

# **A MACHINE LEARNING MODEL FOR WEATHER FORECASTING**

## **A PROJECT REPORT**

**Submitted by**

**Ranbeer Singh Boparai (102303663)**

**Karanveer Singh (102303670)**

**Manjot Singh (102303672)**



**THAPAR INSTITUTE**  
OF ENGINEERING & TECHNOLOGY  
(Deemed to be University)

## **CERTIFICATE**

**Certified that this project report “A MACHINE LEARNING MODEL FOR WEATHER FORECASTING” is the bonafide work of “Ranbeer Singh Boparai, Karanveer Singh, and Manjot Singh” who carried out the project work under my supervision.**

**Dr. Ashima Khosla  
Lab Instructor**

## ACKNOWLEDGEMENT

We would like to express our sincere gratitude to our project guide and lab instructor, **Dr. Ashima Khosla**, for providing us with the invaluable opportunity to work on the project titled “*A Machine Learning Model for Weather Forecasting.*” Her guidance and support throughout the course of this work enabled us to explore the subject deeply and gain meaningful insights.

We also extend our heartfelt thanks to our parents and friends for their constant encouragement and assistance, which played a crucial role in the successful completion of this project within the given timeframe.

Date : 08-05-2025

Ranbeer Singh Boparai  
Karanveer Singh  
Manjot Singh

102303663

102303670

102303672

## TABLE OF CONTENTS

	<b>TITLE</b>	
	<b>ABSTRACT</b>	<b>5</b>
	<b>BACKGROUND</b>	<b>5</b>
	<b>OBJECTIVE</b>	<b>5</b>
<b>1.0</b>	<b>INTRODUCTION</b>	<b>6</b>
	1.1 Introduction	6
	1.2 Problem Statement	7
	1.3 Machine Learning	8
	1.4 Use of Algorithms	8
<b>2.0</b>	<b>METHODOLOGY</b>	<b>9</b>
<b>3.0</b>	<b>EXPERIMENTATION</b>	<b>11</b>
<b>4.0</b>	<b>RESULT AND DISCUSSION</b>	<b>12</b>
	5.1 Multiple Linear Regression	12
	5.2 Decision Tree Regression	13
	5.3 Random Forest Regression	14
<b>5.0</b>	<b>CONCLUSION</b>	<b>15</b>

## **ABSTRACT**

Traditionally, climate assessment has been performed reliably by treating the environment as a liquid. The current wind condition is being observed. The future state of the environment is recorded by understanding thermodynamics and the numerical position of the liquid elements. Nevertheless, this traditional arrangement of differential conditions as observed by physical models is at times unstable under oscillating effects and uncertainties when estimating the underlying states of air. This indicates an insufficient understanding of environmental variations, so it limits climate forecasts to 10-day periods because climate projections are essentially unreliable. But machine learning is moderately hearty for most barometric destabilizing effects compared to traditional techniques. Another favorable position of machine learning is that it does not depend on the physical laws of environmental processes.

## **Background**

For the current situation, India observatory conducts traditional weather forecasting. There are four common methods to predict the weather. The first method is the climatology method that is reviewing weather statistics gathered over multiple years and calculating the averages. The second method is an analog method that is to find a day in the past with weather similar to the current forecast. The third method is the persistence and trends method that has no skill to predict the weather because it relies on past trends. The fourth method is numerical weather prediction the is making weather predictions based on multiple conditions in the atmosphere such as temperatures, wind speed, high-and low-pressure systems, rainfall, snowfall, and other conditions. So, there are many limitations of these traditional methods. Not only it forecasts the temperature in the current month at most, but also it predicts without using machine learning algorithms. Therefore, our project is to increase the accuracy and predict the weather in the future for at least one month by applying machine learning techniques

## **Objective (Brief)**

Purpose of this project is to predict the temperature using different algorithms like linear regression, random forest regression, and Decision tree regression. The output value should be numerically based on multiple extra factors like maximum temperature, minimum temperature, cloud cover, humidity, and sun hours in a day, precipitation, pressure and wind speed.

## **1. INTRODUCTION**

Weather prediction is the task of predicting the atmosphere at a future time and a given area. This has been done through physical equations in the early days in which the atmosphere is considered fluid. The current state of the environment is inspected, and the future state is predicted by solving those equations numerically, but we cannot determine very accurate weather for more than 10 days and this can be improved with the help of science and technology.

Machine learning can be used to process immediate comparisons between historical weather forecasts and observations. With the use of machine learning, weather models can better account for prediction inaccuracies, such as overestimated rainfall, and produce more accurate predictions. Temperature prediction is of major importance in a large number of applications, including climate-related studies, energy, agricultural, medical, or etc.

There are numerous kinds of machine learning calculations, which are Linear Regression, Polynomial Regression, Random Forest Regression, Artificial Neural Network, and Recurrent Neural Network. These models are prepared dependent on the authentic information gave of any area. Contribution to these models is given, for example, if anticipating temperature, least temperature, mean air weight, greatest temperature, mean dampness, and order for 2 days. In light of this Minimum Temperature and Maximum Temperature of 7 days will be accomplished.

## PROBLEM STATEMENT

Weather forecasting plays a critical role in numerous sectors, including agriculture, transportation, disaster management, and energy planning. Accurate and reliable weather predictions can help mitigate the risks associated with adverse weather events, improve operational efficiency, and support informed decision-making. However, traditional forecasting methods—such as climatology, analog approaches, persistence and trends, and numerical weather prediction—suffer from several limitations. These techniques are either overly simplistic, fail to account for non-linear atmospheric dynamics, or are limited by their dependence on physical models and equations. As a result, their forecasting accuracy significantly diminishes beyond a short-term window of about 7 to 10 days. In a country like India, where weather patterns are highly variable and data-driven decision-making is increasingly essential, these limitations can pose serious challenges.

To address these issues, the use of machine learning (ML) offers a promising alternative. ML models can learn complex patterns and relationships from historical data, making them more resilient to the chaotic and non-linear nature of atmospheric processes. Unlike traditional models that rely heavily on physical laws and deterministic equations, ML algorithms are data-driven and can generate accurate predictions without explicitly modeling the underlying physics. Furthermore, ML techniques can efficiently handle large datasets and multiple input features, thereby enhancing the model's capability to generate more nuanced forecasts over longer periods.

This project aims to design and implement a machine learning-based model for temperature prediction, focusing on the city of Kanpur, India. The model leverages historical weather data collected from 2009 to 2020, encompassing parameters such as maximum and minimum temperatures, humidity, cloud cover, sun hours, wind speed, precipitation, and atmospheric pressure. These features serve as inputs for various regression algorithms including Multiple Linear Regression, Decision Tree Regression, and Random Forest Regression. The dataset, sourced from Kaggle and extracted using the [worldweatheronline.com](https://worldweatheronline.com) API, offers a robust foundation for training and evaluating these algorithms.

The primary goal of this project is to assess the performance of different ML models in predicting temperature with higher accuracy than conventional methods. The experimentation phase involves splitting the dataset into training and testing subsets, extracting relevant features, and tuning the models to optimize performance. Performance metrics such as Mean Absolute Error (MAE),  $R^2$  score, and confusion matrix analysis are used to evaluate the effectiveness of each algorithm. The outcome of this project is expected to not only identify the most accurate model among the chosen algorithms but also demonstrate the viability of using machine learning for long-term weather forecasting.

Ultimately, this project seeks to contribute to the ongoing shift from conventional meteorological modeling to more advanced, data-centric approaches. By integrating machine learning into weather prediction, it becomes possible to offer more precise and timely information to stakeholders across various sectors, thereby enabling proactive measures against weather-related uncertainties.

## **Machine Learning**

Machine learning is relatively robust to perturbations and does not need any other physical variables for prediction. Therefore, machine learning is a much better opportunity in the evolution of weather forecasting. Before the advancement of Technology, weather forecasting was a hard nut to crack. Weather forecasters relied upon satellites, data model's atmospheric conditions with less accuracy. Weather prediction and analysis have vastly increased in terms of accuracy and predictability with the use of the Internet of Things, for the last 40 years. With the advancement of Data Science, Artificial Intelligence, Scientists now do weather forecasting with high accuracy and predictability.

### **USE OF ALGORITHMS:**

There are different methods of foreseeing temperature utilizing Regression and a variety of Functional Regression, in which datasets are utilized to play out the counts and investigation. To Train, the calculations 80% size of information is utilized and 20% size of information is named as a Test set. For Example, if we need to anticipate the temperature of Kanpur, India utilizing these Machine Learning calculations, we will utilize 8 Years of information to prepare the calculations and 2 years of information as a Test dataset. The as opposed to Weather Forecasting utilizing Machine Learning Algorithms which depends essentially on reenactment dependent on Physics and Differential Equations, Artificial Intelligence is additionally utilized for foreseeing temperature: which incorporates models, for example, Linear regression, Decision tree regression, Random forest regression. To finish up, Machine Learning has enormously changed the worldview of Weather estimating with high precision and predictivity. What's more, in the following couple of years greater progression will be made utilizing these advances to precisely foresee the climate to avoid catastrophes like typhoons, Tornados, and Thunderstorms.



## 2. METHODOLOGY

The dataset utilized in this arrangement has been gathered from Kaggle which is “Historical Weather Data for Indian Cities” from which we have chosen the data for “Kanpur City”. The dataset was created by keeping in mind the necessity of such historical weather data in the community. The datasets for the top 8 Indian cities as per the population. The dataset was used with the help of the [worldweatheronline.com](https://worldweatheronline.com) API and the `wwo_hist` package. The datasets contain hourly weather data from 01-01-2009 to 01-01-2020. The data of each city is for more than 10 years. This data can be used to visualize the change in data due to global warming or can be used to predict the weather for upcoming days, weeks, months, seasons, etc.

Note: The data was extracted with the help of [worldweatheronline.com](https://worldweatheronline.com) API and we cannot guarantee the accuracy of the data.

The main target of this dataset can be used to predict the weather for the next day or week with huge amounts of data provided in the dataset. Furthermore, this data can also be used to make visualization which would help to understand the impact of global warming over the various aspects of the weather like precipitation, humidity, temperature, etc.

In this project, we are concentrating on the temperature prediction of Kanpur city with the help of various machine learning algorithms and various regressions. By applying various regressions on the historical weather dataset of Kanpur city we are predicting the temperature like first we are applying Multiple Linear regression, then Decision Tree regression, and after that, we are applying Random Forest Regression.

Table 2.1: Historical Weather Dataset of Kanpur City

	maxtempC	mintempC	totalSnow_cm	sunHour	uvIndex	uvIndex.1	moon_illumination	moonrise	moonset	sunrise	...	WindChillC	WindGustKmph	cloud
date_time														
2009-01-01 00:00:00	24	10	0.0	8.7	4	1	31	09:56 AM	09:45 PM	06:57 AM	...	11	21	
2009-01-01 01:00:00	24	10	0.0	8.7	4	1	31	09:56 AM	09:45 PM	06:57 AM	...	12	22	
2009-01-01 02:00:00	24	10	0.0	8.7	4	1	31	09:56 AM	09:45 PM	06:57 AM	...	12	23	
2009-01-01 03:00:00	24	10	0.0	8.7	4	1	31	09:56 AM	09:45 PM	06:57 AM	...	12	23	
2009-01-01 04:00:00	24	10	0.0	8.7	4	1	31	09:56 AM	09:45 PM	06:57 AM	...	14	19	

## Project Report: Machine Learning Model for Weather Forecasting

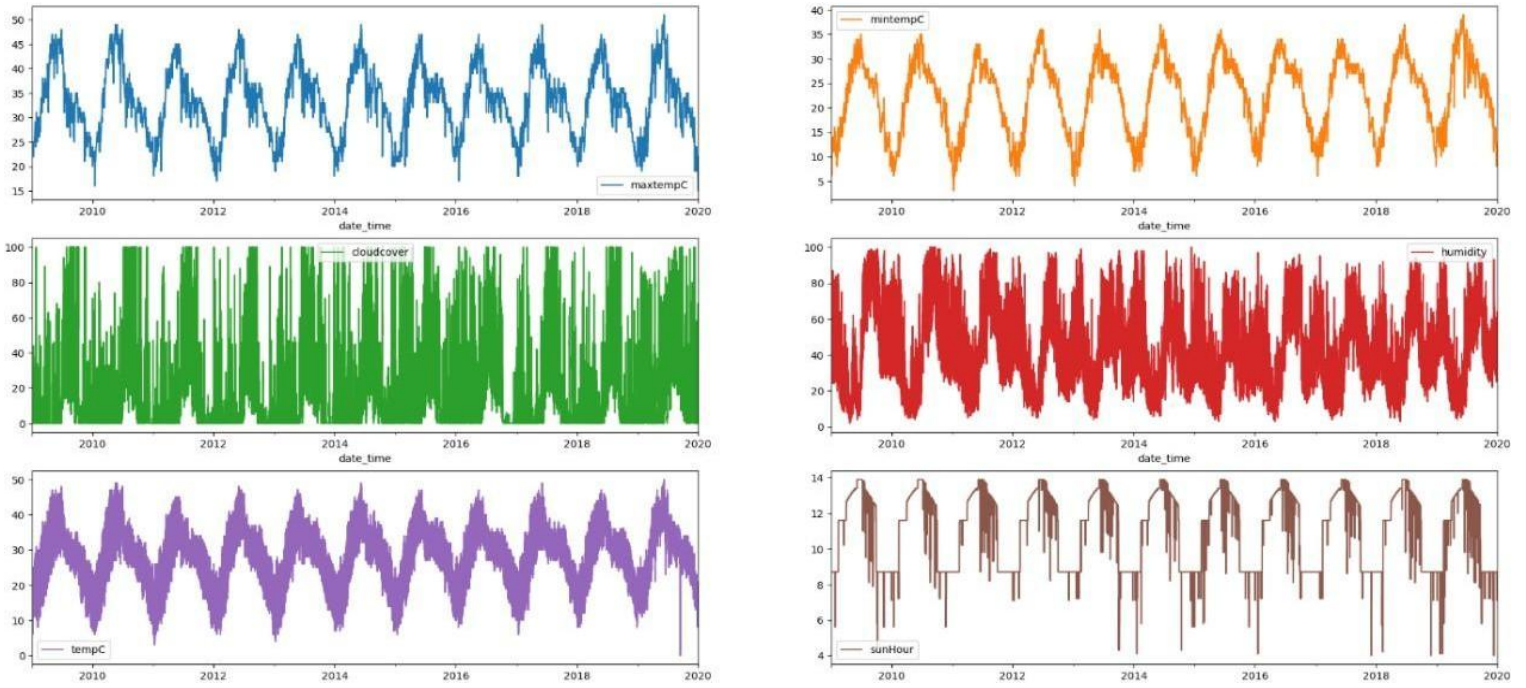


Figure 2.1: Plot for each factor for 10 years

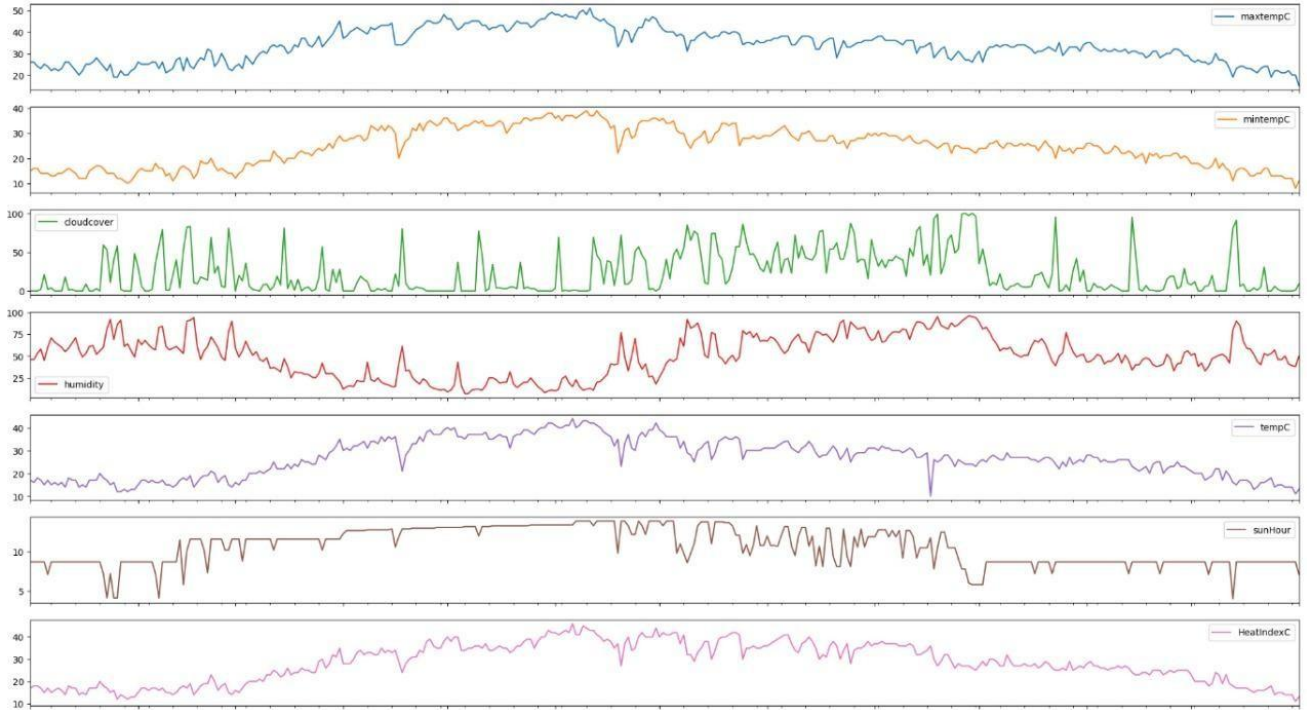
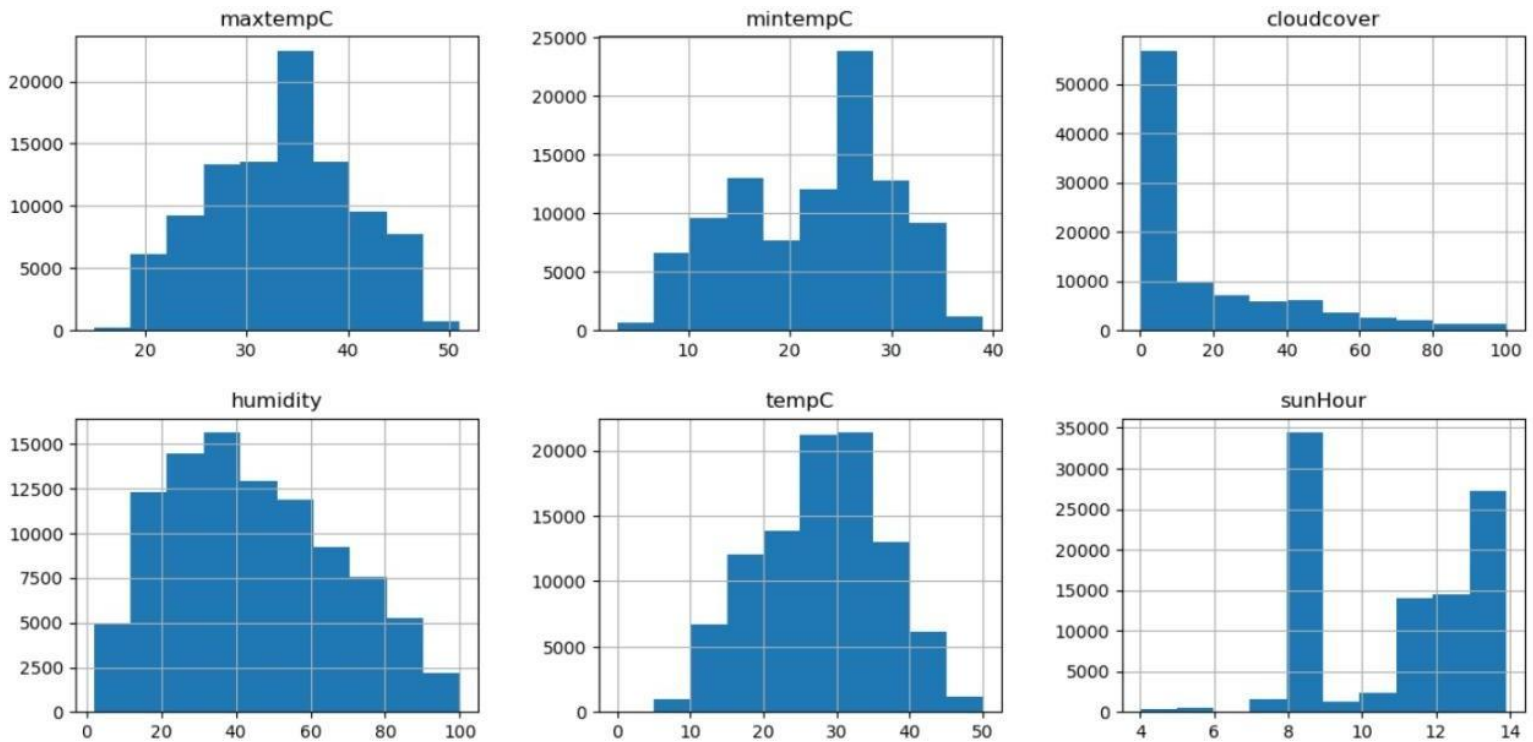


Figure 2.2: Plot for each factor for 1 year

### 3. EXPERIMENTATION

The record has just been separated into a train set and a test set. Each information has just been labeled. First, we take the trainset organizer. We will train our model with the help of histograms and plots. The feature so extracted is stored in a histogram. This process is done for every data in the train set. Now we will build the model of our classifiers. The classifiers which we will take into account are Linear Regression, Decision Tree Regression, and Random Forest Regression. With the help of our histogram, we will train our model. The most important thing in this process is to tune these parameters accordingly, such that we get the most accurate results. Once the training is complete, we will take the test set. Now for each data variable of the test set, we will extract the features using feature extraction techniques and then compare its values with the values present in the histogram formed by the train set. The output is then predicted for each test day. Now in order to calculate accuracy, we will compare the predicted value with the labeled value. The different metrics that we will use confusion matrix, R2 score, etc.



## 4. RESULTS AND DISCUSSION

The results of the implementation of the project are demonstrated below.

### Multiple Linear Regression:

This regression model has high mean absolute error, hence turned out to be the least accurate model. Given below is a snapshot of the actual result from the project implementation of multiple linear regression.

	Actual	Prediction	diff
date_time			
2013-07-10 08:00:00	34	34.89	-0.89
2015-11-04 20:00:00	25	24.57	0.43
2015-09-21 09:00:00	34	35.08	-1.08
2017-02-16 11:00:00	28	25.22	2.78
2012-07-21 01:00:00	28	28.04	-0.04
...	...	...	...
2019-03-30 09:00:00	37	33.55	3.45
2015-11-12 12:00:00	32	30.36	1.64
2019-12-31 05:00:00	8	9.13	-1.13
2019-08-02 17:00:00	35	35.92	-0.92
2019-10-22 08:00:00	26	25.77	0.23

19287 rows × 3 columns

**Decision Tree Regression:**

This regression model has medium mean absolute error, hence turned out to be the little accurate model. Given below is a snapshot of the actual result from the project implementation of multiple linear regression.

	Actual	Prediction	diff
date_time			
2013-07-10 08:00:00	34	34.0	0.0
2015-11-04 20:00:00	25	24.0	1.0
2015-09-21 09:00:00	34	34.0	0.0
2017-02-16 11:00:00	28	27.0	1.0
2012-07-21 01:00:00	28	28.0	0.0
...	...	...	...
2019-03-30 09:00:00	37	32.0	5.0
2015-11-12 12:00:00	32	32.0	0.0
2019-12-31 05:00:00	8	9.0	-1.0
2019-08-02 17:00:00	35	35.0	0.0
2019-10-22 08:00:00	26	26.0	0.0

19287 rows × 3 columns



**Random Forest Regression:**

This regression model has low mean absolute error, hence turned out to be the more accurate model. Given below is a snapshot of the actual result from the project implementation of multiple linear regression.

	Actual	Prediction	diff
date_time			
2013-07-10 08:00:00	34	33.92	0.08
2015-11-04 20:00:00	25	24.84	0.16
2015-09-21 09:00:00	34	34.25	-0.25
2017-02-16 11:00:00	28	27.00	1.00
2012-07-21 01:00:00	28	27.99	0.01
...	...	...	...
2019-03-30 09:00:00	37	32.79	4.21
2015-11-12 12:00:00	32	31.91	0.09
2019-12-31 05:00:00	8	8.81	-0.81
2019-08-02 17:00:00	35	34.98	0.02
2019-10-22 08:00:00	26	26.32	-0.32

19287 rows × 3 columns

## 5. CONCLUSION

All the machine learning models: linear regression, various linear regression, decision tree regression, random forest regression were beaten by expert climate determining apparatuses, even though the error in their execution reduced significantly for later days, demonstrating that over longer timeframes, our models may beat genius professional ones.

Linear regression demonstrated to be a low predisposition, high fluctuation model though polynomial regression demonstrated to be a high predisposition, low difference model. Linear regression is naturally a high difference model as it is unsteady to outliers, so one approach to improve the linear regression model is by gathering more information. Practical regression, however, was high predisposition, demonstrating that the decision of the model was poor and that its predictions can't be improved by the further accumulation of information. This predisposition could be expected to the structure decision to estimate temperature dependent on the climate of the previous two days, which might be too short to even think about capturing slants in a climate that practical regression requires. On the off chance that the figure was rather founded on the climate of the past four or five days, the predisposition of the practical regression model could probably be decreased. In any case, this would require significantly more calculation time alongside retraining of the weight vector  $w$ , so this will be conceded to future work.

Talking about Random Forest Regression, it proves to be the most accurate regression model. Likely so, it is the most popular regression model used, since it is highly accurate and versatile. Below is a snapshot of the implementation of Random Forest in the project.

Weather Forecasting has a major test of foreseeing the precise outcomes which are utilized in numerous ongoing frameworks like power offices, air terminals, the travel industry focuses, and so forth. The trouble of this determining is the mind-boggling nature of parameters. Every parameter has an alternate arrangement of scopes of qualities.

# References

## Data collection from:

### Kaggle:

Historical weather data for Indian cities [Data set]. Kaggle. Retrieved from <https://www.kaggle.com/datasets>

### World Weather Online API:

<https://www.worldweatheronline.com/>

## Algorithm Knowledge from:

### Random forest Regressor:

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

### Linear Regression:

<https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-linear-regression/>

### Descion Tree:

<https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-linear-regression/>