsupervised learning. This is done to classify the data which is collected using various sources.

In the proposed set-up, the dataset is collected from https://wrd.guj.nic.in/dam/hour_reports_h.php and the various input data like Design gross Storage, Rule Level, Present Gross Storage, Outflow River, Cum.Rainfall, gate-Position-Nos, Scheme, FRL, Present-Water Level(m), Inflow, Outflow Canal, Type of Gate, Opening, Using this data, a relationship between the level of the reservoir and the gate opening of the reservoir can be established which can be used to train the random Forest Regression Model.

3.2 Architectural Framework / Conceptual Design

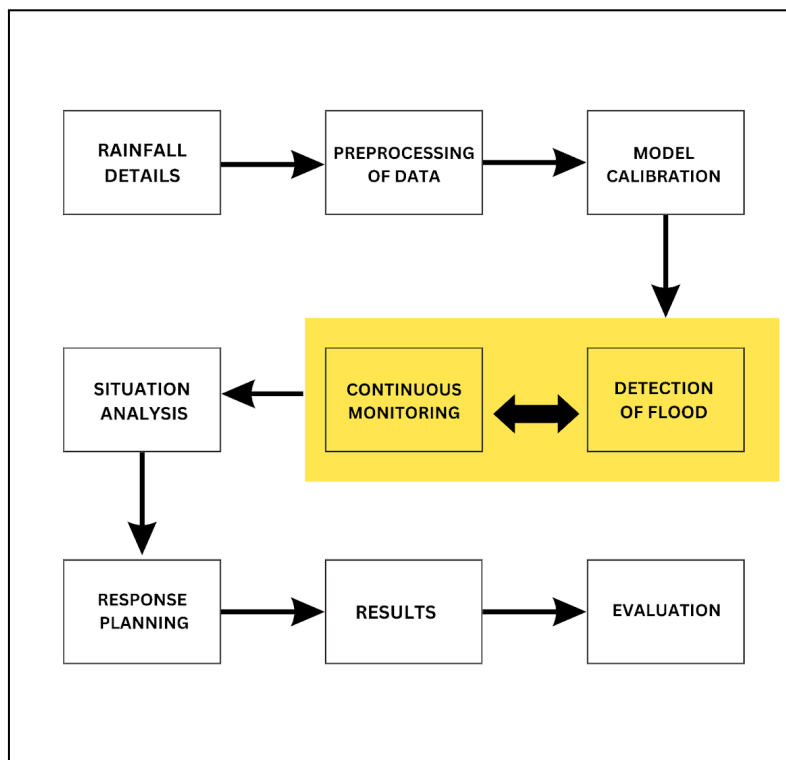


Fig.3.1.Architectural Framework

The architecture flow of our project goes as follows:

Rainfall details :

The rainfall details are fetched from https://wrd.guj.nic.in/dam/hour_reports_h.php which is a government website

The parameters fetched and used are:

1. Scheme Id: This is a unique identifier for each scheme associated with the dam.
2. Scheme Name: The name assigned to each scheme or section of the dam.
3. Design Gross Storage (MCM): The total storage capacity of the dam, typically measured in Million Cubic Meters (MCM). It represents the maximum amount of water the dam can hold when filled to capacity.
4. FRL (Full Reservoir Level) (m): The maximum level to which the reservoir behind the dam can be filled. It indicates the highest water level before spillage occurs.

5. Rule Level (01-08) (m): Rule levels refer to different predetermined water levels within the reservoir, each corresponding to specific operational rules or guidelines for managing water releases.
6. Present Water Level (m): The current water level within the reservoir, measured in meters.
7. Present Gross Storage (MCM): The current gross storage of water in the reservoir, measured in Million Cubic Meters (MCM).
8. Percentage Storage: The percentage of the total storage capacity currently occupied by water.
9. Inflow (Cusecs): The rate of water flowing into the reservoir, typically measured in cubic feet per second (Cusecs).
10. Outflow River (Cusecs): The rate of water being released from the reservoir into the river downstream, measured in cubic feet per second (Cusecs).
11. Outflow Canal (Cusecs): The rate of water being released from the reservoir into irrigation canals, measured in cubic feet per second (Cusecs).
12. Cumm. Rainfall (mm): The cumulative rainfall received in the area surrounding the reservoir, typically measured in millimeters.
13. Type: Indicates whether the scheme is associated with a river (G) or an underground (UG) water source.
14. Gate Position Nos.: The number of gates or openings available for controlling water flow in the scheme.
15. Opening (m): The current opening or position of the gates, measured in meters.

These parameters provide valuable insights into the current status and operational aspects of the Sardar Sarovar Dam, facilitating effective management of water resources and flood control measures in the region.

Preprocessing:

Preprocessing of the data for overflow prediction using machine learning typically involves several steps to prepare the data for model training. Here's how the preprocessing could be done for overflow prediction using the provided data:

1. **Data Cleaning:** Check for any missing or inconsistent values in the dataset and handle them appropriately. This involves imputing missing values or removing rows with missing data.
2. **Feature Selection:** Identify the relevant features (parameters) that are most likely to influence overflow prediction. Based on domain knowledge or through feature importance analysis, selecting a subset of features for model training.
3. **Normalization/Scaling:** Scaling the numerical features to a similar range to ensure that no single feature dominates the model training process. Techniques like Min-Max scaling or Standardization (Z-score normalization) are applied.
4. **Encoding Categorical Variables:** Categorical variables such as the "Type" of scheme (G for river, UG for underground), encoded into numerical values using techniques like one-hot encoding or label encoding.
5. **Feature Engineering:** Creating new features if necessary based on domain knowledge or insights from the data. Calculating the difference between the current water level and the full reservoir level could be a relevant feature for overflow prediction.
6. **Train-Test Split:** Split the dataset into training and testing sets. The training set will be used to train the machine learning model, while the testing set will be used to evaluate its performance. We have split our dataset into 80% for training and 20% for testing.
7. **Handling Time-Series Data:** The data included a time component (e.g., cumulative rainfall over time), consider incorporating time-series analysis techniques such as lagging or rolling window statistics to capture temporal patterns.
8. **Handling Imbalanced Data:** The imbalanced dataset (i.e., one class significantly outweighs the other), employed techniques like oversampling, undersampling, or using weighted classes during model training to address the imbalance.
9. **Feature Scaling:** Applying feature scaling techniques to ensure that all features have a similar scale. This step helps the model converge faster during training.

10. Data Splitting: Splitting the dataset into training, validation, and testing sets. The training set is used to train the model, the validation set is used to tune hyperparameters, and the testing set is used to evaluate the model's performance on unseen data.

By performing these preprocessing steps, the data is transformed into a format suitable for training machine learning models for overflow prediction. This prepared dataset is used to train various machine learning algorithms such as decision trees, random forests, support vector machines, or neural networks for accurate prediction of overflow events at the Sardar Sarovar Dam..

Model Calibration

1. Calibration Process Overview: The calibration process involves fine-tuning model parameters to optimize predictive performance.
2. Hyperparameter Tuning: Hyperparameter tuning entails adjusting parameters such as learning rates or regularization strengths to enhance model accuracy.
3. Performance Evaluation Metrics: Various metrics like accuracy, precision, recall, and F1 score are employed to assess the effectiveness of the calibrated model.
4. Cross-Validation Techniques: Cross-validation methods like k-fold and leave-one-out validation ensure robustness by validating the model on multiple subsets of the data.
5. Iterative Adjustment: Iteratively adjusting parameters based on validation results refines the model's performance gradually.
6. Application to Machine Learning Models: The calibrated parameters are applied to machine learning algorithms like logistic regression or random forests.
7. Validation and Satisfactory Performance: The model's performance is validated against unseen data, ensuring satisfactory results before deployment.
8. Final Parameter Application: Once validated, the calibrated parameters are applied to the model for accurate flood prediction.

Detection of flood and continuous monitoring:

1. Real-Time Data Acquisition: Continuous monitoring involves the acquisition of real-time data from various sources such as weather stations, river gauges, and satellite imagery.
2. Threshold Setting: Threshold levels for water levels and rainfall are established based on historical data and local conditions to trigger flood detection.
3. Data Processing: Real-time data streams are processed using algorithms to detect anomalies or sudden changes indicative of potential flooding events.
4. Machine Learning Algorithms: Machine learning algorithms such as neural networks or decision trees are employed to analyze incoming data and detect patterns associated with flooding.
5. Early Warning Systems: Detected anomalies trigger early warning systems, which alert authorities and residents in at-risk areas to take necessary precautions.
6. Automated Alerts: Automated alerts are generated and disseminated through various communication channels such as SMS, mobile apps, and sirens to notify the population about potential flood threats.

Response and planning system :

1. Advises the disaster relief team to take the appropriate action
2. The system assesses the level of risk based on the processed data, historical records, and predefined thresholds for different environmental parameters.
3. After an event, the system conducts post-event analysis to evaluate the effectiveness of response actions and identify areas for improvement in future planning and execution.
4. Based on the analysis and feedback received, the system adapts and improves its response strategies and plans to enhance preparedness for future events.

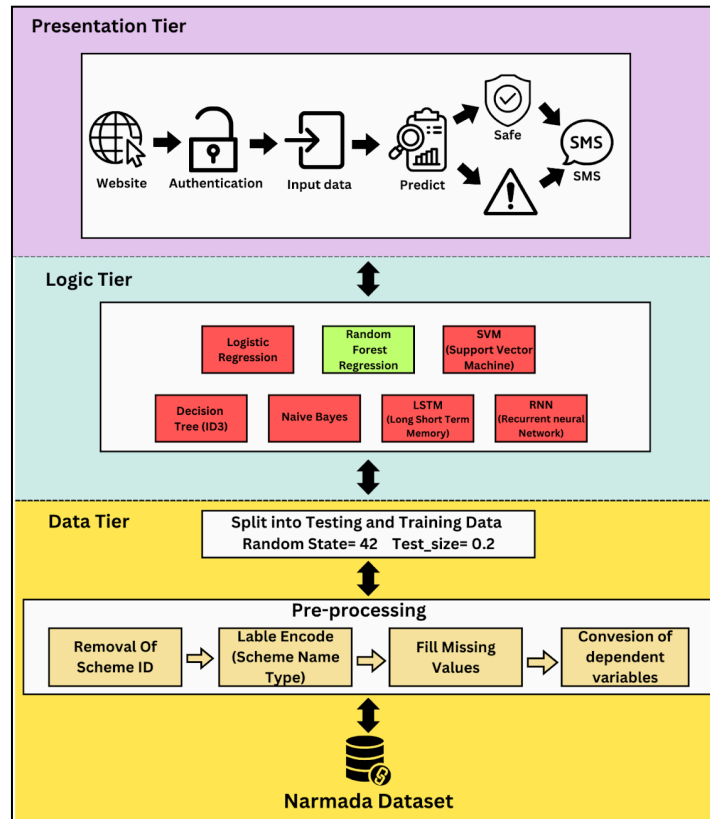


Figure.3.2.Methodology

1. Data Tier:

Narmada Dataset: This is the primary dataset containing information about the Sardar Sarovar Dam and its associated schemes. It includes attributes such as scheme name, design gross storage, FRL (Full Reservoir Level), present water level, present gross storage, percentage storage, inflow, outflow, rainfall, etc.

Preprocessing:

1. **Removal of Scheme ID:** Since Scheme ID serves as an identifier and doesn't contribute to model training, it is removed.
2. **Label Encoding:** Converting categorical variables into numerical format, if needed.
3. **Missing Values Handling:** Filling missing values using techniques like mean, median, or interpolation.
4. **Conversion of Dependent Variables:** Ensuring dependent variables are in a suitable format for model training.
5. **Data Splitting:** Splitting the dataset into training and testing subsets, with a test size of 0.2 and a random state of 2 to ensure reproducibility.

2. Logic Tier:

Machine Learning Models:

1. LSTM (Long Short-Term Memory): A type of recurrent neural network (RNN) capable of learning long-term dependencies. It's suitable for sequence prediction tasks.
2. RNN (Recurrent Neural Network): Similar to LSTM, RNNs are used for sequential data processing.
3. Naive Bayes: A probabilistic classifier based on Bayes' theorem with the assumption of independence between features.
4. ID3 (Iterative Dichotomiser 3): A decision tree algorithm that uses information gain to split the data.
5. SVM (Support Vector Machine): A supervised learning algorithm used for classification and regression tasks by finding the hyperplane that best separates classes.
6. Random Forest: An ensemble learning method that constructs a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction Model Selection: Random Forest is chosen as the suitable model for overflow prediction based on its performance metrics and suitability for the problem domain.

3. Presentation Tier:

Website Interface:

1. Authentication: User authentication to ensure authorized access to the system.
2. Input Data: Interface for users to input data related to the dam and its associated parameters.
3. Prediction Result: Displaying the predicted overflow status based on the input data using the trained Random Forest model.
4. Alert-SMS: Sending alert SMS to relevant stakeholders in case of predicted overflow events, ensuring timely response and mitigation efforts.

3.3 Algorithm and Process Design

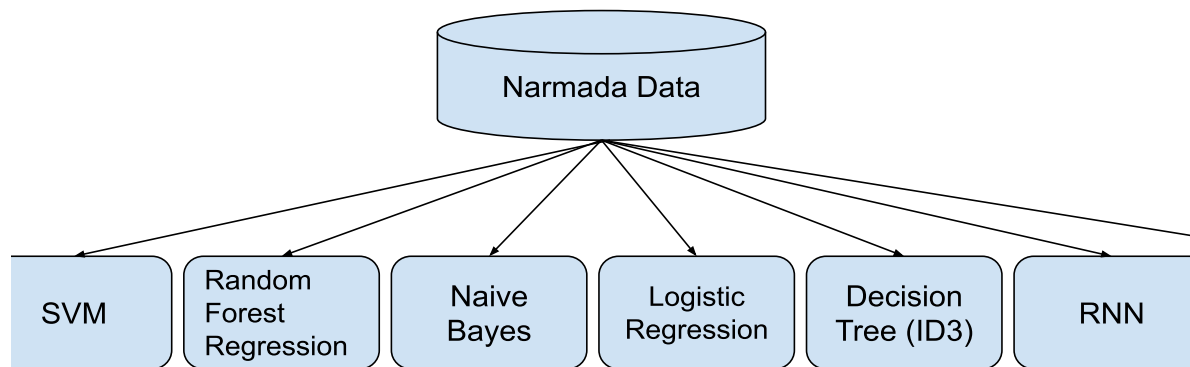


Figure 3.3 Models Used

The following algorithms were implemented on the preprocessed dataset :

1. Support Vector Machine (SVM)
2. Recurrent Neural Networks (RNNs)
3. Random Forest
4. Logistic Regression
5. Long Short-Term Memory (LSTM)
6. Naive Bayes
7. Decision Tree (ID3)

Different Machine Learning algorithms are implemented to check their accuracy and precision

1. Support Vector Machine (SVM):
 - a. SVM aims to find the hyperplane that best separates data points into different classes. In the case of regression, it finds the hyperplane that best fits the data points.
 - b. Methodology: SVM works by finding the hyperplane with the maximum margin, i.e., the maximum distance between the hyperplane and the nearest data points (support vectors). It can handle both linear and non-linear classification and regression tasks using different kernel functions such as linear, polynomial, or radial basis function (RBF) kernels.

2. Recurrent Neural Networks (RNNs):

- a. RNNs process sequential data by maintaining a hidden state that captures information about previous inputs. The output at each time step is dependent on the current input and the previous hidden state.
- b. Methodology: RNNs are designed to handle sequential data like time series or natural language. They use feedback loops to persist information over time, making them suitable for tasks where context or temporal dependencies are important. However, they suffer from vanishing gradient problems when dealing with long sequences.

3. Random Forest:

- a. Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode (classification) or mean prediction (regression) of individual trees.
- b. Methodology: Random Forest builds multiple decision trees based on random subsets of the training data and features. It reduces overfitting by averaging predictions from multiple trees. Random Forest is versatile, robust to overfitting, and handles high-dimensional data well.

4. Logistic Regression:

- a. Logistic Regression models the probability of a binary outcome using a logistic function, which maps the input features to a probability between 0 and 1.
- b. Methodology: Logistic Regression is a linear classification algorithm that estimates the probability that a given input belongs to a particular class. It's commonly used for binary classification tasks and can be extended to handle multi-class classification using techniques like one-vs-rest or softmax regression.

5. Long Short-Term Memory (LSTM):

- a. LSTM is a type of RNN with additional gating mechanisms (forget gate, input gate, output gate) that control the flow of information through the network, allowing it to capture long-term dependencies.
- b. Methodology: LSTM addresses the vanishing gradient problem in traditional RNNs by maintaining a cell state that can retain information over long sequences. It's widely used for tasks involving long-range dependencies, such as speech recognition, language translation, and time series prediction.

6. Naive Bayes:

- a. Naive Bayes applies Bayes' theorem with the "naive" assumption of feature independence, meaning each feature contributes independently to the probability of the outcome.
- b. Methodology: Naive Bayes is a probabilistic classifier that calculates the likelihood of a class given the input features using conditional probability. Despite its simplicity and assumption of feature independence, Naive Bayes often performs well in practice, especially for text classification tasks.

7. Decision Tree (ID3):

- a. Decision Tree recursively splits the dataset based on the feature that provides the maximum information gain or minimum impurity at each node.
- b. Methodology: ID3 (Iterative Dichotomiser 3) is one of the earliest decision tree algorithms. It selects the best attribute to split the data at each node based on information gain, which measures the reduction in entropy or impurity. Decision Trees are interpretable, easy to understand, and can handle both numerical and categorical data.

3.4 Methodology applied

3.4.1 Algorithm

Process Design:

Random forest operates by constructing a multitude of decision at training time and outputting the class that's the mode of the classes or mean prediction

A random forest is a meta-estimator (i.e. it combines the result of multiple predictions), which aggregates many decision trees with some helpful modifications:

1. The number of features that can be split at each node is limited to some percentage of the total (which is known as the hyper-parameter). This limitation ensures that the ensemble model does not rely too heavily on any individual feature and makes fair use of all potentially predictive features.
2. Each tree draws a random sample from the original data set when generating its splits, adding a further element of randomness that prevents overfitting.

Algorithm:

1. Data Preparation:

- a. Ensuring dataset is prepared with the necessary preprocessing steps such as removal of Scheme ID, label encoding, handling missing values, and splitting into training and testing sets.
- b. Making sure the target variable (percentage storage) is separated from the predictor variables.

2. Model Training:

- a. Training the Random Forest model using the training dataset.

In Python, we have used libraries like scikit-learn.

Code:

```
from sklearn.ensemble import RandomForestRegressor
# Instantiate the model
rf_model = RandomForestRegressor()
# Train the model
rf_model.fit(X_train, y_train) # X_train: predictor variables, y_train: target variable
```

3. Prediction:

- a. Once the model is trained, use it to make predictions on the testing dataset.

Code:

```
# Make predictions  
y_pred = rf_model.predict(X_test) # X_test: predictor variables of the testing dataset
```

4. Overflow Prediction:

Calculate the percentage storage threshold (90%) based on the maximum storage capacity of the dam.

Code:

```
threshold = 0.9 * max_percentage_storage_capacity
```

5. Compare the predicted percentage storage values with the threshold to identify overflow events.

Code:

```
overflow_indices = [i for i, pred in enumerate(y_pred) if pred >= threshold]
```

3.4.2 Dataset :

In our real-time study, we extracted data from Narmada river reservoirs, sourced from https://wrd.guj.nic.in/dam/hour_reports_h.php. This dataset comprised essential parameters such as Scheme Id, Scheme Name, Design Gross Storage, FRL, Rule Level, Present Water Level, Present Gross Storage, Percentage Storage, Inflow (Cusecs), Outflow River (Cusecs), Outflow Canal (Cusecs), Cumulative Rainfall (mm), Type, and Gate Position Numbers with their corresponding openings (m). After preprocessing, we stored this refined dataset in a CSV file.

Our objective was to predict overflow conditions for the current day using various machine learning algorithms. We applied Support Vector Machine (SVM), Recurrent Neural Networks (RNNs), Random Forest, Logistic Regression, Naive Bayes, and Decision Tree (ID3) algorithms to the dataset. By evaluating the accuracy and precision of each algorithm, we aimed to determine the most suitable model for our specific dataset.

This approach involved a rigorous analysis of the reservoir data, leveraging advanced machine learning techniques to ensure accurate predictions regarding overflow conditions. The algorithms were chosen based on their appropriateness for handling the dataset's characteristics and complexities, allowing for a comprehensive evaluation of their performance. The results of this analysis were vital for decision-making and ensuring the effective management of reservoir overflow situations.

3.4.3 Preprocessing:

1. Removal of SchemeID:
 - a. Reason: SchemeID might not significantly contribute to the prediction of reservoir overflow conditions. Removing it simplifies the dataset.
 - b. Implementation: Drop the 'SchemeID' column from the dataset.
2. Label Encoding (SchemeName, Type):
 - a. Reason: Machine learning models work with numerical data. Label encoding converts categorical variables like 'SchemeName' and 'Type' into numerical values without introducing the ordinal relationship between categories.
 - b. Implementation: Use label encoding techniques to convert 'SchemeName' and 'Type' columns into numerical values.
3. Filling Missing Values:
 - a. Reason: Missing data can hinder the performance of machine learning models. Filling or imputing missing values ensures that all features have complete data.
 - b. Implementation: Use techniques like mean, median, or interpolation to fill missing values in relevant columns.
4. Conversion of Dependent Variable into 0 and 1:
 - a. Reason: For binary classification tasks like predicting overflow conditions (0 for no overflow, 1 for overflow), converting the dependent variable into binary values is essential.
 - b. Implementation: Use threshold values to convert continuous 'Overflow' into binary values (0 and 1) based on predefined criteria.

This preprocessing process ensures that the dataset is cleaned, standardized, and ready for further analysis and model training. It addresses issues such as missing data and categorical variables, making the dataset suitable for machine learning algorithms to generate accurate predictions of reservoir overflow conditions.

3.5 Hardware & Software Specifications

3.5.1 Software Specifications

1. **Python 3.*:** Python 3.* refers to the latest versions in the Python 3 series, a popular and versatile programming language known for its readability and ease of use, used for various applications from web development to data analysis.
2. **Tensorflow:** TensorFlow is an open-source machine learning library developed by Google, widely used for building and training deep learning models, making it a vital tool in artificial intelligence and data science.
3. **NumPy:** NumPy is a fundamental Python library for numerical and scientific computing, enabling efficient handling of large arrays and matrices, essential for data manipulation and scientific tasks.
4. **Pandas:** Pandas is a Python library specializing in data manipulation and analysis, offering data structures like data frames and various functions for tasks such as data cleaning and transformation.
5. **Scikit-Learn (Sklearn):** Scikit-Learn is a widely-used machine learning library in Python, providing a user-friendly platform for developing and evaluating machine learning models across various applications.
6. **Tkinter:** Tkinter is Python's standard GUI library, allowing developers to create desktop applications with graphical user interfaces in a straightforward manner.
7. **Google Colab:** Google Colab is a cloud-based Jupyter notebook environment that enables collaborative and interactive Python coding, widely used in data science and machine learning for its access to free cloud resources, including GPUs and TPUs.
8. **Django:** Django is a back-end server side web framework. Django is free, open source and written in Python. Django makes it easier to build web pages using Python.
9. **Twilio API :** Twilio API is a cloud communications platform that enables developers to integrate voice, SMS, and messaging functionality into their applications using simple HTTP APIs. It allows for programmable communication solutions, such as sending text messages, making voice calls, and creating chatbots, enhancing user engagement and experience.

3.5.2 Hardware Specification

- 1. Processor:** Intel(R) Core(TM) i3-1005G1 CPU @ 1.20GHz 1.19 GHz
- 2. Installed RAM:** 8.00 GB (7.79 GB usable)
- 3. Operating System:** Windows 11 Home Single Language
- 4. Version:** 22H2
- 5. OS Build:** 22621.2283

3.6 Experiment and Results for Validation and Verification

3.6.1 Support Vector Machine

The Support Vector Machine (SVM) is a supervised machine learning model capable of both classification and regression tasks. It employs the kernel trick to transform data and, based on these transformations, establishes an optimal boundary (hyperplane) to differentiate between possible outputs. This boundary is determined using support vectors, which are points nearest to the line. The distance between these support vectors and the line forms a margin, and the wider this margin, the better. The optimal hyperplane is the one that forms the widest margin. Consequently, it can be observed that SVM is inherently linear in nature, which can be considered a drawback, as it implies that SVM may not perform effectively for nonlinear data.

In Support Vector Machine (SVM), the decision boundary is determined by a linear function of the form:

$$f(x)=w^Tx+b$$

Where:

1. $f(x)$: represents the decision function that determines the class label of the input x .
2. x : represents the input feature vector.
3. w : represents the weight vector, which determines the orientation of the decision boundary.
4. b : represents the bias term, which shifts the decision boundary away from the origin.

For binary classification, the class label y is determined based on the sign of $f(x)$:

1. If $f(x) \geq 0$, then $y=1$ (positive class).
2. If $f(x) < 0$, then $y=-1$ (negative class).

In the context of our dataset:

1. $f(x)$: This represents the decision function that determines the likelihood of reservoir overflow.
2. x : These are the input features (independent variables) from our dataset.
3. w : This represents the weight vector, which determines the orientation of the decision boundary. Each element of w corresponds to a feature in x , indicating its importance in determining the decision boundary.
4. b : This represents the bias term, which shifts the decision boundary away from the origin along the direction defined by w .

1. $f(x)$: Decision function determining the likelihood of reservoir overflow.
2. x : Independent variables/features from our dataset.
3. w : Weight vector, indicating the importance of each feature in determining the decision boundary.
4. b : Bias term, shifting the decision boundary away from the origin.

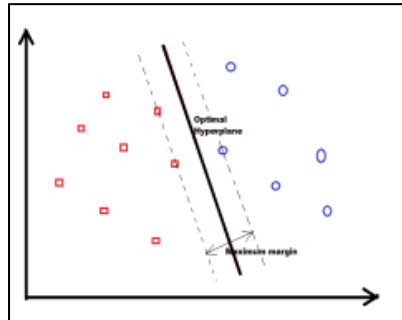


Figure 3.4.1 SVM Working

Accuracy: 0.9411764705882353					
Classification Report:					
	precision	recall	f1-score	support	
0	1.00	0.83	0.91	6	
1	0.92	1.00	0.96	11	
accuracy			0.94	17	
macro avg	0.96	0.92	0.93	17	
weighted avg	0.95	0.94	0.94	17	

Figure 3.4.2 Classification Report of SVM

3.6.2 Logistic Regression

Logistic regression serves as a supervised algorithm primarily employed for classification tasks. It is particularly suitable when the target variable is categorical and labeled. Logistic regression utilizes a simple regression formula for prediction and subsequently applies a threshold (such as 0.5) to classify the predicted value as either 0 or 1. The regression formula employed is:

$$y = mx + b.$$

1. $P(y=1|x)$: Probability of reservoir overflow.
2. x : Independent variables from our dataset.
3. m : Coefficients or weights assigned to each independent variable.
4. b : Intercept or bias term.

Accuracy: 0.7647058823529411					
Classification Report:					
	precision	recall	f1-score	support	
0	0.60	1.00	0.75	6	
1	1.00	0.64	0.78	11	
accuracy			0.76	17	
macro avg	0.80	0.82	0.76	17	
weighted avg	0.86	0.76	0.77	17	

Figure 3.4.3 Classification Report of Logistic Regression

3.6.3 Random Forest Regression

Random Forest is a supervised classification algorithm that is applicable to both classification and regression problems. Its advantages encompass the ability to handle missing values, model classification for categorical values, and mitigate overfitting even when utilizing a larger number of trees in the forest [8]. The mathematical formula for Random Forest can be expressed as:

$$ni_j = W_j C_{j-left(j)} C_{left(j)} - W_{right(j)} C_{right(j)}$$

Accuracy: 0.9411764705882353				
Classification Report:				
	precision	recall	f1-score	support
0	1.00	0.83	0.91	6
1	0.92	1.00	0.96	11
accuracy			0.94	17
macro avg	0.96	0.92	0.93	17
weighted avg	0.95	0.94	0.94	17

Figure 3.4.4 Classification Report of Random Forest Regression

3.6.4 Decision Tree

The Decision Tree is a supervised machine learning algorithm primarily utilized for classification tasks. It adopts a tree-like structure wherein each node serves as a test on a specific attribute, each branch represents an outcome from that node, and each leaf node signifies the class label corresponding to the path it leads to. The significance of a node in decision-making diminishes as we descend from the top to the bottom of the tree. Various methods exist for constructing a decision tree, with one of the most significant being CART (Classification and Regression Tree), which employs the Gini index as a metric.

1) CART (Classification and Regression Tree) which uses the Gini index as a metric.

Formula for Gini Index:

$$\text{gini}_A(D) = \frac{|D_1|}{D} \text{gini}(D_1) + \frac{|D_2|}{D} \text{gini}(D_2)$$

2) ID3 (Iterative Dichotomiser which uses Entropy Function and Information Gain as metrics).

Formula for Information Gain:

$$I(P_i, N_i) = \frac{-p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

Formula for Entropy:

$$\Sigma \left(\frac{P_i + N_i}{p+n} \right) I(P_i, N_i)$$

For both CART (Classification and Regression Tree) and ID3 (Iterative Dichotomiser), the parameters in the provided dataset would be used to determine the splits at each node of the decision tree based on the impurity measures (Gini index for CART and Entropy/Information Gain for ID3). Let's map the parameters to each variable in the impurity measures:

1. For Gini Index (used in CART):

D : Represents the dataset at a particular node of the decision tree.

D1 and D2 : Represent the subsets of the dataset D after splitting on a particular feature.

|D1| and |D2| : Represent the number of samples in subsets D1 and D2, respectively.

gini(D) : Represents the Gini impurity of the dataset D.

The parameters in our dataset, such as 'Design Gross Storage', 'FRL', 'Rule Level', 'Present Water Level', 'Present Gross Storage', 'Inflow', 'Outflow River', 'Outflow Canal', 'Cumm. Rainfall', 'Type', 'Gate Position Nos.', and 'Opening', would be used to calculate the Gini impurity for each node during the construction of the decision tree.

2. For Information Gain and Entropy (used in ID3):

P_i : Represents the proportion of samples in class i (positive or negative) at a particular node of the decision tree.

N_i : Represents the total number of samples in the dataset at a particular node.

The parameters in our dataset would be used to calculate the proportions of positive and negative class labels (or target values) at each node, which are then used to compute Information Gain or Entropy for selecting the best split.

Accuracy: 0.9411764705882353				
Classification Report:				
	precision	recall	f1-score	support
0	1.00	0.83	0.91	6
1	0.92	1.00	0.96	11
accuracy			0.94	17
macro avg	0.96	0.92	0.93	17
weighted avg	0.95	0.94	0.94	17

Figure 3.4.5 Classification Report of Decision Tree

3.6.5 Naive Bayes

Naive Bayes is a probabilistic machine learning algorithm widely applicable in various classification tasks. It stands out for its ease of implementation and quick prediction capabilities, owing to its probabilistic nature. Based on Bayes' theorem, it excels particularly in scenarios with high-dimensional inputs. The theorem is represented as follows:

$$P(A|B)=P(B|A).P(A)/P(B)$$

Here A and B are two events and $P(A|B)$ is the conditional probability that event A occurs, given that event B has occurred. $P(B|A)$ is the conditional probability that event B occurs, given that event A has occurred. $P(A)$ and $P(B)$: Probability of A and B without regard to each other.

The Naive Bayes model offers several advantages, such as rapid computation for both training and prediction, straightforward probabilistic predictions, and ease of interpretation.

1. A: Represents the class label or target variable that you are trying to predict using the given features. In this case, since you are interested in predicting the percentage storage, A would represent the percentage storage variable.
2. B: Represents the set of features or independent variables used to predict the class label A. In this case, B would represent all the other parameters/features in our dataset except for the percentage storage. These features include 'Scheme Name', 'Design Gross Storage', 'FRL', 'Rule Level', 'Present Water Level', 'Present Gross Storage', 'Inflow', 'Outflow River', 'Outflow Canal', 'Cumm. Rainfall', 'Type', 'Gate Position Nos.', and 'Opening'.

Accuracy: 0.9411764705882353				
Classification Report:				
	precision	recall	f1-score	support
0	0.86	1.00	0.92	6
1	1.00	0.91	0.95	11
accuracy			0.94	17
macro avg	0.93	0.95	0.94	17
weighted avg	0.95	0.94	0.94	17

Figure 3.4.6 Classification Report of Naive Bayes

3.6.6 LSTM Neural Network

A conventional RNN operates with a single hidden state passed through time, posing challenges in learning long-term dependencies. To overcome this, LSTM (Long Short-Term Memory) networks introduce a memory cell, capable of retaining information over extended periods. This design enables LSTM networks to effectively capture long-term dependencies in sequential data, rendering them suitable for tasks like language translation, speech recognition, and time series forecasting. Moreover, LSTMs can be integrated with other neural network architectures, such as Convolutional Neural Networks (CNNs), to analyze images and videos.

Activation at time t :

$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j\tilde{h}_t^j$$

Update gate:

$$z_t^j = \sigma(W_z x_t + U_z h_{t-1})^j$$

Candidate activation:

$$\tilde{h}_t^j = \tanh(Wx_t + U(rt \otimes h_{t-1}))^j$$

Reset gate:

$$r_t^j = \sigma(W_r x_t + U_r h_{t-1})^j$$

Where:

Activation at time t (h_t^j)

The hidden state at time t for the j -th unit or neuron in the LSTM network.

Update gate(z_t^j)

The update gate at time t for the j -th unit or neuron in the LSTM network. It controls how much of the previous cell state (h_{t-1}^j) is retained and how much of the new candidate value (\tilde{h}_t^j) is added to the current cell state.

Candidate activation(\tilde{h}_t^j):

The candidate activation at time t for the j -th unit or neuron in the LSTM network. It represents the new information that could be added to the cell state at time t .

Reset gate (r_t^j):

The reset gate at time t for the j -th unit or neuron in the LSTM network. It controls how much of the previous hidden state (h_{t-1}^j) is forgotten or reset when computing the candidate activation (\tilde{h}_t^j)

Activation at time t :

- a. Represents the hidden state at time t for the j -th unit or neuron in the LSTM network.
- b. In the context of our dataset, h_{tj} would represent the hidden state of the LSTM network, which could be interpreted as the internal representation learned by the LSTM based on the input parameters.

Update gate(z_t^j)

- c. Represents the update gate at time t for the j -th unit or neuron in the LSTM network.
- e. Controls how much of the previous cell state (h_{t-1}^j) is retained and how much of the new candidate value (\tilde{h}_t^j) is added to the current cell state.
- f. In the context of our dataset, (z_t^j) would be computed using the provided parameters ('Design Gross Storage', 'FRL', 'Rule Level', 'Present Water Level', 'Present Gross Storage', 'Inflow', 'Outflow River', 'Outflow Canal', 'Cumm. Rainfall', 'Type', 'Gate Position Nos.', and 'Opening') and the weights and biases of the LSTM network.

Candidate activation(\tilde{h}_t^j):

- Represents the candidate activation at time t for the j -th unit or neuron in the LSTM network.
- Represents the new information that could be added to the cell state at time t .
- In the context of our dataset,
- (\tilde{h}_t^j) would be computed using the provided parameters ('Design Gross Storage', 'FRL', 'Rule Level', 'Present Water Level', 'Present Gross Storage', 'Inflow', 'Outflow River', 'Outflow Canal', 'Cumm. Rainfall', 'Type', 'Gate Position Nos.', and 'Opening') and the weights and biases of the LSTM network.

Reset gate (r_t^j):

- Represents the reset gate at time t for the j -th unit or neuron in the LSTM network.
- Controls how much of the previous hidden state is forgotten or reset when computing the candidate activation. Not explicitly mentioned in the provided LSTM formula, but it's a common component in LSTM architectures.

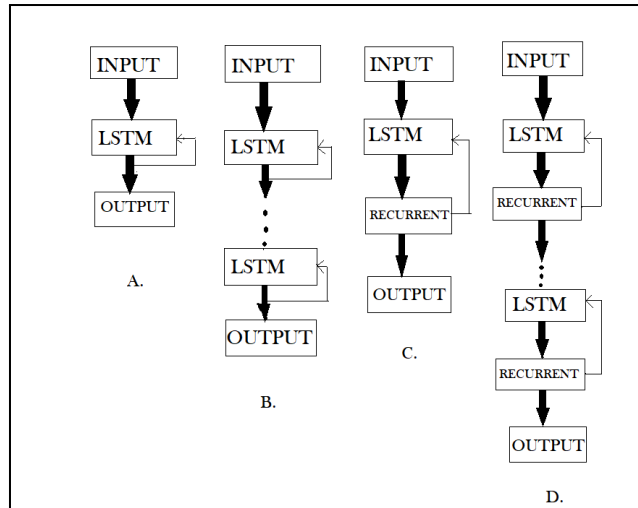


Figure 3.4.7 Working of LSTM

Accuracy: 0.8823529411764706				
Classification Report:				
	precision	recall	f1-score	support
0	0.83	0.83	0.83	6
1	0.91	0.91	0.91	11
accuracy			0.88	17
macro avg	0.87	0.87	0.87	17
weighted avg	0.88	0.88	0.88	17

Figure 3.4.8 Classification Report of LSTM

3.6.7 Prediction Using RNN

A Recurrent Neural Network (RNN) is an algorithm equipped with internal memory, allowing it to retain information from previous inputs. It is specifically designed for handling sequential data. Unlike Feed Forward Neural Networks, where information flows only in one direction, RNN incorporates feedback loops, enabling it to consider both the current input and the knowledge acquired from past inputs when making decisions. This characteristic sets RNN apart from other algorithms. RNN applies weights to both the current input and the previous input, adjusting these weights through gradient descent and backward propagation. The mathematical formulas for RNN are as follows:

$$o^t = f(h^t, \theta)$$

$$h^t = k(h^{t-1}, x^t, \theta)$$

Where, o^t is the output produced at time t ,

h^t is the state of hidden layers at time t ,

x^t is the input given at time t ,

θ indicates the weights and biases for that network

1. Output at time t (o_t):

- Represents the output produced at time t .

- In the context of our dataset, o_t would be the output prediction or classification result generated by the RNN at time t .

2. State of hidden layers at time t (h_t):

- Represents the state of hidden layers at time t .

- In the context of our dataset, h_t would be the hidden state or internal representation learned by the RNN at time t .

3. Input given at time t (x_t):

- Represents the input provided at time t .

- In the context of our dataset, x_t would be the input features or parameters given at time t , which include 'Design Gross Storage', 'FRL', 'Rule Level', 'Present Water Level', 'Present Gross Storage', 'Inflow', 'Outflow River', 'Outflow Canal', 'Cumm. Rainfall', 'Type', 'Gate Position Nos.', and 'Opening'.

4. Function f and k :

- f represents the function that computes the output at time t based on the hidden state h_t .
- k represents the function that updates the hidden state h_t based on the previous hidden state h_{t-1} and the input x_t .
- The specific functions f and k depend on the architecture and parameters of the RNN, including the weights and biases.

Therefore, in the RNN formula applied to our dataset, each parameter contributes to the computation of the output and hidden state at each time step. The RNN learns temporal dependencies and makes predictions based on the input parameters, with the output at each time step being influenced by the hidden state and the input at that time step.

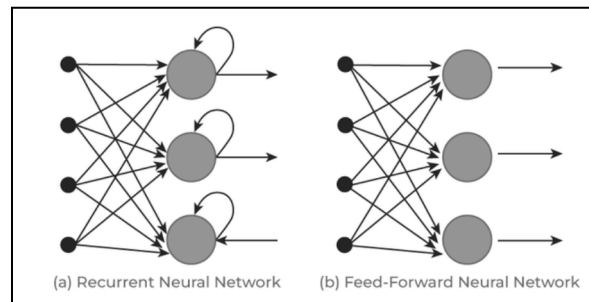


Figure 3.4.9 Working of Recurrent Neural Network

Accuracy: 0.7058823529411765				
Classification Report:				
	precision	recall	f1-score	support
0	0.67	0.33	0.44	6
1	0.71	0.91	0.80	11
accuracy			0.71	17
macro avg	0.69	0.62	0.62	17
weighted avg	0.70	0.71	0.67	17

Figure 3.4.10 Classification report of Recurrent Neural Network

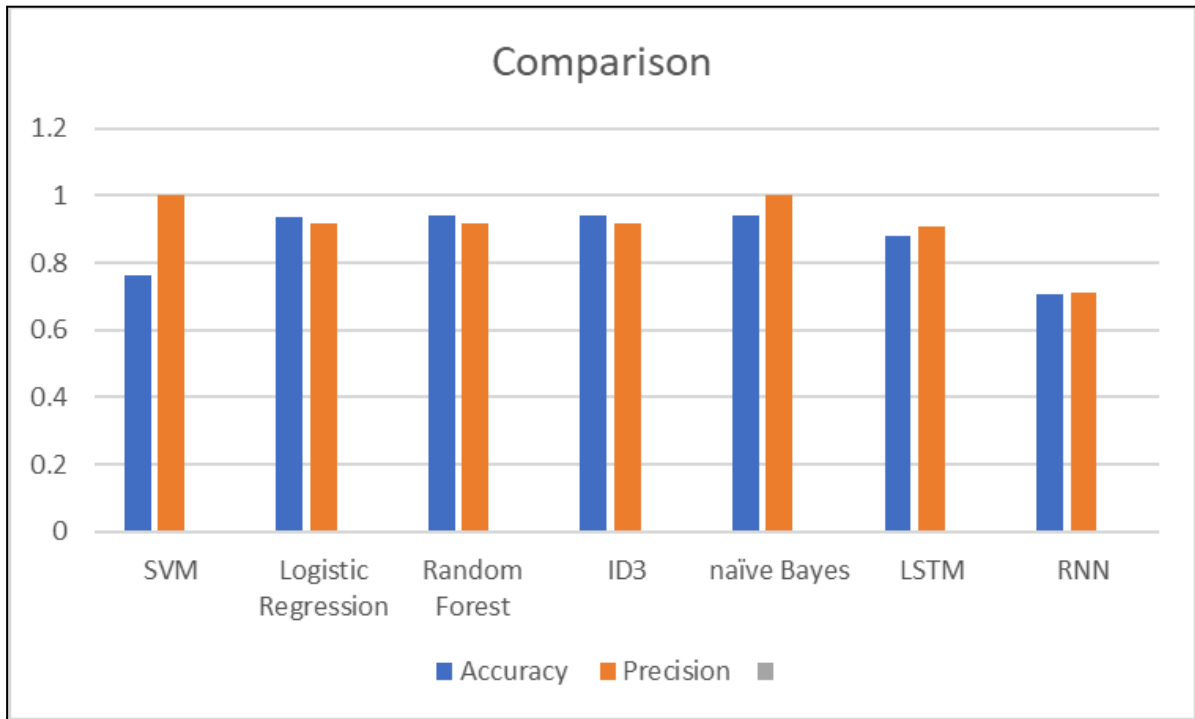


Figure 3.5 Comparison of different machine learning algorithms

Based on the experiment results for flood prediction models, it's evident that Random Forest Regression consistently outperforms other models in terms of both accuracy and precision. Support Vector Machine achieved an accuracy of 76.4% with a precision of 1.0, while Logistic Regression, Decision Tree (ID3), and Naive Bayes attained an accuracy of 94.1% with a precision of 0.92. LSTM achieved an accuracy of 88.2% with a precision of 0.91, and RNN scored an accuracy of 70.5% with a precision of 0.71.

The choice of Random Forest Regression as the preferred model for flood prediction can be attributed to several factors. Firstly, Random Forest Regression demonstrates consistently high accuracy and precision across various experiments, indicating its robustness and reliability in predicting flood occurrences. Random Forest Regression is known for its ability to handle large datasets with high dimensionality, making it well-suited for analyzing complex flood prediction scenarios that involve numerous input parameters. Random Forest Regression inherently addresses issues such as overfitting and data noise, thanks to its ensemble learning approach, which aggregates predictions from multiple decision trees. Moreover, Random Forest Regression provides insights into feature importance, allowing stakeholders to identify key variables influencing flood occurrences and optimize mitigation strategies accordingly. Overall, the superior performance, versatility, and interpretability of Random Forest Regression make it the ideal choice for implementing flood prediction modules, ensuring accurate and timely risk assessment and management.

3.7 Result Analysis and Discussion

We have developed a website using Python Flask. A nodal officer logs in using their credentials and is shown only the schemes under their supervision. They input the real-time data of the scheme and get the prediction of whether the water level is safe or an overflow will occur. In both the cases, a notification via SMS will be sent to the desired user (i.e. the villagers living on the riverbank of that scheme)

The image shows a web form titled "Nodal Officer Login". It contains two input fields: "Username:" with the text "admin" entered, and "Password:" with five dots indicating a masked password. Below these fields is a blue button labeled "Login". The form is enclosed in a black rectangular border.

Figure 3.7.1 Nodal Officer Login

The Nodal Officer Login feature is a critical component of our website's authentication system, implemented using Firebase for robust security and seamless user experience. Through this feature, nodal officers are provided secure access to the platform, ensuring that only authorized personnel can interact with sensitive data and functionalities. One of the key aspects of this login mechanism is its role-based access control, where each nodal officer is granted access only to the schemes assigned to them. This granular access control ensures data confidentiality and integrity by limiting users' scope to only the relevant schemes they are responsible for. By leveraging Firebase's authentication capabilities, we guarantee that nodal officers can securely authenticate themselves, access their designated schemes, and perform necessary actions within their authorized scope, thereby enhancing efficiency, accountability, and data security within our system.

Resilient Rivers - Water Level Overflow Predictor

Design Gross Storage (MCM) 7414.29	FRL (m) 105.16
Rule Level (0-10) (m) 105.16	Present Water Level (m) 104.53
Present Gross Storage (MCM) 7048.69	Inflow (Cusecs) 39525.0
Outflow River (Cusecs) 21264.0	Outflow Canal (Cusecs) 800
Cumm. Rainfall (mm) 1382	Type of Gate: Gated
Gate Position Nos. 0	Opening (m) 0.00
Scheme: Likal	Predict

Figure 3.7.2 Predictor page

After successfully logging in, the nodal officer is seamlessly directed to the predictor page, where they are prompted to input real-time data pertaining to the dam's status. This predictor page serves as a pivotal tool for assessing the current condition of the dam and determining the likelihood of overflow events. The officer inputs crucial parameters such as water levels, storage capacities, inflow rates, outflow rates, and cumulative rainfall data into designated fields. Leveraging this real-time data, the system employs sophisticated algorithms to analyze the current state of the dam and assess the risk of overflow events. By incorporating predictive modeling techniques, historical data analysis, and real-time sensor inputs, the predictor page provides actionable insights into the probability of dam overflow occurrences. This proactive approach empowers the nodal officer to make informed decisions regarding dam management and mitigation strategies, thereby ensuring the safety of downstream communities and infrastructure. Additionally, the system can generate alerts or notifications if the risk of overflow surpasses predefined thresholds, enabling timely intervention and response measures. Overall, the predictor page serves as a pivotal tool for proactive dam management, facilitating effective decision-making and risk mitigation strategies in real time.

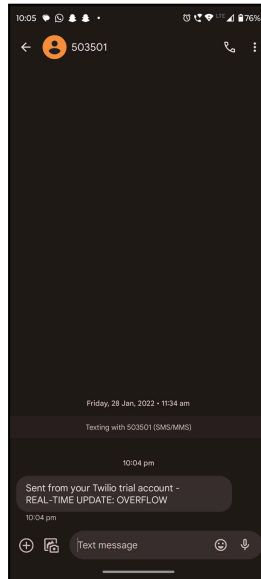


Figure 3.7.3 SMS sent via API

Following the assessment on the predictor page, if there's a determination of potential dam overflow, a crucial step is to alert registered users about the impending risk. This is achieved through the integration of the Twilio API, a powerful communication platform that facilitates SMS notifications. Once the system detects an elevated risk of overflow based on the real-time data inputs and predictive analysis, it triggers an automated process to send SMS alerts to pre-registered users. These users could include local authorities, emergency responders, residents in downstream areas, and other relevant stakeholders.

Utilizing the Twilio API, the system sends concise and informative SMS messages to registered users, notifying them of the imminent risk of dam overflow. The messages typically include crucial information such as the current dam status, the projected overflow risk level, recommended safety precautions, and instructions for evacuation if necessary. By promptly disseminating this vital information via SMS alerts, the system plays a pivotal role in enhancing community preparedness and ensuring timely response measures in the event of a dam overflow emergency.

The integration of Twilio API enables seamless communication with registered users, providing them with timely updates and actionable information to safeguard lives and property in vulnerable areas downstream of the dam. This proactive approach to emergency communication enhances overall disaster resilience and reinforces public safety measures in the face of potential natural disasters.

3.8 Conclusion Future Work

3.8.1 Conclusion

The proposed reservoir flood warning system uses advanced flood prediction models and real-time data integration to issue timely alerts and warnings. The system aims to enhance disaster preparedness and minimize flood risks. The system seeks to create a resilient approach to flood management. The successful implementation of the system will contribute to improved public safety, disaster response, and sustainable water resource planning.

3.8.2 Future Work

The system will be deployed on a website and will be available for the nodal officer of the reservoir or dam. The system can be integrated with other state-specific systems to share information and ensure everyone is on the same page. A mobile app may be developed to allow users to receive flood warnings and information on their smartphones. The system will be used to collect data on flood events and their impacts to improve flood prediction models. The system will be used to educate the public about flood risk and preparedness through public awareness campaigns and school programs.

References

- [1] Zhaoli Wang , Chengguang Lai , Xiaohong Chen , Bing Yang , Shiwei Zhao , Xiaoyan Bai, Flood hazard risk assessment model based on random forest, Year of Publication: August 2015.
- [2] Sunmin Lee, Jeong-Cheol Kima, Hyung-Sup Junga, Mounng Jin Leecand Saro Lee, Spatial prediction of flood susceptibility using random-forest and boosted-tree models in Seoul metropolitan city, Korea, Year of Publication: 10 Apr 2017.
- [3] Chinmayee Kinage, Sejal Mandora, Abhishek Kalgutkar, Sunita Sahu , Amruta Parab, Performance Evaluation of Different Machine Learning Based Algorithms for Flood Prediction and Model for Real Time Flood Prediction, Year of Publication: 2019, 5th International Conference On Computing, Communication, Control And Automation (ICCUBEA).
- [4] Nikunj K. Mangukiya; Darshan J. Mehta; Raj Jariwala, Flood frequency analysis and inundation mapping for lower Narmada basin, India, Year of Publication: 1 February 2022.
- [5] Minister of State for Jal Shakti, Shri Bishweswar Tudu, Flood Forecasting and Early Warning System, Publication Date : 08 AUG 2022 6:02PM by PIB Delhi
- [6] Press Information Bureau, Government of India, Ministry of Water Resources, River Development and Ganga Rejuvenation, Flood Control Schemes by Indian Government, Year of Publication: January 10, 2023
- [7] Darshan Mehta , Jay Dhabuwala , Sanjaykumar M. Yadav , Vijendra Kumar , Hazi M. Azamathulla , Improving flood forecasting in Narmada river basin using hierarchical clustering and hydrological modeling, Year of Publication: December 2023
- [8] William Mobley, Antonia Sebastian, Russell Blessing, Wesley E. Highfield, Laura Stearns Samuel D. Brody, Institute for a Disaster Resilient Texas Texas A&M University at Galveston, Quantification of Continuous Flood Hazard using Random Forrest Classification and Flood Insurance Claims at Large Spatial Scales: A Pilot Study in Southeast Texas, October 28, 2020
- [9] Bhanu Kanwar, Development of flood prediction models using machine learning techniques, summer 2022
- [10] A Jaya Prakash, Indian Institute of Technology Kharagpur, Development of an Automated Method for Flood Inundation Monitoring, Flood Hazard and Soil Erosion Susceptibility Assessment Using Machine Learning and AHP-MCE Techniques, June 26th, 2023.
- [11] Jun Yan; Jiaming Jin; Furong Chen; Guo Yu; Hailong Yin; Wenjia Wang, Journal of Hydroinformatics (2018) 20 (1): 221–231, Urban flash flood forecast using support vector machine and numerical simulation
- [12] September 2020 International Research Journal of Multidisciplinary Technovation 2(5):8-14 DOI:10.34256/irjmt2052 License CC BY 4.0, Tabassum Ullah, CHRIST (Deemed to be University) Gnana Prakasi O.S.

- [13] Jaeyeong Lee ¹ andByunghyun Kim ², Scenario-Based Real-Time Flood Prediction with Logistic Regression
- [14] Water Level Prediction Model Applying a Long Short-Term Memory (LSTM)–Gated Recurrent Unit (GRU) Method for Flood Prediction
- [15] Using Multi-Factor Analysis to Predict Urban Flood Depth Based on Naive Bayes

Annexure

Industry/Inhouse: **Project Evaluation Sheet 2023-24** **Class: D12A**

Title of the Project (Group no): 6. *Resilient rivers- empowering flood free futures.*

Group Members: *Aanya Rajee (50), Ishita Marathe (41), Anyan Rajee (51)*

	Engineering Concepts & Knowledge (5)	Interpretation of Problem & Analysis (5)	Design / Prototype (5)	Interpretation of Data & Dataset (3)	Modern Tool Usage (5)	Societal Benefit, Safety Consideration (2)	Environment Friendly (2)	Ethics (2)	Team work (2)	Presentation Skills (3)	Applied Engg & Mgmt principles (3)	Life - long learning (3)	Professional Skills (5)	Innovative Approach (5)	Total Marks (50)
Review of Project Stage I	04	04	04	03	04	02	02	02	02	03	03	03	04	05	45
Comments:	<i>Back propagation, Nodal zone, Notifications</i>														

Dr. G. Bhatia
Name & Signature Reviewer1

	Engineering Concepts & Knowledge (5)	Interpretation of Problem & Analysis (5)	Design / Prototype (5)	Interpretation of Data & Dataset (3)	Modern Tool Usage (5)	Societal Benefit, Safety Consideration (2)	Environment Friendly (2)	Ethics (2)	Team work (2)	Presentation Skills (3)	Applied Engg & Mgmt principles (3)	Life - long learning (3)	Professional Skills (5)	Innovative Approach (5)	Total Marks (50)
Review of Project Stage I	04	04	04	03	04	02	02	02	02	03	03	03	04	05	45
Comments:															

Date: 10th FEB, 2024

Priyanka
Name & Signature Reviewer2

Inhouse/ Industry /Innovation/Research: Inhouse **Project Evaluation Sheet 2023 - 24** **Class: D12 A/B/C**

Sustainable Goal: *Climate change*

Title of Project: *Resilient Rivers: Empowering Flood Free Futures*

Group Members: *Prasand Lahane (36), Ishita Marathe (41), Anya Rajee (50), Anyan Rajee (51)*

	Engineering Concepts & Knowledge (5)	Interpretation of Problem & Analysis (5)	Design / Prototype (5)	Interpretation of Data & Dataset (3)	Modern Tool Usage (5)	Societal Benefit, Safety Consideration (2)	Environment Friendly (2)	Ethics (2)	Team work (2)	Presentation Skills (2)	Applied Engg & Mgmt principles (3)	Life - long learning (3)	Professional Skills (3)	Innovative Approach (3)	Research Paper (5)	Total Marks (50)
	05	05	04	03	05	02	02	02	02	02	03	02	03	03	05	48
Comments:	<i>Dr. G. Bhatia</i>															

Dr. G. Bhatia
Name & Signature Reviewer1

	Engineering Concepts & Knowledge (5)	Interpretation of Problem & Analysis (5)	Design / Prototype (5)	Interpretation of Data & Dataset (3)	Modern Tool Usage (5)	Societal Benefit, Safety Consideration (2)	Environment Friendly (2)	Ethics (2)	Team work (2)	Presentation Skills (2)	Applied Engg & Mgmt principles (3)	Life - long learning (3)	Professional Skills (3)	Innovative Approach (3)	Research Paper (5)	Total Marks (50)
	05	04	04	03	05	02	02	02	02	02	03	02	03	03	05	47
Comments:																

Date: 9th March, 2024

Priyanka Kevin Nagda
Name & Signature Reviewer2