

# **CLIMATE EXTREMITIES PROJECTION AND ANALYSIS FOR TAMIL NADU REGION USING MACHINE LEARNING**

**A PROJECT REPORT**

*Submitted by*

***ABEYANKAR G (2020115002)***

***SRINIKETHAN S (2020115091)***

***VISHNU GSK (2020115103)***

*submitted to the Faculty of*

**INFORMATION AND COMMUNICATION ENGINEERING**

*in partial fulfillment for the award of the degree*

*of*

**BACHELOR OF TECHNOLOGY**

*in*

**INFORMATION TECHNOLOGY**



**DEPARTMENT OF INFORMATION SCIENCE AND TECHNOLOGY**

**COLLEGE OF ENGINEERING, GUINDY**

**ANNA UNIVERSITY**

**CHENNAI 600 025**

**MAY 2024**

**ANNA UNIVERSITY**  
**CHENNAI - 600 025**  
**BONA FIDE CERTIFICATE**

Certified that this project report titled “**CLIMATE EXTREMITIES PROJECTION AND ANALYSIS FOR TAMIL NADU REGION USING MACHINE LEARNING**” is the bona fide work of ABEYANKAR G(2020115002), SRINIKETHAN S (2020115091) and VISHNU GSK (2020115103) who carried out project work under my supervision. Certified further that to the best of my knowledge and belief, the work reported herein does not form part of any other thesis or dissertation on the basis of which a degree or an award was conferred on an earlier occasion on this or any other candidate.

**PLACE:**

**CHENNAI**

**DATE: 03/05/2024**

**Ms. G. MAHALAKSHMI**

**TEACHING FELLOW**

**PROJECT GUIDE**

**DEPARTMENT OF IST, CEG**

**ANNA UNIVERSITY**

**CHENNAI 600025**

**COUNTERSIGNED**

**Dr. S. SWAMYNATHAN**

**HEAD OF THE DEPARTMENT**

**DEPARTMENT OF INFORMATION SCIENCE AND TECHNOLOGY**

**COLLEGE OF ENGINEERING, GUINDY**

**ANNA UNIVERSITY**

**CHENNAI 600025**

## ABSTRACT

Aligned with the UN Sustainable Development Goal No.13, "Climate Action," this research initiative seeks to proactively address climate-related challenges. This research work aims to address and amend the projection of meteorological droughts and floods in Tamil Nadu, a region heavily reliant on agriculture for sustenance and food security, over the next three decades. Recognizing the significant impact of these natural calamities on the environment, society, and economy, this project endeavors to refine forecasting accuracy through advanced machine learning methodologies. While conventional indices like the Standardized Precipitation Index (SPI) and Standardized Precipitation-Evapotranspiration Index (SPEI) offer valuable insights, our approach integrates historical SPI and SPEI data with additional meteorological parameters, including relative humidity and solar radiation, for more comprehensive and precise predictions. Objectives include projecting drought and flood extremes for the next 30 years, examining the influence of relative humidity and solar radiation on drought indices, and comparing the efficacy of various machine learning models for extremity determination.

The research employs a technology stack encompassing Climate Data Operators (CDO) and RStudio for data preprocessing, the SPEI package in R for drought indices, and machine learning algorithms such as Support Vector Machines, Random Forest, XGBoost Classifier, and Co-Adaptive Neuro-Fuzzy Inference Systems. Correlation analysis is utilized to unravel variable interdependencies. The tech stack also integrates tools like JupyterLab and Scikit-learn for interactive development and intricate data analysis. This research work not only aids in climate-resilient decision-making but also aligns with broader aspirations for climate action, fostering a sustainable future for Tamil Nadu.

## ACKNOWLEDGEMENT

We wish to record our deep sense of gratitude and profound thanks to our project supervisor **Ms. G. Mahalakshmi**, Teaching Fellow, Department of Information Science and Technology, College of Engineering, Guindy for her keen interest, inspiring guidance, constant encouragement throughout this project.

We are extremely indebted to **Dr. S. Swamynathan**, Project Coordinator and Head of the Department of Information Science and Technology, Anna University, Chennai, for extending the facilities of the Department towards our project and for his unstinting support. We would also like to express our sincere thanks panel of reviewers **Dr. K. Vani**, Professor, **Dr. K. Indra Gandhi**, Associate Professor, **Ms. T. Sindhu**, Teaching Fellow, Department of Information Science and Technology for their valuable suggestions throughout the course of our project.

We thank our parents, family, and friends for bearing with us throughout the course of our project and for the opportunity they provided us in undergoing this course in such a prestigious institution.

**ABEYANKAR G**  
**SRINIKETHAN S**  
**VISHNU GSK**

# TABLE OF CONTENTS

	<b>ABSTRACT</b>	iii
	<b>LIST OF FIGURES</b>	vii
	<b>LIST OF SYMBOLS AND ABBREVIATIONS</b>	viii
<b>1</b>	<b>INTRODUCTION</b>	<b>1</b>
	1.1 INTRODUCTION	1
	1.2 OBJECTIVES	2
	1.3 PROBLEM STATEMENT	5
	1.4 OVERVIEW	5
	1.5 TECH STACK AND LIBRARIES USED	7
<b>2</b>	<b>LITERATURE SURVEY</b>	<b>9</b>
	2.1 LITERATURE REVIEW	9
	2.2 SUMMARY OF THE LITERATURE SURVEY	13
<b>3</b>	<b>SYSTEM DESIGN</b>	<b>14</b>
	3.1 SYSTEM ARCHITECTURE	14
	3.1.1 Data Collection	16
	3.1.2 Preprocessing	16
	3.1.3 Machine Learning Model Implementation	17
	3.1.4 Data Analysis and Correlation	17
	3.1.5 Time Series Analysis	17
	3.1.6 Forecasting of Cyclonic Occurrences	18
	3.1.7 System Components	19
<b>4</b>	<b>IMPLEMENTATION</b>	<b>20</b>
	4.1 DATA COLLECTION AND PREPROCESSING	20
	4.2 STANDARDIZED PRECIPITATION INDEX	21
	4.2.1 Features and Parameters	22
	4.2.2 Calculation in R Studio	22
	4.3 POTENTIAL EVAPOTRANSPIRATION	22
	4.3.1 Thornthwaite PET	23
	4.3.2 Hargreaves Method	23
	4.3.3 Spei	24
	4.4 MACHINE LEARNING MODELS USED	26

4.4.1	Support Vector Machine	26
4.4.2	Random Forest	28
4.4.3	XGBoost	30
4.4.4	K-Nearest Neighbors	31
4.4.5	Artificial Neural Networks	33
4.5	TIME SERIES ANALYSIS	33
4.5.1	SARIMAX Model	34
4.5.2	Prophet Model	35
4.5.3	LSTM Model	36
<b>5</b>	<b>RESULTS AND PERFORMANCE ANALYSIS</b>	<b>39</b>
5.1	PERFORMANCE EVALUATION	39
5.1.1	Accuracy	39
5.1.2	Precision	40
5.1.3	F1-Score	40
5.2	MODELS ANALYSIS	41
5.2.1	SVM	41
5.2.2	Random Forest	42
5.2.3	KNN	42
5.2.4	XGBoost	42
5.2.5	ANN	43
5.2.6	Model Accuracy Comparison	43
5.2.7	Results of the Time Series Analysis Models	45
<b>6</b>	<b>CONCLUSION AND FUTURE WORK</b>	<b>48</b>
6.1	CONCLUSION	48
6.2	FUTURE WORK	49
	<b>REFERENCES</b>	<b>50</b>

## LIST OF FIGURES

3.1	Proposed System Architecture	15
5.1	Confusion Matrix of SVM	41
5.2	Confusion Matrix of Random Forest	42
5.3	Confusion Matrix of KNN	42
5.4	Confusion Matrix of XGBoost	43
5.5	Accuracy of ANN	43
5.6	F1 Score, Precision and Recall Comparison	44
5.7	Model Accuracy Comparison	44
5.8	Lowest AIC Score for SARIMAX(1,0,0)(0,0,3)12	45
5.9	Comparison of Actual vs Forecasted SARIMAX Data	45
5.10	Cross Validation Results of Prophet Algorithm	46
5.11	Comparison of Actual vs Forecasted Prophet Data	46
5.12	Comparison of Actual vs Projected Data Validation Set (2015-2024)	47
5.13	Comparison of Actual vs Projected Data on Test Set (2024-2050)	47

## LIST OF ABBREVIATIONS

AI	Artificial Intelligence
ANN	Artificial Neural Network
ARIMA	Autoregressive Integrated Moving Average
CDO	Climate Data Operators
CNN	Convolutional Neural Network
KNN	K-Nearest Neighbors
LSTM	Long short-term memory
ML	Machine Learning
PET	Potential Evapo-Transpiration
SARIMAX	Seasonal Autoregressive Integrated Moving Average with Exogenous Factors
SPI	Standard Precipitation Index
SPEI	Standardized Precipitation-Evapotranspiration Index
SVM	Support Vector Machine



# **CHAPTER 1**

## **INTRODUCTION**

### **1.1 INTRODUCTION**

Enhancing the projection of meteorological droughts and floods in Tamil Nadu, where agriculture is pivotal for livelihoods and food security, is the primary focus. Recognizing the profound impacts of these natural disasters on the environment, society, and the economy, the project aims to improve the accuracy of forecasting tools using advanced machine learning models. While traditional indices like the Standardized Precipitation Index and Standardized Precipitation Evapotranspiration Index offer valuable insights, this project integrates historical SPI and SPEI data with additional meteorological variables, such as relative humidity and solar radiation, to provide more comprehensive and precise predictions. Motivated by the UN Sustainable Development Goal No.13, "Climate Action" the project addresses the need for proactive measures to mitigate climate-related challenges

Time Series Analysis refers to the analysis of the sequence of data points to find significant statistics and other characteristics. These data points, collected at constant time intervals, are dependent on the previous observations. Time Series Analysis is a crucial part of statistical studies and business analysis, forecasting stock market trends, economic forecasting, sales forecasting, etc. It consists of different models such as autoregressive (AR) models, moving average (MA) models, and autoregressive integrated moving average (ARIMA) models. Its main objective focuses on understanding the inherent structure and function of the time series to predict or forecast data points.

The objectives include projecting drought and flood extremities over the next 30 years, analysing the impact of relative humidity and solar radiation on drought indices, and comparing the accuracy of various ML models for extremity determination. The project utilizes a diverse technology stack, incorporating Climate Data Operators (CDO) and RStudio for data preprocessing, the SPEI package in R for drought indices, and machine learning models such as Support Vector Machines, Random Forest, XGBoost Classifier, and Co-Adaptive Neuro-Fuzzy Inference Systems. Correlation analysis is employed to understand variable interdependencies. The tech stack also includes tools like JupyterLab and Scikit-learn for interactive development and complex data analysis. By strategically leveraging cutting-edge technologies, the project not only contributes to climate-resilient decision-making but also aligns with the broader goal of climate action for a sustainable and resilient future in Tamil Nadu.

## 1.2 OBJECTIVES

### **Phase 1:**

**Data Acquisition and Preprocessing:** Gather historical meteorological and hydrological data from diverse sources such as IMD NetCDF files and TAMU global weather dataset.

Perform rigorous data cleaning, quality control, and feature engineering to ensure data integrity and suitability for analysis.

### **Extremities Classifier Model Development:**

Conduct comprehensive feature importance analysis to identify the most influential meteorological and hydrological factors contributing to

droughts and floods.

Employ machine learning algorithms (e.g., Random Forest, Support Vector Machines) to develop models capable of classifying past climate events (1985- 2014) into distinct severity classes based on the identified features.

Leverage deep learning techniques (e.g., Convolutional Neural Networks, Recurrent Neural Networks) to further refine the model architecture and improve classification accuracy.

### **Model Validation and Testing:**

Rigorously test and evaluate the developed Extremities Classifier Model in 5 key locations across Tamil Nadu, carefully assessing its performance and generalizability.

### **Phase 2:**

#### **Time Series Forecasting of Input Parameters:**

Implement an ensemble model (e.g., SARIMAX, RNN-LSTM, ARIMA, Exponential Smoothing) by strategically combining the strengths of various algorithms to generate robust and accurate forecasts of key meteorological and hydrological parameters for the next 30 years.

#### **Regional Expansion and Model Adaptation:**

Expand the model's coverage to encompass 25 districts in Tamil Nadu.

Develop and apply tailored ensemble models for each district,

accounting for regional weather patterns and unique characteristics.

### **Extremities Prediction and Validation:**

Utilize the forecasted meteorological and hydrological data as input to the Extremities Classifier Model to predict the likelihood of droughts and floods occurring within the designated timeframe.

Employ time series forecasting on the target variable (drought/flood occurrence) to validate the model's predictive accuracy, ensuring the reliability of the generated forecasts.

### **Short-termed prediction of cyclonic events:**

LSTM models have been employed to predict future values of critical data points. Leveraging the power of long short-term memory networks, accurate forecasts of important data trends have been achieved, facilitating proactive decision-making and strategic planning. Additionally, classifiers have been deployed to identify the presence of cyclonic events in forthcoming data.

By combining LSTM prediction with classifier detection, the project aims to provide a comprehensive tool for anticipating both the trajectory of key parameters and the occurrence of significant weather events, enabling timely and informed actions to mitigate potential risks.

### **1.3 PROBLEM STATEMENT**

Tamil Nadu, a region highly dependent on agriculture and susceptible to meteorological extremes, faces recurrent challenges of droughts and floods. Traditional indices like SPI and SPEI offer valuable insights into these climatic events, but their predictive accuracy can be enhanced. The problem at hand is to develop and implement an advanced machine learning approach that integrates historical SPI and SPEI data with additional meteorological parameters, including relative humidity, solar radiation, and temperature, to create more precise and robust projections of future SPEI values. Increased frequency and intensity of droughts and floods pose a significant threat to various sectors in Tamil Nadu, including water security, agriculture, and overall livelihood. Lack of long-term predictions for these events significantly hinders the effectiveness of preparedness and mitigation efforts, leaving communities vulnerable to their adverse impacts. This project seeks to improve the accuracy of drought and flood predictions in Tamil Nadu, addressing the critical need for timely and accurate information to guide climate-resilient decision-making and advance disaster preparedness efforts.

### **1.4 OVERVIEW**

This project tackles the critical challenge of drought and flood prediction in Tamil Nadu through the application of advanced machine learning techniques and a comprehensive technology stack. The project's foundation is built on input data spanning from 1985 to 2014, encompassing daily precipitation, maximum and minimum temperatures, sourced from the Indian Meteorological Department (IMD). The data preprocessing phase relies on the Climate Data Operators (CDO) in a Linux environment to compute monthly temperature averages and rainfall sums, generating a structured data frame with the help of Bash scripting and Python.

Evapotranspiration Index (SPEI) and the Standardized Precipitation Index (SPI), are computed in RStudio, incorporating Hargreaves and Thornthwaite potential evapo transpiration (PET) methods. This dataset is seamlessly integrated with data on solar radiation and relative humidity, procured from Texas A&M University's Soil and Water Assessment Tool (SWAT) repository. The project's core analysis relies on a diverse technology stack, employing Python-based machine learning models such as Support Vector Machines (SVM), Random Forest, XGBoost Classifier, and Co-Adaptive Neuro-Fuzzy Inference Systems (Co-ANFIS) for predictive modeling and classification tasks. Additionally, correlation analysis is conducted to unravel the inter dependencies among variables, enhancing our comprehension of their roles in drought and flood dynamics. In conclusion, this project not only empowers climate-resilient decision-making but also aligns with the UN Sustainable Development Goal No. 13, "Climate Action," by fostering preparedness and adaptation in Tamil Nadu through the strategic utilization of cutting-edge technologies.

Long Short-Term Memory (LSTM) is a type of recurrent neural network (RNN) architecture used in the field of deep learning. Unlike standard feedforward neural networks, LSTM has feedback connections which make it a "general purpose computer" – it can not only process single data points, but also entire sequences of data.

An LSTM network consists of various memory blocks known as cells, and there are three types of gates within each cell: input, forget, and output gate. Each gate can be imagined as a conventional artificial neuron, i.e., a logistic regression-like unit, enabling or blocking the information travelling to the next neuron.

SARIMAX, which stands for Seasonal AutoRegressive Integrated Moving Averages with eXogenous regressors, is a statistical model used for time series forecasting. This model captures the dependencies among the observations and the seasonal changes in the trends, while also accounting for the impact of additional variables. It's an extension of the SARIMA (Seasonal AutoRegressive Integrated Moving Averages) model, with the 'X' signifying the inclusion of exogenous variables. SARIMAX models can be used in a variety of fields such as economics and weather forecasting.

## 1.5 TECH STACK AND LIBRARIES USED

**CDO** : Climate Data Operators (CDO) is a powerful and versatile software tool designed for handling and analyzing climate and meteorological data. It provides a comprehensive set of command-line utilities for manipulating, transforming, and processing climate-related datasets in various file formats, including NetCDF. CDO is particularly valuable for researchers, scientists, and climatologists working with large and complex climate datasets, as it simplifies tasks like data extraction, calculation of statistical parameters, and re-gridding.

**RStudio**: RStudio is an integrated development environment (IDE) specifically designed for the R programming language. It provides a user-friendly and comprehensive platform for data analysis, statistical modeling, and data visualization. RStudio offers a range of features and tools that make R more accessible and efficient for researchers, data scientists, and statisticians.

**SPEI Package in R:** The Standardized Precipitation Evapotranspiration Index (SPEI) is a widely used drought index, and it is available as an R package. The SPEI package in R allows researchers to calculate and analyze SPEI values based on precipitation and potential evapotranspiration data.

**JupyterLab:** Jupyter Lab is a popular web-based interactive development environment (IDE) that comes integrated with the Anaconda distribution of Python. It is designed to enhance the Python coding experience and facilitate data analysis, scientific computing, and interactive documentation.

**Scikit-learn:** It is a famous Python library for working with complex data. Scikit-learn is an open-source library that supports machine learning. It supports variously supervised and unsupervised algorithms like linear regression, classification, clustering, etc.



## **CHAPTER 2**

### **LITERATURE SURVEY**

In this chapter, an exhaustive exploration is undertaken to furnish comprehensive insights into the application of research publications for the purpose of projecting climate extremes.

The discourse delves into the intricacies of employing scholarly works to inform and refine the projections pertaining to climatic extremities.

#### **2.1 LITERATURE REVIEW**

Mokhtarzad, M. et.al. [1] proposed a solution on drought forecasting and compares the effectiveness of artificial intelligence techniques, specifically Artificial Neural Network (ANN), Adaptive Neuro-Fuzzy Inference System (ANFIS), and Support Vector Machine (SVM). Using data from the Bonjouré meteorological station, the study tests these models for predicting the Standardized Precipitation Index (SPI) at 3-month time scales. Input parameters include temperature, humidity, and seasonal precipitation, with SPI as the output. The results indicate that SVM is the most accurate model for drought forecasting, as confirmed by nonparametric inference tests, outperforming ANN and ANFIS.

Soh, Y.W. et.al [2] presents an study on drought forecasting using the Wavelet-ARIMA-ANN (WAANN) and Wavelet-Adaptive Neuro-Fuzzy Inference System (WANFIS) models for different time scales at the Langat River Basin. It employs wavelet decomposition for data preprocessing and

evaluates model performance using various metrics. The research concludes that the WAANN model is more accurate for short-term and mid-term drought forecasting, while WANFIS performs well in mid-term predictions.

Danandeh Mehr et.al.[3]introduces a novel hybrid Random Forest (RF) model, GARF, which combines Genetic Algorithm (GA) with RF to enhance forecasting accuracy. GARF is applied to model and forecast a multitemporal drought index (SPEI-3 and SPEI-6) at two meteorology stations in Ankara, Turkey, and is compared to classic RF, standalone extreme learning machine (ELM), and Bat algorithm-optimized hybrid ELM (Bat-ELM). GARF outperforms benchmark models, achieving up to 30% and 40% improvement in forecasting accuracy for SPEI-3 and SPEI-6, respectively, particularly during the testing period. This study highlights the potential of GARF to significantly enhance classic RF techniques in drought prediction.

Nguyen et.al [4] proposed an solution on the use of the Adaptive Neuro-Fuzzy Inference System (ANFIS) for drought forecasting in Khanhhoa Province, Vietnam, employing the Standardized Precipitation Index (SPI) and Standardized Precipitation Evapotranspiration Index (SPEI). Input variables from sea surface temperature anomalies (SSTA) events in NinoW and Nino4 zones are used to forecast drought across different time scales. The study finds that ANFIS models with SSTA events as input variables outperform those using precipitation, providing high accuracy and reliability, particularly in long-term forecasting (SPEI-12). This research demonstrates the potential of ANFIS for effective drought prediction.

Dhangar et.al [5] gave an insights on the impact of drought on agriculture, ecosystems, and livelihoods. It introduces the Standardized Precipitation Evapo transpiration Index (SPEI), which considers both rainfall and evapo transpiration for drought monitoring. Utilizing data spanning 34

years (1980-2014) from the India Meteorological Department (IMD), gridded precipitation and temperature data are used to calculate SPEI for mapping drought in India. The SPEI is applied to normal and drought years (1985, 1987, 2002, 2004, 2009, and 2014), categorizing drought severity. The study finds that SPEI aligns with the Standardized Precipitation Index (SPI) and effectively identifies drought-affected regions. A strong correlation ( $r$ ) of 0.63-0.79 is observed between SPI and SPEI for various study years. Additionally, the comparison of SPEI with the Vegetation Health Index (VHI) underscores SPEI's utility for drought monitoring and assessment in India, positioning as drought index.

Praveen et.al[6] proposed solution to utilizes regional climate models, specifically the PRECIS model from the UK Met Office Hadley Centre, to provide detailed climate change forecasts for Thiruvallur, South India, corresponding to the IPCC-SRES A1B emission scenario for the period 2040-2070, with a reference to the base period of 1970-2000. The findings reveal a notable increase in mean maximum and minimum temperatures, along with a slight decrease in precipitation. Historical data analysis using the IMD method of Percent Deviation indicates moderate to mild drought occurrences, with 1974, 1980, 1982, and 1999 being prominent moderate drought years. The Standardized Precipitation Index (SPI-12) analysis identifies the extremely severe drought year of 1974, supported by Pearson's correlation analysis showcasing a significant positive correlation (at the 0.05 level) between the climate model outputs and historical drought data.

Mahalingam et.al [7] addresses the recurring challenge of drought in Tamil Nadu, an agrarian region, which significantly hampers societal and economic development. The analysis covers the period from 1981 to 2017, focusing on high-resolution (17 Km) precipitation data. Rainfall deviations from long-term averages during the Southwest Monsoon (SWM), Northeast

Monsoon (NEM), and annually are computed at the district level to identify drought-prone areas. The results highlight a higher frequency of drought events during the SWM compared to the NEM, with districts like Tirunelveli, Theni, Tiruppur, Karur, and Tiruchirapalli being particularly vulnerable during the SWM. Special attention and drought mitigation strategies are recommended for districts such as Kanyakumari, Salem, and Thiruvallur, which face deficit rainfall during the critical NEM crop-growing season.

Nitesh et.al [8] research delves into using LSTM networks for accurate weather forecasting. This study systematically evaluates LSTM's capability to model and predict meteorological time series data such as temperature, humidity, and atmospheric pressure. The literature notes the superiority of LSTM over traditional models due to its ability to learn long-term dependencies and handle nonlinear data patterns inherent in weather sequences. Key findings suggest improved forecast accuracy with LSTM, emphasizing its potential for operational meteorological applications.

In the study of LSTM and Climate data et.al [9], traditional statistical downscaling methods for climate data, which are often inadequate for precise local forecasts, are contrasted with a novel approach using a Super-Resolution Convolutional Long Short Term Memory Neural Network. This advanced model leverages Earth System Model outputs and enhances resolution from  $(1.25^\circ \times 0.9^\circ)$  to  $(0.25^\circ \times 0.25^\circ)$ , focusing on monthly precipitation in China. The Convolutional LSTM outperforms traditional methods in terms of mean squared error, relative bias, and correlation, indicating superior predictive accuracy and robustness for localized climate impacts.

## **2.2 SUMMARY OF THE LITRATURE SURVEY**

In conclusion, this literature survey underscores the significance of accurate drought prediction and monitoring. Advanced indices like SPEI and innovative models such as GARF and ANFIS enhance the precision of forecasting. Regional climate models, like PRECIS, reveal the effects of climate change, showing rising temperatures and shifting precipitation patterns. The strong correlation between SPI and SPEI demonstrates their reliability in drought assessment. Notably, SPEI's capacity for long-term drought prediction is an advantage. These insights emphasize the need for targeted drought mitigation strategies to protect agriculture and socioeconomic well-being.

## **CHAPTER 3**

### **SYSTEM DESIGN**

This chapter consists of technical architecture for Climate Extremities projection and analysis for Tamil Nadu region which is designed to effectively integrate various components and processes to enable accurate classification of extremities and study of influence of additional parameters on extremities.

#### **3.1 SYSTEM ARCHITECTURE**

The Figure 3.1 illustrates a system for short-term weather forecasting and long-term climate extremities projection. It begins with the acquisition and extraction of atmospheric and oceanic data, used to classify and predict cyclone intensity over the next 10 days, ensuring timely preparedness. Simultaneously, it incorporates variable features like temperature, rainfall, SPI, humidity, and solar radiation, along with constant features such as soil composition and land elevation. These inputs are used in a time series forecasting module employing techniques like SARIMAX, RNN-LSTM, ARIMA, and Exponential Smoothing to predict values over 30 years.

The long-term component projects climate extremities like droughts and floods using forecasted parameters processed by an extremities classifier module. This system provides a detailed outlook for the next 30 years in the Region of Interest (ROI). By integrating short-term and long-term forecasting, the system enhances predictive accuracy and reliability, essential for effective disaster preparedness and climate resilience planning.

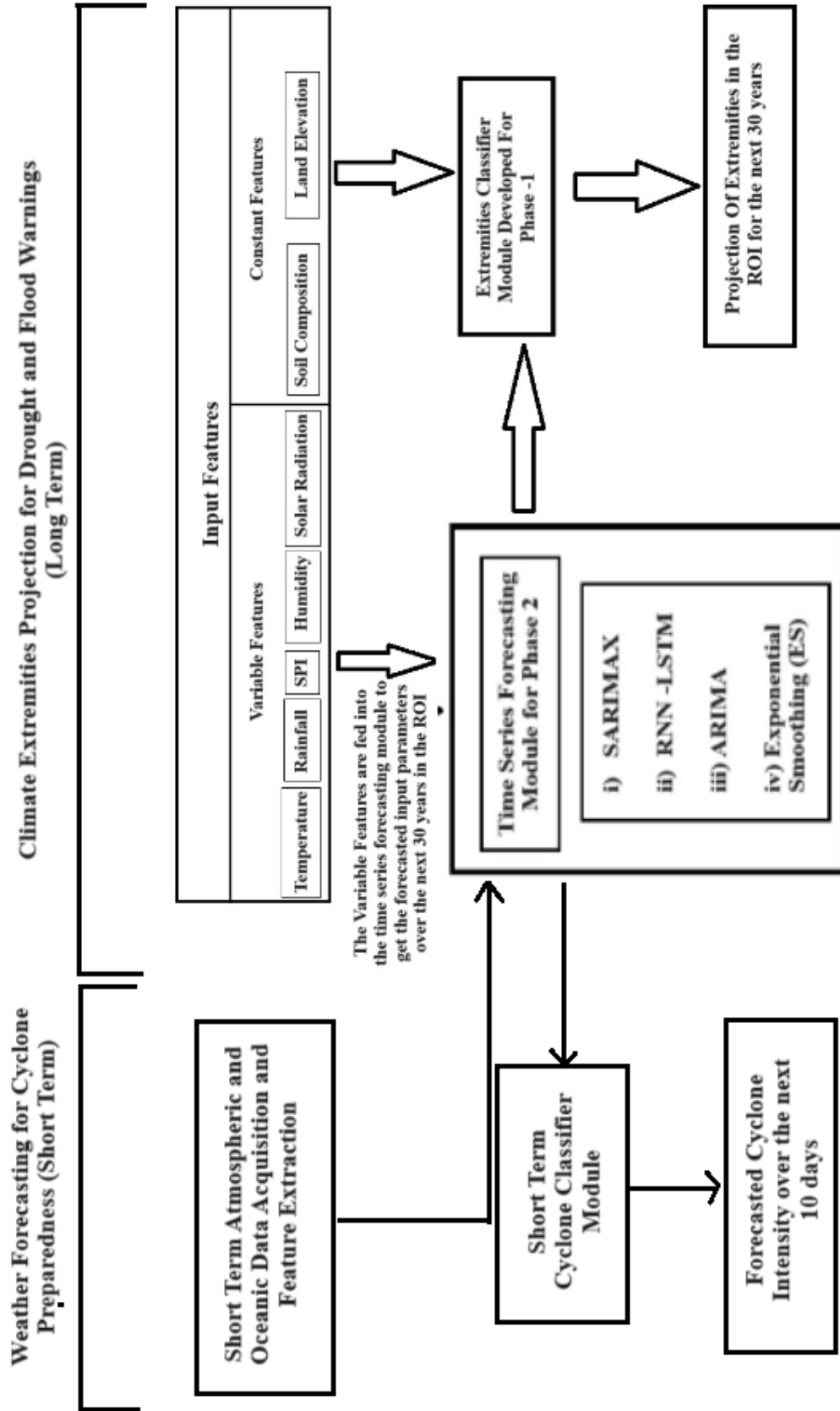


Figure 3.1: Proposed System Architecture

### **3.1.1 Data Collection**

This section describes the various methodologies employed in the collection and preprocessing of the data , also describes the types of the data used.

#### **Meteorological Data**

Historical meteorological data from the India Meteorological Department (IMD), including precipitation, maximum temperature, and minimum temperature for the Tamil Nadu region from 1985 to 2014 in netCDF format was gathered.

#### **Oceanic and Atmospheric Data**

Historical Oceanic and atmospheric from the Earth System Grid Foundation, which includes Wind Speed, Sea Level Pressure, Relative Humidity, Precipitation and Sea Surface Temperature.

#### **Additional Data**

Data pertaining to solar radiation and relative humidity was collected from the Texas A&M University Soil and Water Assessment Tool (SWAT) dataset repository.

### **3.1.2 Preprocessing**

Combination of the meteorological data with the solar radiation and relative humidity datasets based on their latitude and longitude coordinates.



Utilization of Climate Data Operators (CDO) in Linux and Python scripts to calculate monthly averages and sums of temperature and rainfall data.

Creation of a structured dataframe with year, month, precipitation, T-min, and T-max values using bash scripting, CDO, and Python.

Calculation of drought indices, including the Standardized Precipitation Evapotranspiration Index (SPEI) and Standardized Precipitation Index (SPI), using RStudio, incorporating Hargreaves and Thornthwaite PET methods.

### **3.1.3 Machine Learning Model Implementation**

Implementation of machine learning models in Python, including Support Vector Machines (SVM), Random Forest, XGBoost, and ANN, to project future SPEI values and classify data into extremities (drought or not).

### **3.1.4 Data Analysis and Correlation**

Study the correlation between meteorological variables (e.g., relative humidity, solar radiation) and drought indices (SPI and SPEI) to gain insights into the driving factors of drought and flood events.

### **3.1.5 Time Series Analysis**

Time series analysis involves analyzing data points over time to understand patterns and make forecasts. Several models such as ARIMA, SARIMAX, FB-Prophet and LSTMs were used to project the features with a higher accuracy. SARIMAX extends ARIMA to handle

seasonality and exogenous variables, and was evaluated using expanding window cross-validation and AIC. Prophet decomposes data into trend, seasonal, holiday, and error components for effective forecasting, assessed with cross-validation. LSTM models, a type of RNN, capture long-term dependencies in sequential data like time series. Meteorological data from 1950-2014 was fitted into LSTM with a window size of 24, evaluated using RMSE. LSTM outperformed SARIMAX, yielding an RMSE of 0.9415. In conclusion, SARIMAX, Prophet, and LSTM offer powerful techniques for time series analysis, each with its strengths and suitability for different data and forecasting needs, providing varied approaches to model temporal patterns, trends, and forecast future values.

### **3.1.6 Forecasting of Cyclonic Occurrences**

This module focuses on analyzing oceanic and atmospheric data collected from a specific latitude and longitude within the Bay of Bengal spanning from 2000 to 2004.

To project future data for key features influencing cyclone formation in the region. To achieve this, Long Short-Term Memory (LSTM) models are employed to generate future projections for each relevant parameter. These LSTM models are chosen for their effectiveness in capturing complex temporal dependencies within sequential data.

Subsequently, the projected future values are input into classifier models, particularly XGBoost, trained to predict cyclonic occurrences. By integrating LSTM-based forecasting with XGBoost classification, the aim is to determine the likelihood of cyclone formation based on projected atmospheric and oceanic conditions.

In short-term projections of cyclonic occurrence features, LSTM models exhibited superior performance over ARIMA models, with an RMSE value of 0.276772. This indicates the effectiveness of LSTM in accurately predicting cyclonic events within a short timeframe.

### **3.1.7 System Components**

Data Processing: Linux, CDO, Python, and RStudio for data collection, integration, and preprocessing.

Drought and Flood Index Calculation: RStudio for SPI and SPEI computation.

Machine Learning: Python for implementing machine learning models (SVM, Random Forest, XGBOOST, Co-ANFIS).

Data Analysis and Visualization: Python libraries (Matplotlib, Seaborn) for data analysis and visualization.

## **CHAPTER 4**

### **IMPLEMENTATION**

The implementation phase of this report represents the crucial juncture where theoretical frameworks and conceptual foundations are translated into actionable insights. This chapter is dedicated to elucidating the methodologies employed for data collection, with a keen focus on extracting meaningful information to gauge societal progress.

#### **4.1 DATA COLLECTION AND PREPROCESSING**

In this study, an extensive analysis of climatic variables over a 30-year period was conducted, focusing on specific latitudinal and longitudinal coordinates in India. The dataset encompasses crucial meteorological elements, including maximum and minimum temperatures, precipitation, relative humidity, and solar radiation. The data, sourced from IMD for temperature and precipitation and TAMU for relative humidity and solar radiation, has been meticulously structured into a comprehensive table. Each row represents a day, while columns capture the various climate parameters. Through statistical measures such as mean, median, and standard deviation, key insights into the central tendencies of these variables were derived over the three-decade timeframe. To enhance understanding, visualizations such as line charts and bar graphs were employed to illustrate trends and variations. The interpretation of these patterns is crucial in understanding the complex interplay of climate factors at the specified geographical coordinates in India. This climate dataset serves as a foundational component of the broader analysis pursued, complementing other variables such as SPI values, PET models, and

extremities projections. The integration of these diverse datasets contributes to a holistic understanding of the climatic dynamics in the region.

## 4.2 STANDARDIZED PRECIPITATION INDEX

The Standardized Precipitation Index (SPI) is a widely used drought index that standardizes precipitation values to facilitate the comparison of drought conditions across different time scales and locations. SPI is particularly valuable in assessing the severity and duration of meteorological droughts.

The SPI is calculated through a three-step process: Data Preparation, Parameter Estimation, and Index Computation.

**Data Preparation:** Accumulate precipitation data over a specified time scale i.e. 12 months here. Fit a probability distribution to the accumulated precipitation data.

**Parameter Estimation:** Estimate the parameters of the chosen probability distribution. Commonly used distributions include the gamma distribution for monthly data and the Pearson Type III distribution for longer time scales. Transform the precipitation data to the standard normal distribution using the estimated parameters.

**Index Computation:** The transformed data is the SPI. Negative SPI values indicate below-average precipitation, suggesting drought conditions, while positive values indicate above-average precipitation.

### **4.2.1 Features and Parameters**

Precipitation Data: Monthly or longer time-scale precipitation values.

Time Scale: The time period over which precipitation is accumulated (e.g., 1 month, 3 months).

Probability Distribution: The choice of the distribution depends on the characteristics of the precipitation data. Commonly used distributions include gamma and Pearson Type III.

### **4.2.2 Calculation in R Studio**

In R Studio, the calculation of SPI involves utilizing SPEI package that provide functions to estimate parameters and compute the index. The spi function typically takes precipitation data, a time scale, and a probability distribution as inputs.

## **4.3 POTENTIAL EVAPOTRANSPIRATION**

Potential Evapotranspiration (PET) represents the maximum possible rate of water evaporation and transpiration under optimal climatic conditions. The Thornthwaite method is a commonly used approach to estimate PET, considering temperature as a key factor.

#### 4.3.1 Thornthwaite PET

The equation (4.1) is used for the calculation of PET

$$PET = 1.6 * (10 * T/2)^a \quad (4.1)$$

where:

$T$  is the average monthly temperature in degrees Celsius.  $I$  is the heat index, calculated as the sum of monthly temperature anomalies raised to the 1.514 power.  $a$  is a coefficient that varies based on latitude.

#### Required Features And Parameters

- Temperature Data: Monthly average temperatures are required for the Thornthwaite PET calculation.
- Latitude: The geographical latitude of the location to determine the coefficient 'a'.
- Time Period: The calculations are typically done on a monthly basis.

#### 4.3.2 Hargreaves Method

The Hargreaves method is a widely used approach for estimating PET, relying on temperature data as a key determinant.

The equation (4.2) is used for the calculation of PET using Hargreaves method

$$PET = 0.0023 * (T_{max} - T_{min}) * (T_{avg} + 17.8) * \sqrt{T_{max} - T_{min}} \quad (4.2)$$

where :

$T_{max}$  is the daily maximum temperature,

$T_{min}$  is the daily minimum temperature,

$T_{avg}$  is the average temperature

### **Required Features and Parameters**

**Temperature Data:** Daily maximum and minimum temperatures are necessary for Hargreaves PET calculations. **Time Period:** The calculations are typically done on a daily basis. **Solar Radiation:** Although not explicitly present in the formula, solar radiation data can enhance PET estimates.

#### **• 4.3.3 Spei**

The Standardized Precipitation Evapotranspiration Index (SPEI) is a drought index that combines precipitation and potential evapotranspiration (PET) to assess drought conditions. Unlike the Standardized Precipitation Index (SPI), which focuses solely on precipitation, SPEI incorporates both precipitation and temperature data, providing a more comprehensive measure of drought.

### **Calculation Method**

#### **Calculate PET:**

Potential Evapotranspiration (PET) represents the maximum amount



of water that could evaporate under optimal climatic conditions. It is typically calculated using methods like Thornthwaite, Hargreaves,

### **Calculate the Water Balance**

The water balance is computed as the difference between precipitation (P) and PET.

### **Calculate Accumulated Water Balance**

The accumulated water balance is determined for specific time scales (e.g., 1 month, 3 months, 6 months).

### **Standardize the Accumulated Water Balance**

The standardized value is obtained by normalizing the accumulated water balance using the mean and standard deviation over a reference period. the equation (4.3) is used to calculate the SPEI

$$SPEI = (W - w) / \zeta_w \quad (4.3)$$

### **Features Required for Calculation**

**Precipitation Data (P):** Monthly or seasonal precipitation data is a fundamental input for SPEI calculations.

**Potential Evapotranspiration Data (PET):** Calculated using methods such as Thornthwaite, Hargreaves, or Penman-Monteith.

**Time Scale:** The time scale defines the duration over which the water balance is accumulated (e.g., 1 month, 3 months).

**Reference Period:** The reference period is used to calculate the mean and standard deviation for standardization. It's typically a long-term historical period.

**Interpretation:**

- Positive SPEI Values indicate wetter-than-average conditions.
- Negative SPEI Values indicate drier-than-average conditions and the severity of drought increases with more negative values.

## **4.4 MACHINE LEARNING MODELS USED**

This section covers the major machine learning models that are used on the proposed work

### **4.4.1 Support Vector Machine**

Support Vector Machine (SVM) is a supervised machine learning algorithm that is used for classification and regression tasks. Here are some key notes on SVM:

SVM aims to find a hyperplane in an N-dimensional space (where N is the number of features) that distinctly classifies data points into different classes.

The margin is the distance between the hyperplane and the nearest data point from either class. SVM seeks to maximize this margin.

Support vectors are the data points that lie closest to the decision boundary (hyperplane) and influence the positioning of the hyperplane.

SVM can handle non-linear decision boundaries through the use of a kernel function, which implicitly maps the input data into a higher-dimensional space.

### **Types of SVM**

1. Linear SVM: Assumes a linear decision boundary. It is effective when the data is linearly separable
2. Non-linear SVM: Utilizes kernel functions (e.g., polynomial, radial basis function) to handle non-linear decision boundaries.
3. SVM Parameters:

C (Cost Parameter): Controls the trade-off between having a smooth decision boundary and classifying training points correctly. A smaller C encourages a wider margin, but some points may be misclassified. Kernel: Determines the type of decision boundary. Common kernels include linear, polynomial, and radial basis function (RBF). Gamma (for RBF Kernel): Defines the influence of a single training point. Higher values result in a more localized decision boundary.

#### **4.4.2 Random Forest**

Random Forest is an ensemble learning method that operates by constructing a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees.

At the core of Random Forest are decision trees, which are simple models that make decisions based on the values of input features.

Random Forest builds multiple decision trees by resampling the training data with replacement (bootstrap samples). Each tree is trained on a different subset of the data.

For each split in the decision tree, a random subset of features is considered. This introduces diversity among the trees, making the model more robust and less prone to overfitting.

In classification tasks, the final prediction is made by a majority vote of the individual trees. In regression tasks, it's the average prediction of all trees.

#### **Advantages**

Random Forests generally provide high accuracy, making them suitable for various tasks.

They are less prone to overfitting, thanks to the ensemble of diverse trees.

Random Forests can provide insights into feature importance, helping in feature selection.

Can handle missing values in the dataset.

### **Limitations**

The ensemble nature of Random Forests makes them less interpretable compared to individual decision trees.

Training and predicting with a large number of trees can be computationally expensive.

### **Features Used For Training**

- Maximum Temperature
- Minimum Temperature
- Precipitation
- Relative Humidity
- Solar Radiation
- SPI
- Temperature Difference
- Month

**Target Variable is the Extremities**

### **4.4.3 XGBoost**

XGBoost (eXtreme Gradient Boosting) is an optimized gradient boosting framework designed for speed and performance. It belongs to the family of boosting algorithms.

The base learners in XGBoost are typically decision trees, which are combined to create a powerful predictive model.

XGBoost employs boosting, where each tree corrects the errors of the preceding one. This results in a strong predictive model.

XGBoost provides a mechanism to evaluate the importance of each feature in making predictions, aiding in feature selection.

#### **Advantages**

- XGBoost is known for its speed and performance, often outperforming other machine learning algorithms.
- Regularization techniques, like L1 and L2 regularization, help control overfitting and improve the generalization of the model.
- XGBoost is designed to be highly parallelizable, making it efficient for large datasets.

#### **Features Used For Training**

- Maximum Temperature
- Minimum Temperature

- Precipitation
- Relative Humidity
- Solar Radiation
- SPI
- Target Variable is **Extremities**

#### 4.4.4 K-Nearest Neighbors

K-Nearest Neighbors (KNN) is a simple and intuitive supervised learning algorithm used for both classification and regression tasks. It belongs to the family of instance-based or lazy learning algorithms.

KNN relies on a distance metric (usually Euclidean distance) to measure the similarity between data points. The "k" in KNN represents the number of nearest neighbors to consider.

For classification, the majority class among the k-nearest neighbors determines the class of the new data point. Commonly used k values are 3, 5, or other odd numbers to avoid ties.

In regression tasks, the average or weighted average of the values of the k-nearest neighbors is used to predict the target value for the new data point.

#### **Advantages**

KNN is straightforward and easy to understand, making it an excellent choice for quick and simple solutions.

KNN does not have a training phase; it memorizes the entire training dataset. This makes it easy to adapt to new data.

KNN can be applied to various types of data, including both numerical and categorical.

### **Limitation**

Predictions can be computationally expensive, especially with large datasets, as it requires calculating distances between the new point and all data points.

KNN is sensitive to the scale of features. It's essential to normalize or standardize features to ensure equal importance.

In high-dimensional spaces, the concept of proximity becomes less meaningful, and KNN may struggle to find meaningful neighbors.

### **Features Used For Training**

Maximum Temperature

Minimum Temperature

Precipitation

Relative Humidity

Solar Radiation

SPI



Target Variable is the **Extremities**

#### **4.4.5 Artificial Neural Networks**

Artificial Neural Networks (ANNs) draw inspiration from the structure and functioning of the human brain. They are composed of interconnected nodes, or artificial neurons, organized in layers.

ANNs consist of an input layer, one or more hidden layers, and an output layer. Each layer contains nodes that process information and transfer it to the next layer.

Neurons are fundamental units in a neural network. Each neuron receives input, processes it, and produces an output.

Weights and biases are parameters that adjust during the training process. They determine the strength of connections between neurons and help in capturing the complexity of relationships in the data.

Activation functions introduce non-linearity to the model, allowing neural networks to learn complex patterns. Common activation functions include ReLU (Rectified Linear Unit) and Sigmoid.

### **4.5 TIME SERIES ANALYSIS**

Time series analysis involves analyzing data points collected or recorded over time, aiming to understand patterns, trends, and make forecasts. Here's how each method contributes to this analysis:

#### 4.5.1 SARIMAX Model

SARIMAX (Seasonal AutoRegressive Integrated Moving Average with eXogenous factors) represents an advanced iteration of the ARIMA model, geared towards comprehensive time series analysis by accommodating seasonal patterns. This extension facilitates the incorporation of exogenous variables, thereby enhancing the model's capability to capture intricate temporal dynamics.

In the evaluation of SARIMAX models, a rigorous methodology is adopted, utilizing Expanding Window Cross Validation with 5 folds alongside the Akaike Information Criterion (AIC). The AIC serves as a pivotal metric for model comparison, striking a balance between predictive accuracy and the complexity of the model. A lower AIC value signifies a superior fit, indicative of an optimal trade-off between precision and parsimony.

The process of determining the optimal SARIMAX configuration entails an exhaustive search, systematically exploring various combinations of ARIMA and seasonal orders. This search spans a range of parameter values encompassing:

p: Autoregressive (AR) order, denoting the number of lag observations included in the model.

d: Differencing (I) order, indicating the number of times differencing is applied to attain stationarity.

q: Moving Average (MA) order, representing the number of lagged forecast errors included in the model.

P: Seasonal AR order, defining the autoregressive order for seasonal components.

D: Seasonal differencing order, specifying the number of seasonal differences applied.

Q: Seasonal MA order, delineating the moving average order for seasonal components.

By systematically iterating through various combinations of these parameters, the SARIMAX model attains optimization, ensuring the most appropriate configuration is selected to accurately capture the underlying patterns within the time series data. This meticulous approach underscores the reliability and robustness of the resultant model, rendering it apt for informed decision-making and forecasting in diverse domains.

#### **4.5.2 Prophet Model**

Prophet, a forecasting model developed by Facebook, is renowned for its adeptness in handling time series data characterized by irregular trends and intricate seasonal patterns. Conceptually akin to a generalized additive model, Prophet employs a sophisticated approach to time series analysis, particularly suited for datasets with nuanced temporal dynamics.

Central to the effectiveness of Prophet is its capability to decompose time series data into distinct components, namely trend, seasonal, holiday, and error components. This decomposition facilitates a granular understanding of the underlying patterns, thereby enabling the model to capture and forecast complex temporal variations with remarkable precision.

**Assessing Accuracy:** In the assessment of forecast accuracy, an expanding window cross-validation approach with 5 folds was used. This rigorous methodology entails training the model on progressively larger subsets of the data in each fold, followed by evaluation on subsequent portions. By averaging the performance across multiple iterations, this approach furnishes a robust estimate of the model's predictive capability on unseen data, thus enhancing its reliability and generalizability.

**Model Equation:** The Prophet equation encapsulates the essence of this decomposition, where the forecast  $y(t)$  is expressed as the sum of trend ( $g(t)$ ), seasonality ( $s(t)$ ), holiday effects ( $h(t)$ ), and an error term ( $e(t)$ ). Each component plays a pivotal role in shaping the overall forecast, with trend reflecting long-term changes, seasonality capturing periodic fluctuations, holiday effects accommodating special occasions, and the error term encompassing stochastic variations specific to individual circumstances.

#### 4.5.3 LSTM Model

Long Short-Term Memory networks (LSTMs), a subset of Recurrent Neural Networks (RNNs), have emerged as a powerful tool for modeling sequential data by effectively capturing both short and long-term dependencies. Unlike traditional RNNs, LSTMs are designed to mitigate the vanishing gradient problem, allowing them to retain important information over extended sequences. This characteristic makes LSTMs particularly well-suited for a wide range of applications, including time series analysis, natural language processing, and speech recognition.

In the context of time series forecasting, LSTMs have demonstrated remarkable performance, surpassing conventional methods such as SARIMAX

(Seasonal Autoregressive Integrated Moving Average with Exogenous Factors). By leveraging their inherent ability to learn complex patterns and relationships within sequential data, LSTM models offer enhanced accuracy and robustness in predicting future values.

The LSTM models developed in this study have showcased their efficacy in both short and long-term projections of critical features. In short-term forecasting, LSTMs excel at capturing immediate trends and fluctuations, providing timely insights into the dynamic behavior of the data.

### **Proposed Methodology:**

**Acquisition and Preprocessing of Meteorological Data:** Historical meteorological data spanning from 1950 to 2014, comprising 780 months of monthly observations, was gathered. Subsequently, rigorous preprocessing procedures were applied to ensure the data's compatibility with the LSTM model architecture.

**Window Size Selection:** To effectively capture the temporal dependencies and seasonality inherent in the dataset, a window size of 24 months was chosen. This window size allows for the consideration of the preceding 24 months' effects on predicting the subsequent data point (i.e., the 25th month).

**Evaluation Metric:** The performance of the LSTM models was assessed using the Root Mean Squared Error (RMSE) metric. RMSE serves as a robust measure to quantify the disparity between predicted and actual values, providing insight into the model's predictive accuracy. The efficacy of the LSTM model was further validated through comparisons between the projected

values and the actual observations within distinct time periods: a validation set spanning from 2015 to 2024 and a separate test set encompassing the years 2025 to 2050.

By leveraging the RMSE metric and rigorous validation procedures, the LSTM models' efficacy in forecasting future meteorological trends was systematically evaluated. This methodological approach lays a solid foundation for leveraging deep learning techniques in predictive modeling tasks, offering insights into future climatic patterns and facilitating informed decision-making.

## CHAPTER 5

### RESULTS AND PERFORMANCE ANALYSIS

The model was further tested with images from outside the Dataset, and was made to predict the disease from this given image of the leaf, the results of this have been attached below depicting the final implementation of the model on several images.

#### 5.1 PERFORMANCE EVALUATION

This section provides the information about the performance of the various models used and the overall comparison of various models used.

##### 5.1.1 Accuracy

Accuracy is one metric for evaluating classification models. Informally, accuracy is the fraction of predictions our model got right. The equation (5.1) depicts formally the definitions of accuracy.

$$Accuracy = \text{Number of correct predictions} / \text{Total number of predictions} \quad (5.1)$$

The equation (5.2) shows that in binary classification, accuracy can also be calculated in terms of positives and negatives as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5.2)$$

### 5.1.2 Precision

Precision is the ratio between the True Positives and all the Positives. The equation (5.3) shows the formula for precision:

$$Precision = TP / (TP + FP) \quad (5.3)$$

where

TP = True Positives,

TN = True Negatives,

FP = False Positives,

FN = False Negatives

### 5.1.3 F1-Score

F1-score is the Harmonic mean of the Precision and Recall. This is easier to work with since now, instead of balancing precision and recall, we can just aim for a good F1-score and that would be indicative of a good Precision and a good Recall value as well. The equation (5.4) depicts the formula for F1-Score:

$$F1 - Score = 2 * Precision * Recall / (Precision + Recall) \quad (5.4)$$



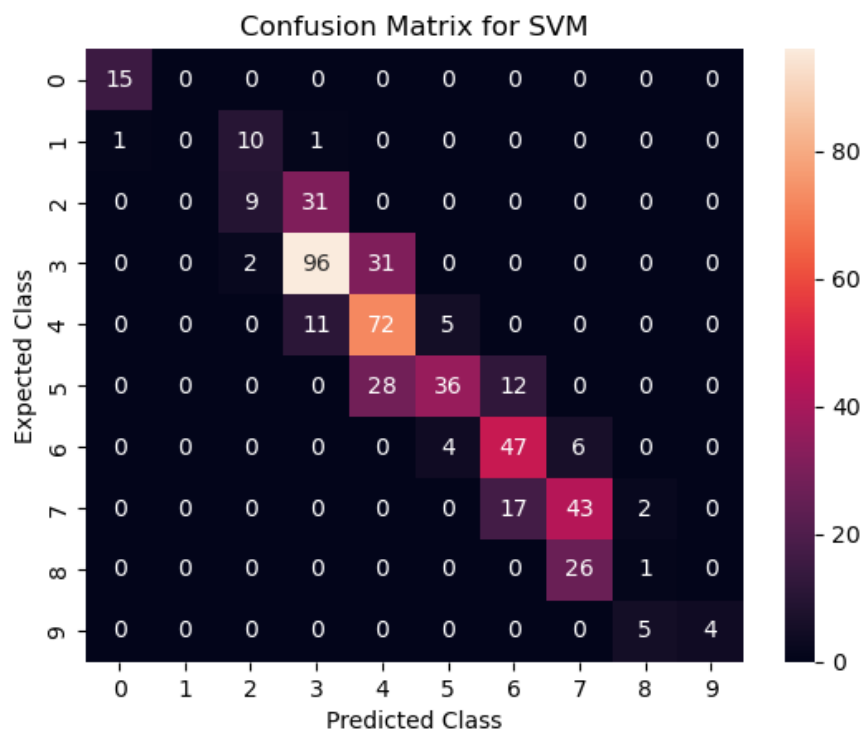
## 5.2 MODELS ANALYSIS

This section gives the information about the analysis of the various machine learning models.

### 5.2.1 SVM

#### Confusion Matrix

Figure 5.1 shows the Confusion Matrix of SVM .



**Figure 5.1: Confusion Matrix of SVM**

### 5.2.2 Random Forest

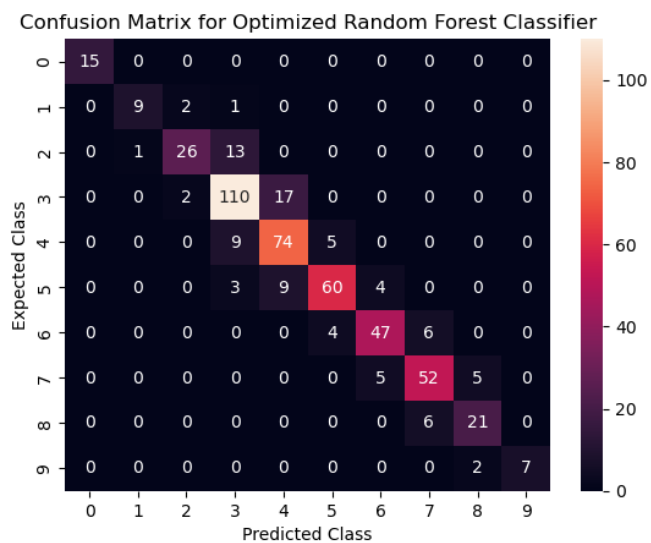


Figure 5.2: Confusion Matrix of Random Forest

### 5.2.3 KNN

#### Confusion Matrix

Figure 5.5 shows the Confusion Matrix of KNN

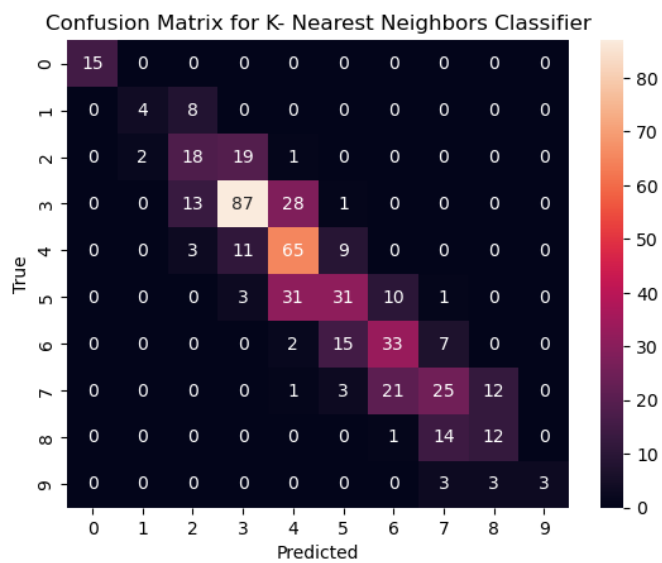


Figure 5.3: Confusion Matrix of KNN

## 5.2.4 XGBoost

### Confusion Matrix of XGBoost

Figure 5.4 shows the Confusion Matrix of XGBoost **Precision**

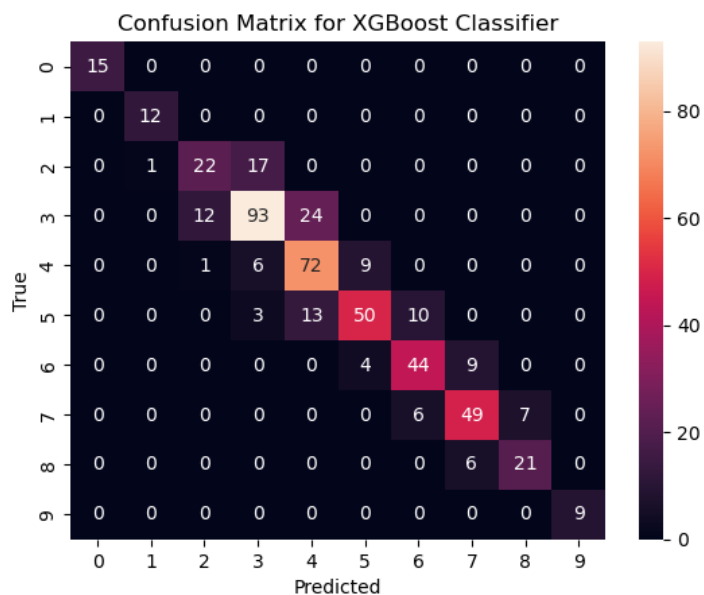


Figure 5.4: Confusion Matrix of XGBoost

## 5.2.5 ANN

```
Epoch 6/500
6/6 [=====] - 0s 4ms/step - loss: 0.2579 - accuracy: 0.8965
Epoch 7/500
6/6 [=====] - 0s 4ms/step - loss: 0.2497 - accuracy: 0.9023
Epoch 8/500
6/6 [=====] - 0s 4ms/step - loss: 0.2492 - accuracy: 0.8907
Epoch 9/500
6/6 [=====] - 0s 4ms/step - loss: 0.2447 - accuracy: 0.8907
Epoch 10/500
6/6 [=====] - 0s 3ms/step - loss: 0.2536 - accuracy: 0.8914
Epoch 11/500
6/6 [=====] - 0s 3ms/step - loss: 0.2331 - accuracy: 0.9089
Epoch 12/500
6/6 [=====] - 0s 3ms/step - loss: 0.2438 - accuracy: 0.9038
Epoch 13/500
...
Epoch 499/500
6/6 [=====] - 0s 3ms/step - loss: 0.2517 - accuracy: 0.8834
Epoch 500/500
6/6 [=====] - 0s 3ms/step - loss: 0.2643 - accuracy: 0.8936
```

Figure 5.5: Accuracy of ANN

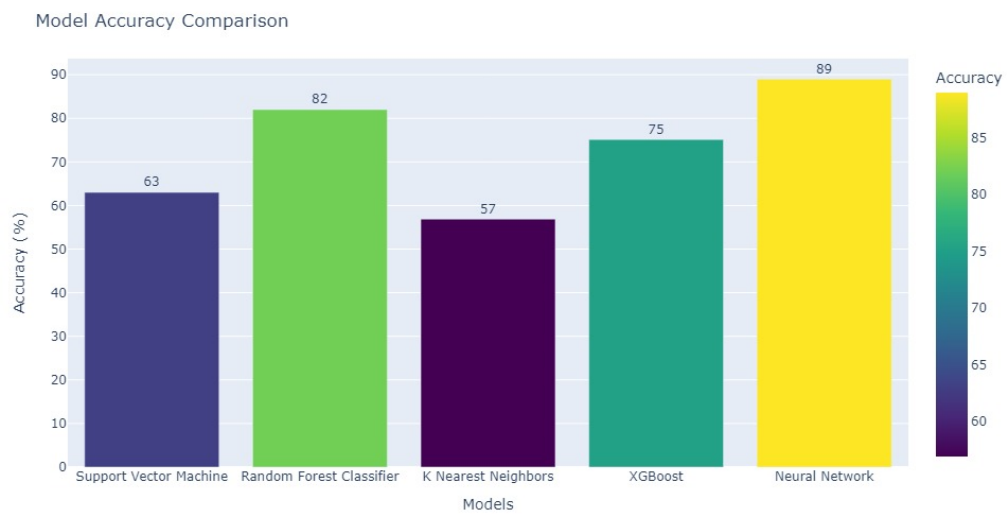
### 5.2.6 Model Accuracy Comparison

Figure 5.6 depicts the comparison of F1 score, precision and recall of the machine learning models used

	F1 Score	Precision	Recall
SVM	59.83	60.71	62.72
Random Forest	81.72	82.23	81.75
KNN	56.49	58.1	56.89
XGBoost	75.15	75.37	75.04

**Figure 5.6: F1 Score, Precision and Recall Comparison**

Figure 5.7 depicts the comparison of accuracy of the machine learning models used



**Figure 5.7: Model Accuracy Comparison**

## 5.2.7 Results of the Time Series Analysis Models

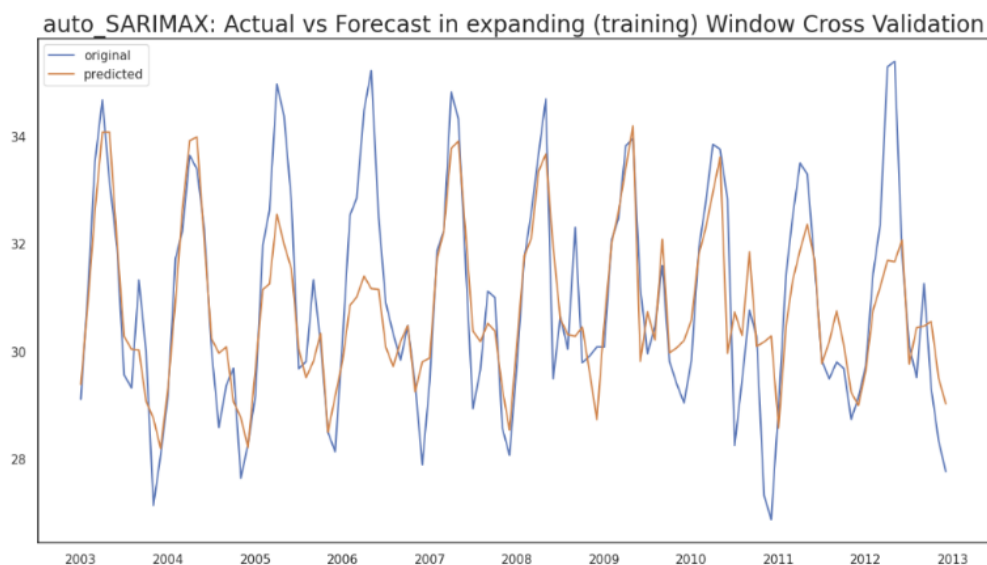
### SARIMAX Model

Figure 5.8 shows the Results of SARIMAX Algorithm on the Validation Dataset(2014-2024) and the best fit on the validation dataset was obtained for  $p=1, d=0, q=0, P=0, D=0, Q=3$

```
Best model: ARIMA(1,0,0)(0,0,3)[12] intercept
Total fit time: 87.934 seconds
Refitting data with previously found best parameters
Best aic metric = 913.2
```

SARIMAX Results			
Dep. Variable:	Max Temperature (C)	No. Observations:	336
Model:	SARIMAX(1, 0, 0)x(0, 0, [1, 2, 3], 12)	Log Likelihood	-449.583
Date:	Mon, 04 Mar 2024	AIC	913.166
Time:	15:22:27	BIC	939.069
Sample:	01-01-1985 - 12-01-2012	HQIC	923.534
Covariance Type:	opg		

**Figure 5.8: Lowest AIC Score for SARIMAX(1,0,0)(0,0,3)12**



**Figure 5.9: Comparison of Actual vs Forecasted SARIMAX Data**

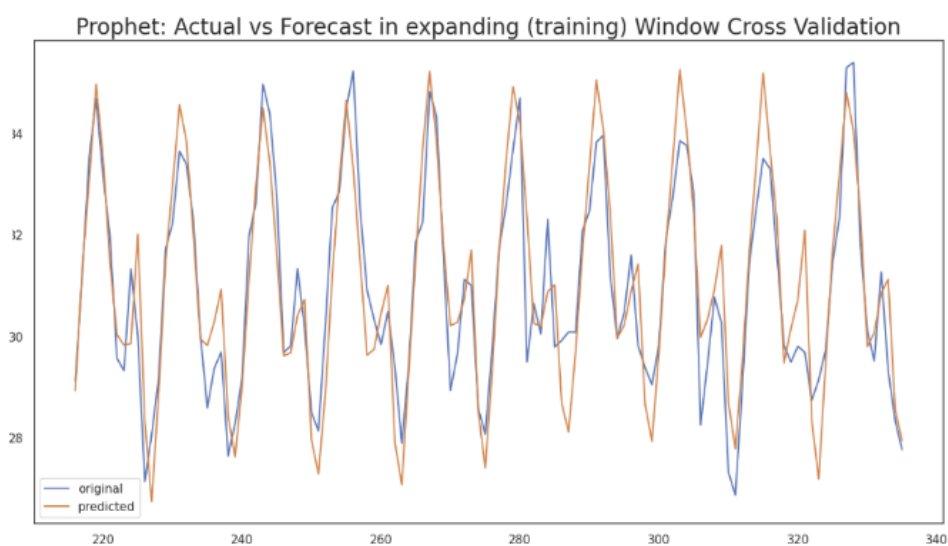
## Prophet Model

Figure 5.10 shows the Results of Prophet Algorithm on the Validation Dataset(2014-2024) and an RMSE of 0.9415 was obtained

### Model Cross Validation Results:

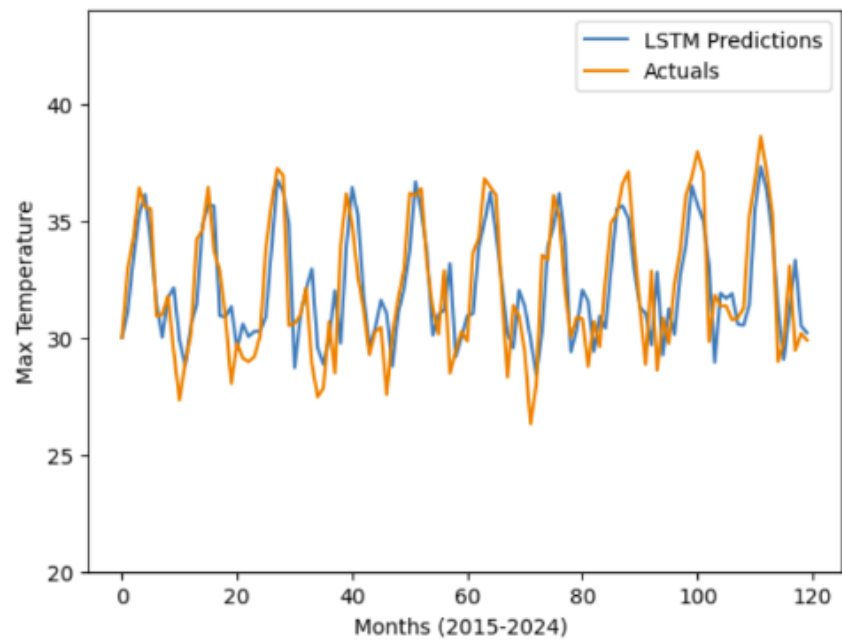
```
MAE (Mean Absolute Error) = 0.76
MSE (Mean Squared Error) = 0.89
MAPE (Mean Absolute Percent Error) = 2%
RMSE (Root Mean Squared Error) = 0.9415
Normalized RMSE (MinMax) = 11%
```

**Figure 5.10: Cross Validation Results of Prophet Algorithm**

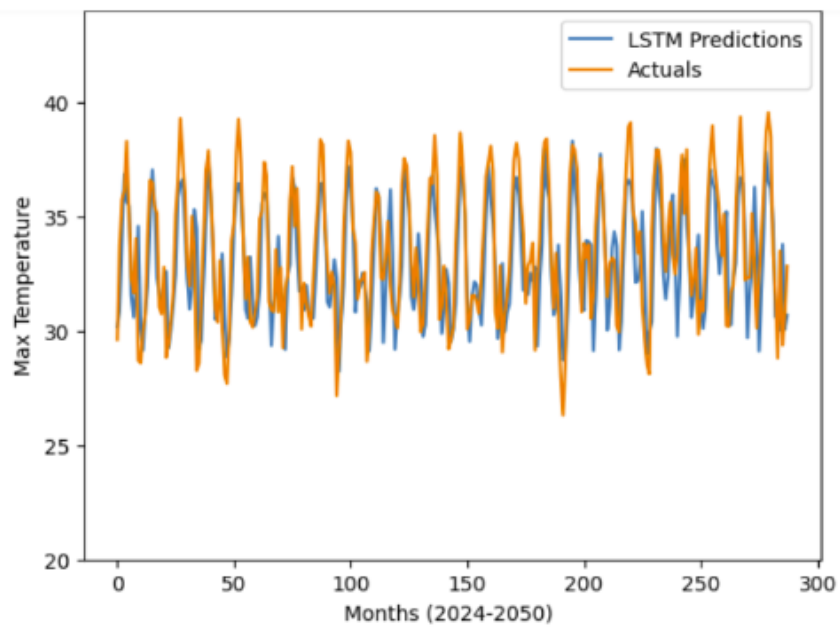


**Figure 5.11: Comparison of Actual vs Forecasted Prophet Data**

### LSTM Model



**Figure 5.12: Comparison of Actual vs Projected Data Validation Set (2015-2024)**



**Figure 5.13: Comparison of Actual vs Projected Data on Test Set (2024-2050)**

## CHAPTER 6

### CONCLUSION AND FUTURE WORK

#### 6.1 CONCLUSION

In conclusion, this project represents a pivotal step towards addressing the complex challenges posed by climate extremities in Tamil Nadu. Through meticulous data analysis, advanced modeling techniques, and comprehensive forecasting, we have laid a solid foundation for understanding and mitigating the adverse effects of droughts and floods.

Phase 1 provided crucial insights into the historical patterns of climate extremities, leveraging machine learning and deep learning models to classify and understand the severity of these events. The inclusion of hydrological features further enhanced the accuracy of extremity classification, setting the stage for more robust forecasting in Phase 2.

Phase 2 builds upon this foundation by expanding the geographical scope to encompass 25 districts in Tamil Nadu, allowing for a more nuanced understanding of regional variations in climate patterns and extremities. The implementation of ensemble forecasting models, including SARIMAX, RNN-LSTM, ARIMA, and Exponential Smoothing algorithms, will enable us to generate accurate projections of climate parameters for the next 30 years. By integrating forecasted data with extremity classification, we aim to validate the performance of our models and provide actionable insights for policymakers, stakeholders, and communities. Furthermore, our analysis of carbon dioxide emissions and their impact on climate extremities adds an important dimension to the project, highlighting the interconnectedness of climate change and



extreme weather events. By evaluating existing mitigative strategies, we seek to identify the most effective approaches for managing climate extremities and enhancing resilience in Tamil Nadu over the coming decades. In essence, this project underscores the importance of proactive measures in addressing climate change and its consequences. By leveraging cutting-edge technology, interdisciplinary collaboration, and a nuanced understanding of local contexts, we strive to empower decision-makers with the knowledge and tools needed to build a more resilient and sustainable future for Tamil Nadu.

## **6.2 FUTURE WORK**

The project aims to integrate advanced technologies and diverse data sources into meteorological research, notably through machine learning, enhancing prediction accuracy. Future advancements, especially in deep learning like neural networks, offer promise for refining forecasts by identifying intricate patterns. Real-time satellite imagery and remote sensing data provide crucial, up-to-the-minute information for disaster preparedness. Collaborations with meteorological institutions, government bodies, and tech firms could establish a nationwide monitoring network, improving early warning systems. Continued technological progress, including innovative sensors and IoT devices, will enhance the precision of forecasts. Beyond technology, the project's impact extends to societal and environmental domains, informing land-use planning, water resource management, and climate-resilient infrastructure. Scaling up and regional collaboration can establish a robust climate resilience framework, fostering collaboration among meteorologists, environmental scientists, policymakers, and data experts. Ultimately, it positions Tamil Nadu as a model for regions facing climate variability, driving meteorological advancements, influencing policy, and fostering sustainable development.

## REFERENCES

- [1] Mokhtarzad and N. Arabasadi A. M.Eskandari F., Jamshidi Vanjani. Drought forecasting by ann, anfis, and svm and comparison of the models. *Environmental Earth Sciences*, pages 1–10, 2017.
- [2] Y.W Soh, Koo, and Huang C.H., and K.F. Y.F., Fung. Application of artificial intelligence models for the prediction of standardized precipitation evapotranspiration index (spei) at langat river basin, malaysia. *Computers and Electronics in Agriculture*, pages 164–173, 2018.
- [3] Danandeh Mehr, Torabi Haghighi A., M. A., Jabarnejad, and Nourani V. Safari, M.J.S. New evolutionary hybrid random forest model for spei forecasting,. *Water*, page 755, 2022.
- [4] V. Nguyen and Nguyen L. Li, Q. Drought forecasting using anfis-a case study in drought-prone area of vietnam. *Paddy and Water Environment*, pages 605–616, 2017.
- [5] Dhangar, Narendra, Vyas, Swapnil, Guhathakurta, Pulak, Mukim, Shweta, Ramamurthy, Chattopadhyay Balasubramanian, and Nabansu. Drought monitoring over india using multi-scalar standardized precipitation evapotranspiration index. *Mausam*, pages 551–560, 2019.
- [6] Ramachandran A. Praveen, Dhanya. Projected warming and occurrence of meteorological droughts—insights from the coasts of south india. *American Journal of Climate Change*, 4, 173-179., pages 173–179, 2015.
- [7] Mahalingam Vengateswari, Geethalakshmi, Vellingiri, K. Bhuvaneswari, Panneerselvam Jagannathan, R., and Shanmugam. District level drought assessment over tamil nadu. *Madras Agricultural Journal*, pages 225–227, 2019.
- [8] Koduru Nitesh, Yanamala Abhiram, Rayapudi Krishna Teja, and S. Kavitha. Weather prediction using long short term memory (lstm) model,. *IEEE*, 2023.
- [9] Eric Chou. Christopher Chou, Junho Park. Generating high-resolution climate change projections using super-resolution convolutional lstm neural networks. *IEEE*, 2021.