# RAINFALL IN INDIA
# IS603

FNU Naveen Kumar, Sudhanshu Jha, Rithvik Pabba, Srinija Dannamanane, Swathi Kotla
IS Department, University Of Maryland, Baltimore County

## ABSTRACT

By reviewing enormous amounts of data, Data Science has grown increasingly popular in identifying hidden patterns, correlations, and other insights. As India is a country where agriculture is the only resource of income for most of population, in which rainfall plays a major role. By predicting the amount of rainfall, we can help farmers to know which kind of crops can be grown time of the year. "The Rainfall in India" dataset is used to analyze and visualize the rainfall in past 45 years (1970-2015).

## INTRODUCTION

To perform the data analysis on the data set for rainfall in India prediction,We have taken two dataset one that distributes data of Indian rainfall with respect to states and districts. At first from each state some districts data were collected, and data cleansing is performed on the data set. Data visualization is used to show the rainfall prediction in different region based on district. Data-mining techniques like Classification, Clustering, Regression are used and tools like Weka for analysis and python for visualization is used.

## DATA CLEANSING

In order to employ data mining methods, data cleansing requires transforming raw data into well-formed data sets. Raw data, is termed as incomplete and structured incorectly. Previuously we collected data for 115 years i.e., 1901-2015, and as there was some missing data in the data set such as inclusion of old region names we analysed and cleansed the data from 1971-2015.



BEFORE CLEANSING



AFTER CLEANSING

## METHODLOGY

```
=== Cross-validation ===
=== Summary ===


Correlation coefficient              0.9999
Mean absolute error                 27.5457
Root mean squared error             42.2526
Relative absolute error              4.7141 %
Root relative squared error          5.0469 %
Total Number of Instances            631
```

REGRESSION  CORELATION COEFFICIENT 0.99
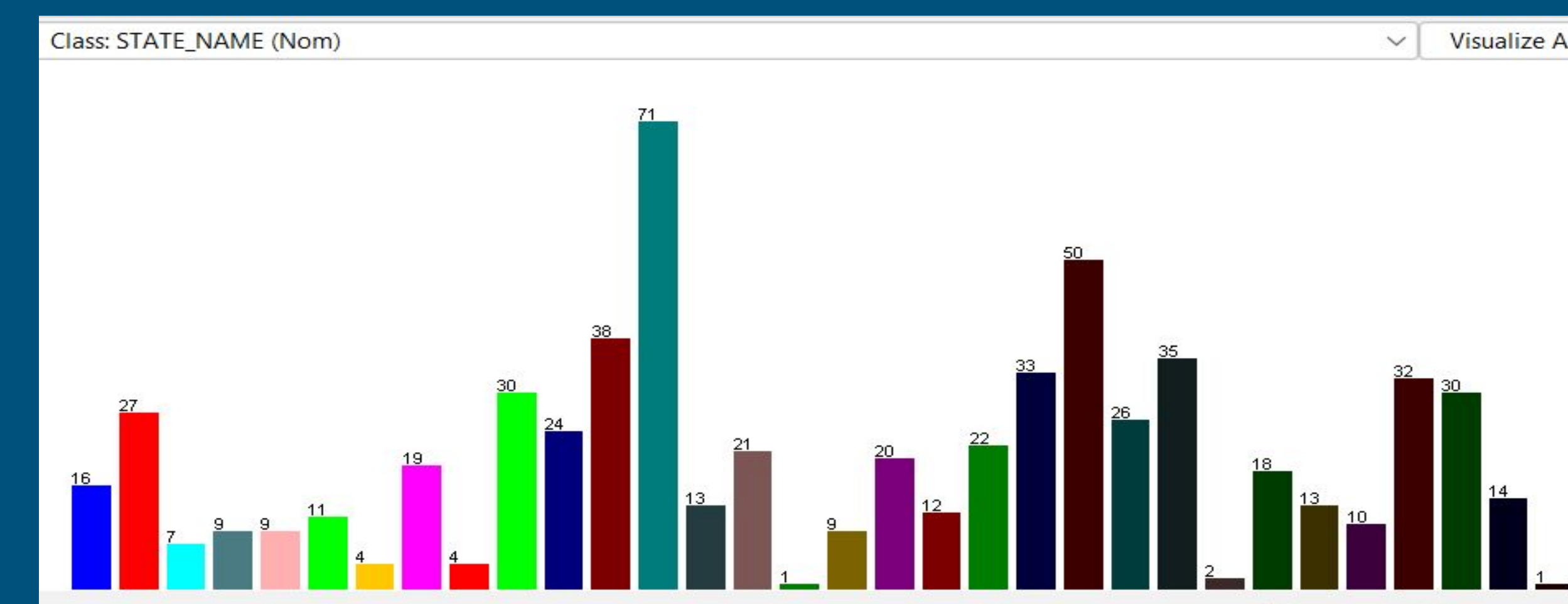
```
=== Model and evaluation on training set ===

Clustered Instances

0        38 (  6%)
1       273 ( 43%)
2        26 (  4%)
3       182 ( 29%)
4        21 (  3%)
5        91 ( 14%)
```
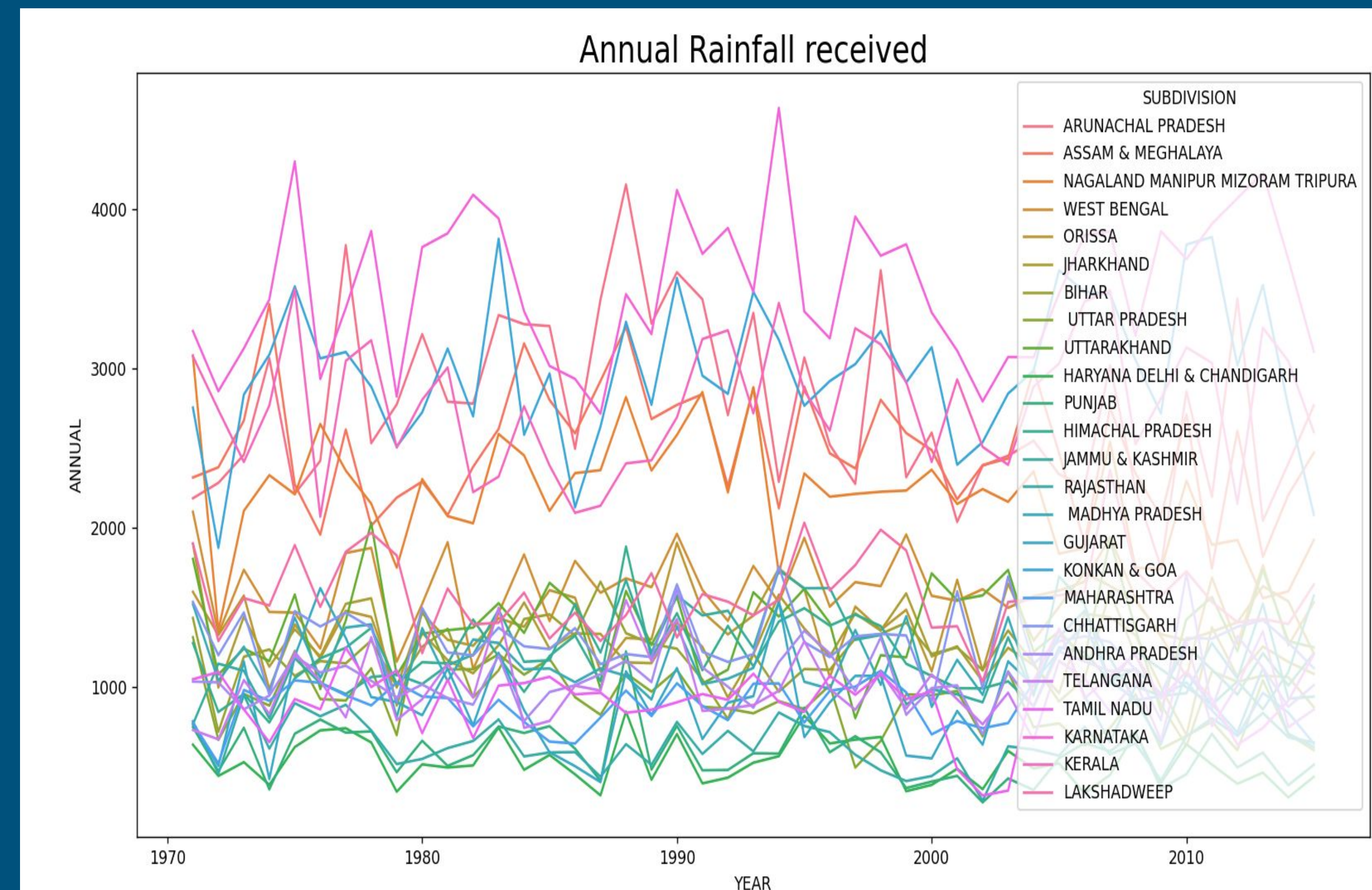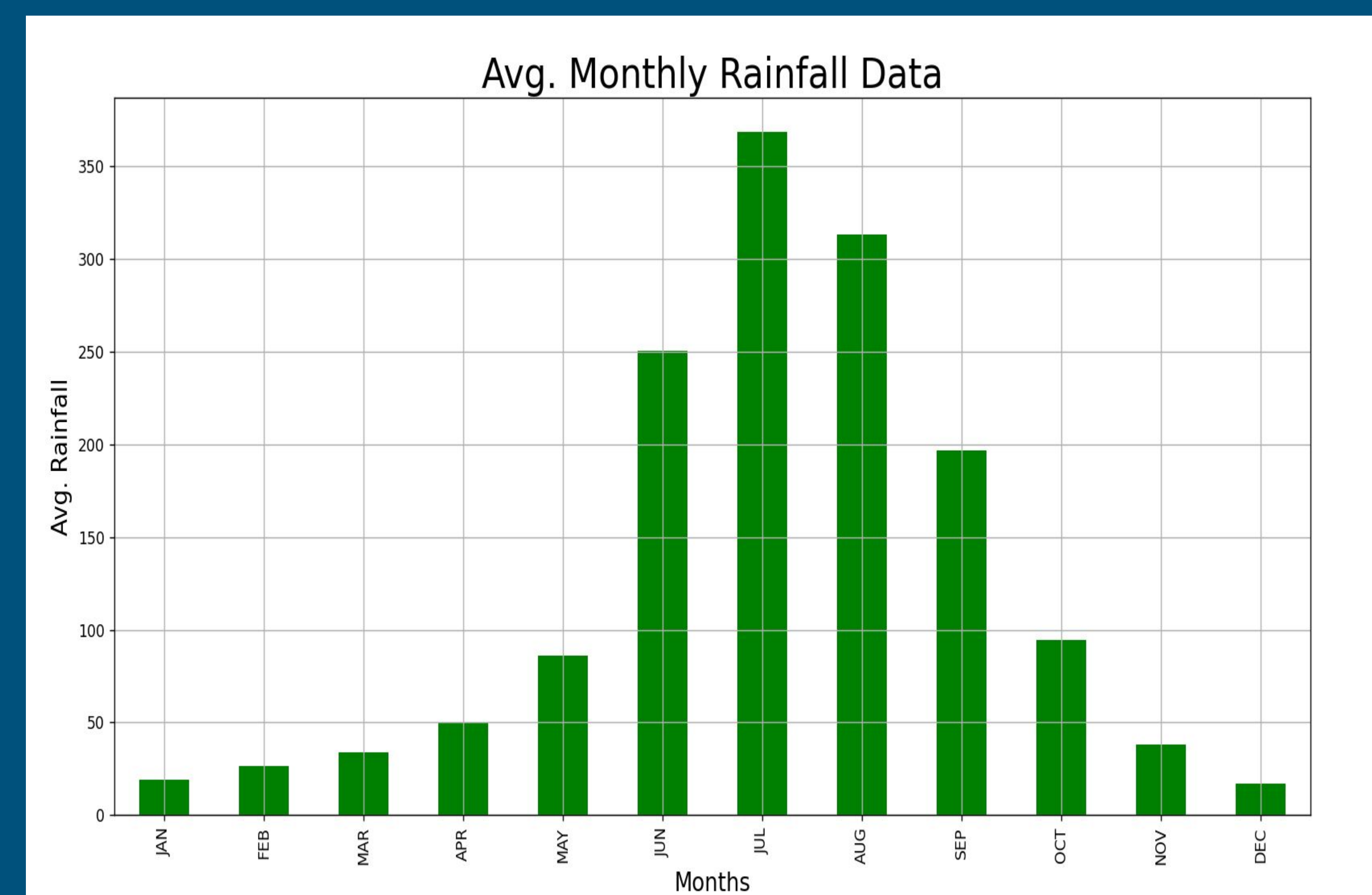
CLUSTERING WITH SIX CLUTER

## RESULT



STATES CLUSTERED BY ITS DISTRICTS

## DATA VISUALIZATION



YEARLY RAINFALL IN EACH STATE



AVERAGE RAINFALL IN EACH MONTH

## CONCLUSION AND FUTURE WORK

The data mining tasks we executed on "Rainfall in India" revealed data that helps farmers estimate rain effectively using techniques such as regression and clustering. Currently machine learning is used in no. industries. As the data increases the complexity of that data will increase and for that we are using tools such as weka and python to better understand the data. In future we are planning to increase our work in Storm predictions and Crop prediction with the rainfall prediction.
Because farmers' survival is strongly based on rainfall, their suicide rates can be forecast using future data analysis.

## REFERENCES

1) https://www.kaggle.com/rajanand/rainfall-in-india
2) Provost, F., & Fawcett, T. (2013). Data science for business: [what you need to know about data mining and data-analytic thinking]. Sebastopol, Calif.: O'Reilly