

PREDICTIVE ANALYSIS ON RAINFALL STATISTICS IN MAHARASHTRA

Sarthak Chudgar

Department of Instrumentation and
Control Engineering
Vishwakarma Institute of Technology
Pune, Maharashtra, India
sarthak.chudgar17@vit.edu

Johnathan Fernandes

Department of Instrumentation and
Control Engineering
Vishwakarma Institute of Technology
Pune, Maharashtra, India
johnathan.fernandes17@vit.edu

Aneesh Poduval

Department of Instrumentation and
Control Engineering
Vishwakarma Institute of Technology
Pune, Maharashtra, India
aneesh.poduval17@vit.edu

Abstract - Rainfall plays an integral role in the lives of millions of people worldwide. This is especially true for an agriculture heavy country such as India. Not only do we depend on rain as a source of fresh water, we also use it to irrigate farms and for rainwater harvesting.

Given our dependence on rain, prior knowledge about the weather forecast is essential for optimal rainwater utilization. This is where predictive analytics comes in.

Keywords: Rainfall, Prediction, Maharashtra, Weather, Rain.

I. INTRODUCTION

Predictive analytics is the 3rd step in big data analytics. It involves the previous two steps, i.e. descriptive and diagnostic analytics, and aims to build a model which can analyze data trends and predict future trends in order to help individuals better prepare themselves.

In this project we will carry out predictive analytics on weather data over the state of Maharashtra to construct a predictive model which will be able to predict the amount of rainfall in millimeters.

While the primary focus of this project is in the field of agriculture, it also plays a vital role in other fields such as disaster management

II. ANALYTICS

A. Descriptive Analytics

This step involves obtaining data of our variable (rainfall, in this case) and possible related factors (weather properties that may or may not affect rainfall).

We then visualize our dependent variable (rainfall) with each of the independent variables (other factors).

We obtain our data from the National Centers for Environmental Prediction (NCEP) website. They have constructed a Climate Forecast System Reanalysis (CFSR) which provides worldwide weather readings from as far back as 1980



Figure 1: NCEP Map GUI

The NCEP website provides a map GUI to choose an area, and then provides all data from that area. While the selection includes 657 stations, Maharashtra only contains about 370 of them. We use Tableau to extract those stations. From the system, we obtain data on the following:

Weather Factor	Units & Remarks
Precipitation	Millimeters
Humidity	Fraction
Location	Latitude and Longitude
Temperature	°C, (Minimum and Maximum)
Wind speed	Meters per second
Solar Coverage	Mega joules per square meter

Table 1: Selected Features

B. Diagnostic Analytics

Diagnostic analytics is closely related to descriptive to the point where they are usually carried out simultaneously. It entails using the previously created visualizations to determine the relations between our independent and dependent variable.

Average Precipitation And Average Relative Humidity By Month

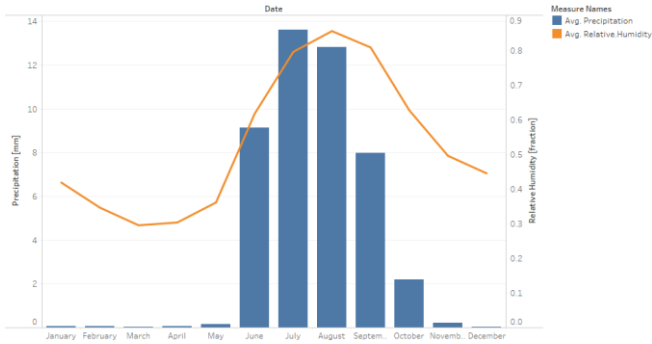


Figure 2: Average Precipitation and Average Humidity by Month

Average Precipitation And Average Maximum Temperature By Month

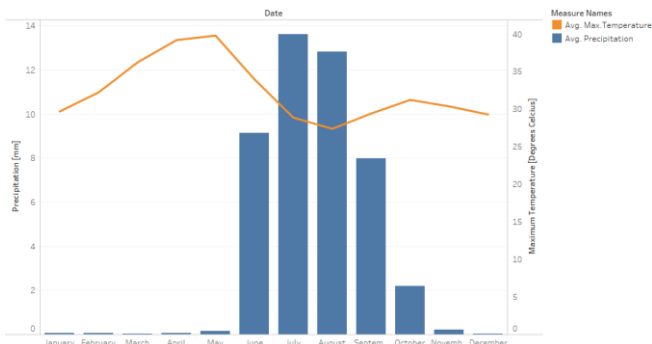


Figure 3: Average Precipitation and Average Maximum Temperature by Month

Average Precipitation And Average Minimum Temperature By Month

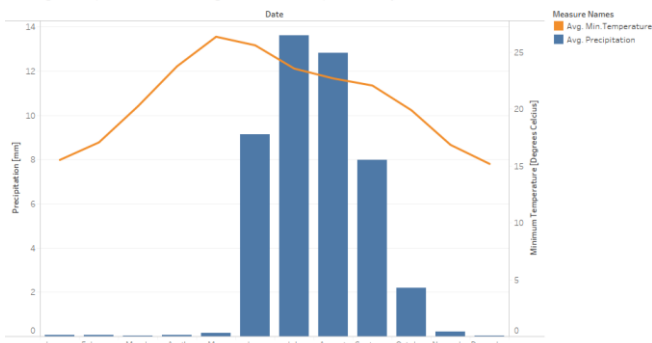


Figure 4: Average Precipitation and Average Minimum Temperature by Month

Average Precipitation And Average Solar Coverage By Month

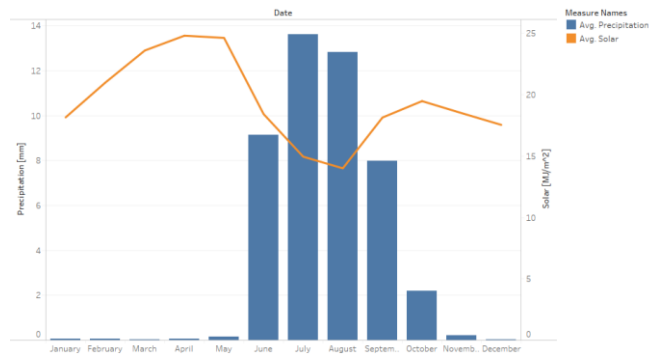


Figure 5: Average Precipitation and Average Solar Coverage by Month

Average Precipitation And Average Wind Speed By Month

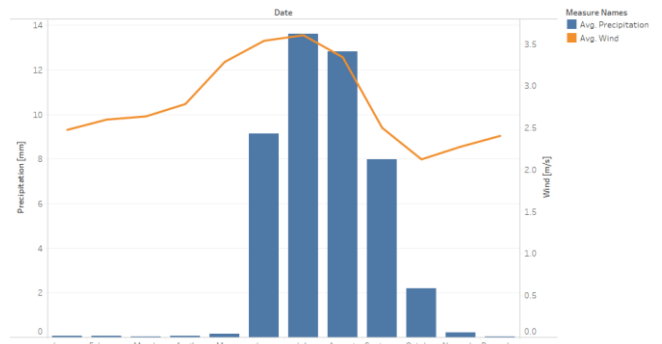


Figure 6: Average Precipitation and Average Wind Speed by Month

We use Tableau, a well-known business insights tool to visualize the data into various easy to understand charts.

C. Predictive Analytics

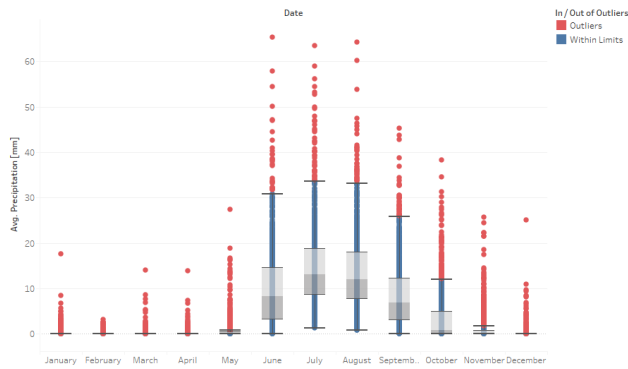
Our third and final step involved importing data into a suitable program to construct a prediction model. Due to the sheer volume of the data considered, normal analysis applications such as MATLAB and Microsoft Excel are not suitable. Hence we use Apache Spark along with the R programming language to import the data. The initial step in building a predictive model is to clean up the data. This involves removing undefined null values and outliers.

```
> colSums(is.na(MH))
      X      Date      Latitude      Longitude
      0       0          0          0
Max. Temperature Min. Temperature wind Relative Humidity
      0       0          0          0
      Solar      Precipitation
      0       0
```

Figure 7: Analyzing null values

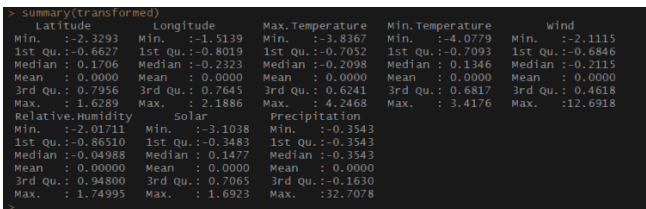
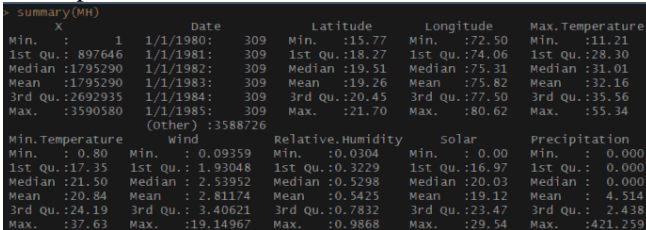
As observed, the dataset contains no null values.

Box And Whisker Plot To Visualize Outliers



Using tableau, we visualize the outliers in a box and whisker plot and remove them. In our plot, (Figure 8) each dot represents a single day.

After cleaning up null values and outliers, we scale and center the data points so as to maintain zero mean and unit variance. This reduces model computing time and improves model prediction.

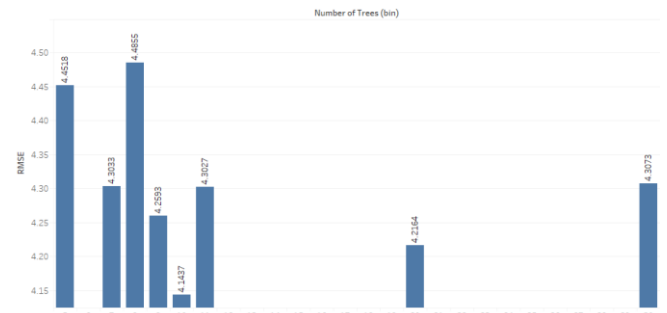


After cleaning up and scaling our data, we proceed to construct a random forest regression model in Spark using data from 1980 to 2010.

III. Prediction Model AND WORKING

After cleaning up and scaling our data, we proceed to construct multiple prediction models in Spark using data from 1980 to 2010. After thorough testing, we determine that a random forest regression model has the best performance. We further test multiple models using differing numbers of trees and compare their root mean square error (in fig. 11) and determine that a model with 10 trees has the least error.

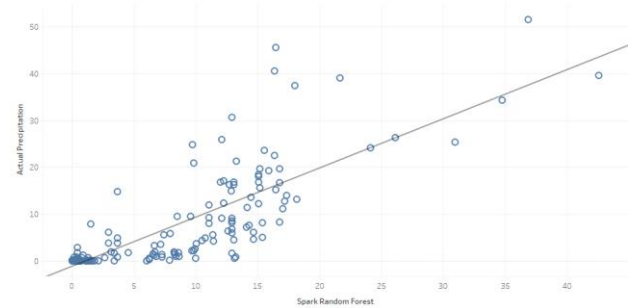
Number of trees vs. RMSE



IV. MODEL RESULTS

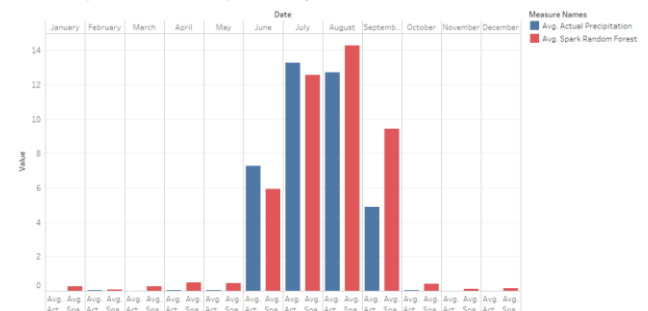
In order to judge the performance of our model, we use testing data from the year 2011, and use this formula to predict the value of precipitation for each day. We first compare the predicted values to the actual values by using a scatter plot in Figure 12.

Actual Precipitation vs Spark Random Forest model predictions



The closer the trend line is to 45°, the better the prediction. We also compare the monthly average predictions in Figure 13.

Actual Precipitation vs KNN model predictions by month



Our model is able to predict rainfall with an average error of ± 1.90 mm.

Rainfall prediction

Latitude
21.58
Longitude
78.12
Max. Temperature (°C)
25.45
Min. Temperature (°C)
22.49
Wind (Kmph)
10.13
Relative Humidity (%)
0.82
Solar Coverage (M.J/m²)
3.17

Predicted Precipitation in mm:
27.4674376444999

Figure 14: GUI

We utilized R and its “Shiny” package to build a basic GUI which allows a user to input the independent values and predicts the output with the click of a button.

V. FUTURE SCOPE

We aim to improve our prediction by utilizing better modelling techniques and market this program in various countries.

We also plan on expanding this project to make it more versatile in terms of automatic data acquisition and user notification.

CONCLUSION

Through this research project we studied the importance of rainfall, its measurement techniques, and their shortcomings and devised a solution to overcome these by integrating predictive analytics techniques from the new and upcoming field of big data analytics.

ACKNOWLEDGMENT

We would like to thank Prof. Vijaykumar Bhanuse for giving us his constant support and this opportunity to work on this project under him. We would also like to thank our honorable director Mr. Rajesh Jalnekar and head of department prof. Dr. Shilpa Sondkar for their inspiration.

BIBLIOGRAPHY

- Dile, Y. T., R. Srinivasan, 2014. Evaluation of CFSR climate data for hydrologic prediction in data-scarce watersheds: an application in the Blue Nile River Basin. *Journal of the American Water Resources Association (JAWRA)* 1-16. DOI: 10.1111/jawr.12182
- Fuka, D.R., C.A. MacAllister, A.T. Degaetano, and Z.M. Easton. 2013. Using the Climate Forecast System Reanalysis dataset to improve weather input data for watershed models. *Hydrol. Proc.* DOI: 10.1002/hyp.10073.
- Lily Ingsrisawang ET. Al. Machine Learning Techniques for Short-Term Rain Forecasting System in the Northeastern Part of Thailand
- National Centers for Environmental Prediction (NCEP) Climate Forecast System Reanalysis (CFSR) globalweather.tamu.edu