# Rainfall Prediction in India

**Mentor :** Dr.Sreeja SR

**project Code :** B21SSR02

**Topic : Rainfall Prediction**



**Department of Computer Science and Engineering**
**Indian Institute of Information Technology Sri City**

# Rainfall Prediction

**Group-20:**
**(Team Members)**

**G.Sai Divya(S20180010059)**

**G.Sai Mani(S20180010057)**

**R.Bhavani(S20180010024)**

# Outline of Presentation:

- Problem Statement

- Challenges/Motivation

- Literature Survey

- Dataset Description

- Missing value methods

- Visualization

- Machine learning techniques

- Work done upto know

- Future Plan
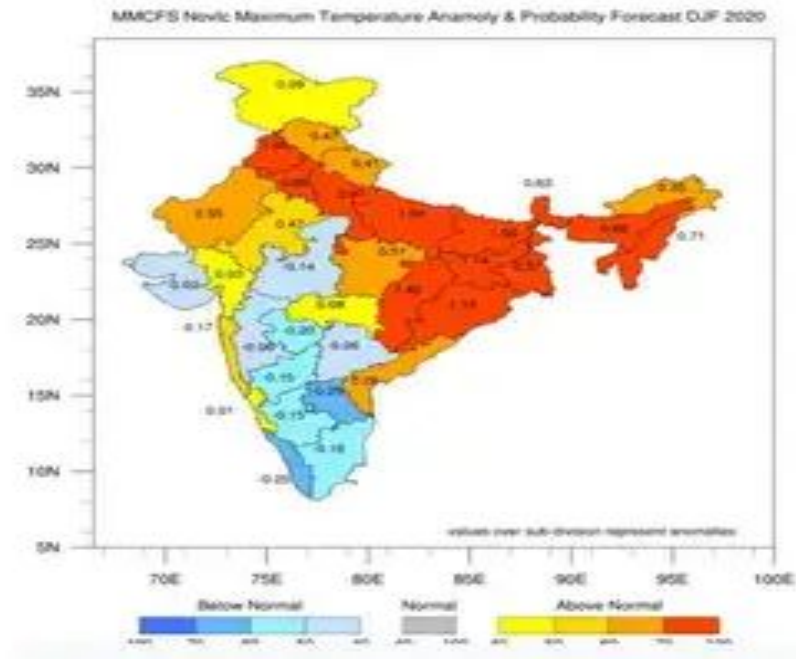
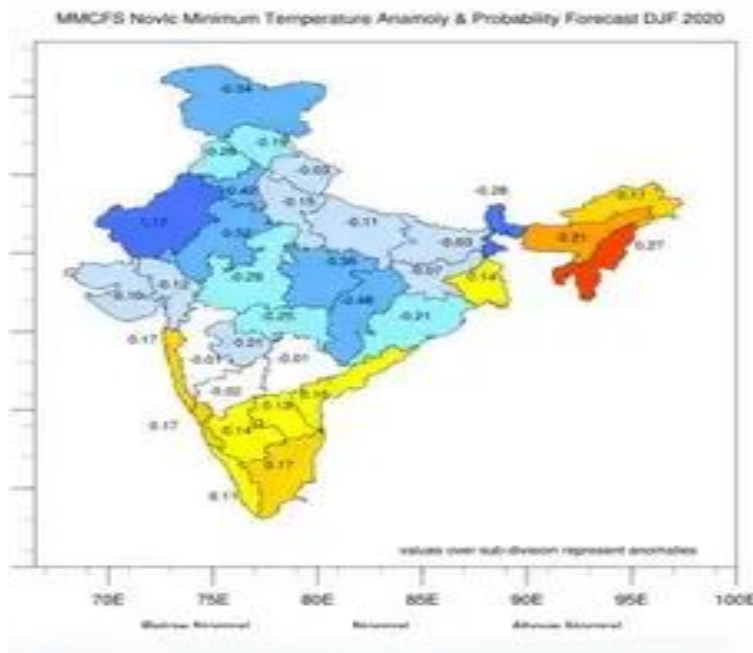- References

# Problem Statement:

- To Study the seasonal variations of rainfall and Annual Average Rainfall Variation of years range from 1901 to 2015 with month wise distribution of rainfall in mm for all major Indian states and Union Territories.

- Predict the Rainfall for upcoming years in India based on dataset we have .

- Display the predictions through website.

# Challenges/Motivation:

- Rainfall is a key part of hydrological cycle and alteration of its pattern directly affect the water resources. The changing pattern of rainfall in consequence of climate change is now concerning issues to water resource managers and hydrologists.

- The drought and flood like hazardous events can be occurred frequently because of the extreme changes of rainfall trend. the amount of soil moisture for crop production is totally determined by the amount of rainfall.

# Challenges/Motivation:

- The monsoon rainfall plays a vital role for agriculture in India. Hence, the research on the climate change or most specifically on the changes of rainfall occurrences and its allocation are the most significant way for sustainable water resource management.

# Literature Survey:

| Authors | Published on, year | Technique | Dataset | Advantage |
|---|---|---|---|---|
| Bushra Praveen et al., IIT INDORE [1] | Scientific Reports<br><br>JUNE 2020 | Artificial Neural Network-Multilayer Perceptron (ANN-MLP) | Indian Rainfall from 1901-2015 | analyzes and forecasts the long-term Spatio-temporal changes |
| Nikhil Oswal., University of Ottawa,canada [2] | arXiv:1910.13827<br><br>October 2019 | Linear classier, Tree based, Distance based, Rule based and Ensemble for original, undersampled, oversampled dataset | https://www.kaggle.com/jsphyg/weather-dataset-rattle-package | comparative study , concentrating on three aspects: modeling inputs, modeling methods, and pre-processing techniques. |
| Emilcy Hernandez et al., Coventry University UK [3] | International Conference on Hybrid Artificial Intelligence Systems. (2016) | Autoencoder and MLP, MSE RMSE errors | meteorological station located in a central area of Manizales, Colombia. (almost 30yrs) | Autoencoder-feature treatment in time series. MLP-prediction task |

# Literature Survey:

| Authors | Published on, year | Technique | Dataset | Advantage |
|---|---|---|---|---|
| Pritpal Singh et al., Tezpur University [4] | _Stochastic Environmental Research and Risk Assessment_ <br><br> 13 Feb 2013 <br><br> @Springer | ANN using BPNN with MLFF | time series data for 140 years (1871– 2010) | a model for prediction of ISMR on monthly and seasonal time scales, its applicability for advance prediction of the seasonal rainfall amounts. |
| Kumar Abhishek et al., NIT Patna [5] | Control and System Graduate Research Colloquium (ICSGRC), 2012 IEEE | BPA, LRN, CBP with different adaptive learning function, training functions | Data from 1960 to 2010. | The multilayered artificial neural network with learning by back-propagation algorithm configuration is the most common in use, due to of its ease in training |

# Dataset Description:

The data set consists of the measurement of rainfall from year 1901-2015 for each state.

- Data consists of 19 attributes (individual months, annual, and seasonal months) for 36 sub divisions.
  [subdivision, year, jan, feb, mar, apr, may, june, july, aug, sep, oct, nov, dec, annual, jan-feb, mar-may , jun-sep, oct-dec]

- The attributes are amount of rainfall measured in mm

**Simple glimpse of data:**

# Dataset Description:

# Missing Values methods:

**Ignore the missing values :**

❖ Missing data under 10% for an individual case or observation can generally be ignored

❖ **Drop the missing values:**

**Dropping a variable**

❖ The number of missing values in a feature is very high in a data, then that feature should be left out of the analysis. [6]

**Imputation by Mean/Mode/Median**

❖ Imputation is the process of substituting the missing data by some statistical methods. [2]

❖ If the missing values in a column or feature are numerical, the values can be imputed by the mean of the complete cases of the variable. [2]

# Missing Values methods:

## Regression Methods

❖ The variables with missing values are treated as dependent variables and variables with complete cases are taken as predictors or independent variables. [7]

❖ The independent variables are used to fit a linear equation for the observed values of the dependent variable. This equation is then used to predict values for the missing data points. [7]

## K-Nearest Neighbour Imputation (KNN)

❖ This method uses k-nearest neighbour algorithms to estimate and replace missing data.

❖ The k-neighbours are chosen using some distance measure and their average is used as an imputation estimate. [8]

❖ This could be used for estimating both qualitative attributes (The most frequent value among the k nearest neighbours) and quantitative attributes (the mean of the k nearest neighbours). [8]

# Applying Missing Values method:

None and NaN are used for indicating missing or null values.

To facilitate this convention, there are several useful methods. They are:

❖ **isnull():** Generate a boolean mask indicating missing values

❖ **notnull():** Opposite of isnull()

❖ **dropna():** Return a filtered version of the data

❖ **fillna():** Return a copy of the data with missing values filled or imputed

So first we will find total missing values using sum:

dataset.isnull().sum()

Then we apply missing values method using fillna method

Again we check the missing value sum to confirm that no missing values in our dataset

# Missing values for the dataset :

❖ **Mean:** Mean (commonly known as average) is equal to the sum of all values in the column divided by the number of values present in the column. In Excel, you can use the AVERAGE() function to compute the mean.

**dataset = dataset.fillna(dataset.mean())**

❖ **pad():** This method is similar to the dataframe.fillna(). method and it fills NA/NaN values using the ffill() method.It returns the object with missing values filled or none if replace=True.

**dataset = dataset.fillna(method ='pad')**

❖ **bfill():** is used to backward fill the missing values in the dataset. It will backward fill the NaN values that are present in the pandas dataframe.

**dataset = dataset.fillna(method ='bfill')**

# Missing values for the dataset:

| Before applying mean method | After applying mean method |
| --- | --- |
| SUBDIVISION    0<br>YEAR           0<br>JAN            4<br>FEB            3<br>MAR            6<br>APR            4<br>MAY            3<br>JUN            5<br>JUL            7<br>AUG            4<br>SEP            6<br>OCT            7<br>NOV           11<br>DEC           10<br>ANNUAL        26<br>Jan-Feb        6<br>Mar-May        9<br>Jun-Sep       10<br>Oct-Dec       13<br>dtype: int64 | SUBDIVISION    0<br>YEAR           0<br>JAN            0<br>FEB            0<br>MAR            0<br>APR            0<br>MAY            0<br>JUN            0<br>JUL            0<br>AUG            0<br>SEP            0<br>OCT            0<br>NOV            0<br>DEC            0<br>ANNUAL         0<br>Jan-Feb        0<br>Mar-May        0<br>Jun-Sep        0<br>Oct-Dec        0<br>dtype: int64 |

```
[27] data = data.fillna(data.mean())
     data.head()
```

| | SUBDIVISION | YEAR | JAN | FEB | MAR | APR | MAY | JUN | JUL | AUG | SEP | OCT | NOV | DEC | ANNUAL | Jan-Feb | Mar-May | Jun-Sep | Oct-Dec |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 0 | ANDAMAN & NICOBAR ISLANDS | 1901 | 49.2 | 87.1 | 29.2 | 2.3 | 528.8 | 517.5 | 365.1 | 481.1 | 332.6 | 388.5 | 558.2 | 33.6 | 3373.2 | 136.3 | 560.3 | 1696.3 | 980.3 |
| 1 | ANDAMAN & NICOBAR ISLANDS | 1902 | 0.0 | 159.8 | 12.2 | 0.0 | 446.1 | 537.1 | 228.9 | 753.7 | 666.2 | 197.2 | 359.0 | 160.5 | 3520.7 | 159.8 | 458.3 | 2185.9 | 716.7 |
| 2 | ANDAMAN & NICOBAR ISLANDS | 1903 | 12.7 | 144.0 | 0.0 | 1.0 | 235.1 | 479.9 | 728.4 | 326.7 | 339.0 | 181.2 | 284.4 | 225.0 | 2957.4 | 156.7 | 236.1 | 1874.0 | 690.6 |
| 3 | ANDAMAN & NICOBAR ISLANDS | 1904 | 9.4 | 14.7 | 0.0 | 202.4 | 304.5 | 495.1 | 502.0 | 160.1 | 820.4 | 222.2 | 308.7 | 40.1 | 3079.6 | 24.1 | 506.9 | 1977.6 | 571.0 |
| 4 | ANDAMAN & NICOBAR ISLANDS | 1905 | 1.3 | 0.0 | 3.3 | 26.9 | 279.5 | 628.7 | 368.7 | 330.5 | 297.0 | 260.7 | 25.4 | 344.7 | 2566.7 | 1.3 | 309.7 | 1624.9 | 630.8 |

# Visualization :

## 1. HeatMap:
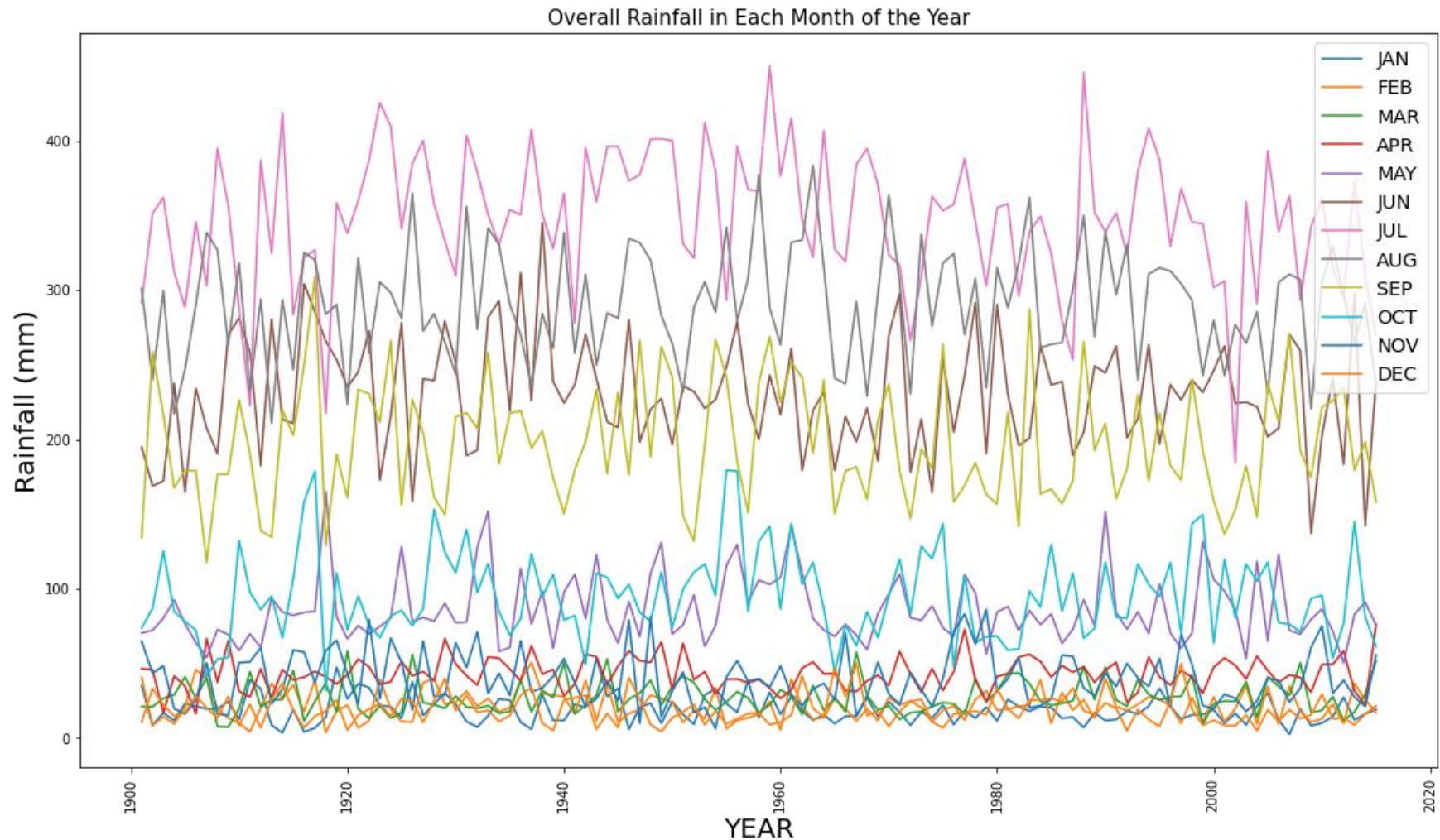


Correlation matrix

# Visualization :

## 2. BarPlot:



Month wise rainfall in each subdivision

# Visualization :

## 2.  BarPlot:



seasonal wise rainfall in each subdivision

## 3.  Line Plot:



Overall Rainfall in Each Month of the Year

# Visualization :

## 3. Line Plot:

# Visualization :

## 3.    Line Plot:



Annual Rainfall Trend in India

## 4. Box Plot:



Annual Rainfall in Subdivisions of India

# Applying Machine Learning Techniques:

**Process:**

**Step1:** Import the rainfall data set.csv file.

**Step2:** Fill the missing values with mean value of the data.

**Step 3:** The data is divided into training set (80%) and testing set (20%).

**Step 4:** Apply machine learning or deep learning algorithms and calculate the errors like Mean Absolute Error, Mean Squared Error.

**Step 5:** Choose the best Model among all models.

The models we used are: Elastic net, Lasso, Ridge, SVR, Random forest regression,Neural network model.

For prediction we formatted data in the way, given the rainfall in the last

three months we try to predict the rainfall in the next consecutive month.

# Applying Machine Learning Techniques:

| Model | Mean Absolute Error | Mean Squared Error |
|-------|---------------------|--------------------|
| Elastic Net Regression | 93.3734 | 142.5711 |
| Ridge Regression | 93.9519 | 141.8863 |
| Lasso Regression | 93.9520 | 141.8863 |
| SVR | 125.1810 | 202.5503 |
| Random Forest Regressor | 91.7085 | 145.7345 |
| Neural Network | 85.1830 | 135.8926 |

We Observe that all the linear models got almost same MAE and MSE.

Neural Network got less errors among all other methods , second better algorithm was Random Forest Regression.

So we want to proceed using neural network model for upcoming work.

# Applying Machine Learning Techniques:

```
[10]  subdivs

       array(['ANDAMAN & NICOBAR ISLANDS', 'ARUNACHAL PRADESH',
              'ASSAM & MEGHALAYA', 'NAGA MANI MIZO TRIPURA',
              'SUB HIMALAYAN WEST BENGAL & SIKKIM', 'GANGETIC WEST BENGAL',
              'ORISSA', 'JHARKHAND', 'BIHAR', 'EAST UTTAR PRADESH',
              'WEST UTTAR PRADESH', 'UTTARAKHAND', 'HARYANA DELHI & CHANDIGARH',
              'PUNJAB', 'HIMACHAL PRADESH', 'JAMMU & KASHMIR', 'WEST RAJASTHAN',
              'EAST RAJASTHAN', 'WEST MADHYA PRADESH', 'EAST MADHYA PRADESH',
              'GUJARAT REGION', 'SAURASHTRA & KUTCH', 'KONKAN & GOA',
              'MADHYA MAHARASHTRA', 'MATATHWADA', 'VIDARBHA', 'CHHATTISGARH',
              'COASTAL ANDHRA PRADESH', 'TELANGANA', 'RAYALSEEMA', 'TAMIL NADU',
              'COASTAL KARNATAKA', 'NORTH INTERIOR KARNATAKA',
              'SOUTH INTERIOR KARNATAKA', 'KERALA', 'LAKSHADWEEP'], dtype=object)
```

```
   meanerror

       array([140.4678507 , 199.28770968, 109.12406223,  94.3462306 ,
              128.72484457,  77.49932806,  81.35422742,  67.86587383,
               77.17628137,  61.79250835,  62.64865462,  88.31259251,
               40.20679799,  46.87334736,  66.61765539,  57.08613747,
               28.17995848,  41.28195253,  63.45461393,  78.16995787,
               80.97948572,  50.93588132, 200.61164137,  54.52411525,
               60.60954418,  64.22243438,  75.99868872,  64.00701032,
               60.46937416,  50.05342925,  54.359456  , 185.69336385,
               49.22969688,  61.56358943, 189.81332926, 104.37353324])
```

# Work done upto now/Future plan:

## Work done upto now:

❖   Importing data
❖   Applying missing value method
❖   Visualization  (Exploratory Data Analysis)
❖   Predicted the best model using MAE, MSE.

## Future plan:

Based on the model we try to predict future values, by changing the data every time like appending the result at the end, as we are taking the data as an array.

Based on the future values we want to do an application

# References:

1.  Praveen, Bushra, et al. "Analyzing trend and forecasting of rainfall changes in India using non-parametrical and machine learning approaches." *Scientific reports* 10.1 (2020): 1-21.

2.  Oswal, Nikhil. "Predicting rainfall using machine learning techniques." *arXiv preprint arXiv:1910.13827* (2019).

3.  Hernández, Emilcy, et al. "Rainfall prediction: A deep learning approach." *International Conference on Hybrid Artificial Intelligence Systems*. Springer, Cham, 2016.

4.  Singh, Pritpal, and Bhogeswar Borah. "Indian summer monsoon rainfall prediction using artificial neural network." *Stochastic environmental research and risk assessment* 27.7 (2013): 1585-1599.

5.  Abhishek, Kumar, et al. "A rainfall prediction model using artificial neural network." *2012 IEEE Control and System Graduate Research Colloquium*. IEEE, 2012.

# References:

6. Kang, Hyun. "The prevention and handling of the missing data." *Korean journal of anesthesiology* 64.5 (2013): 402.

7. Mostafa, Samih M. "Imputing missing values using cumulative linear regression." *CAAI Transactions on Intelligence Technology* 4.3 (2019): 182-200.

8. De Silva, Hiroshi, and A. Shehan Perera. "Missing data imputation using Evolutionary k-Nearest neighbor algorithm for gene expression data." *2016 Sixteenth International Conference on Advances in ICT for Emerging Regions (ICTer)*. IEEE, 2016.

# Thank you!

# ANY QUESTIONS?