

# A Method for Weather Forecasting Using Machine Learning

Prathyusha  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

bobba\_prathyusha@srmap.edu.in

Zakiya  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

pathan\_zakiya@srmap.edu.in

Savya  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

adudotla\_savya@srmap.edu.in

Tejaswi  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

kata\_tejaswi@srmap.edu.in

Neena Alex  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

neena\_alex@srmap.edu.in

Dr. Sobin C C  
Dept.of CSE  
SRM University AP  
Andhra Pradesh  
India

sobin.c@srmap.edu.in

**Abstract**— Agriculture is a sector that plays a crucial role in the economies of many countries around the globe, like India where it contributes 16% of the total economy. Weather forecasting is one of the challenges faced by this sector, due to its dynamic and turbulent nature, the statistical methods fail to provide forecasting at an accurate precision. This paper aims to develop an accurate way to predict the temperature forecast using machine learning techniques especially using Long Short Term memory networks (LSTM). Despite the advances made, there are still significant obstacles to overcome in expanding the use of weather forecasts in the agricultural sector due to the dynamics in climate changes. These include the need for improved model accuracy, quantitative evidence of the utility of climate predictions as instruments for agricultural risk management and addressing major chances of disease incidence which are usually seasonal and depends on parameters like temperature and rainfall. The goal of this study is to forecast parameters that could help farmers to make an informed decision so as to reduce the losses by taking required proactive measures. This paper provides a detailed analysis of weather forecasting techniques and explores future research goals in this field.

**Keywords**— *Weather Forecasting, Machine Learning*

## I. INTRODUCTION

Weather forecasting is one of the crucial and complex tasks that is consummated by meteorologists. Weather forecasting answers the basic questions like what is the weather expected to be tomorrow? Is it going to rain today? One may conclude whether it is going to be sunny, foggy, misty, or cloudy on a particular day and plan their business accordingly. Considering the significance of forecasting in everyday life, meteorologists strive for near accuracy in their predictions.

The agriculture sector is another field that depends on the forecast to a vast extent. Weather accounts for the annual profit or loss of farm production directly or indirectly. In many countries, agriculture is the main source for their economic development. Crop loss can be reduced by making adjustments based on timely and accurate weather forecasts. So, getting a gist of the factors like the amount of humidity precipitation, temperature, upcoming rainfall, or precisions of natural disasters like floods, droughts, storms, hurricanes, etc. helps the farmers to manage their jobs, minimize the damage of property, and selecting crops that are most suited

to the predicted climatic conditions. To successfully foster growth and ensure food security in this changing environment, accurate weather prediction is needed.

There are many techniques and algorithms that are used for predictions. Weather forecasting has time series data and in this paper a temperature prediction using Auto Regressive Integrated Moving Average (ARIMA) model and LSTM is deployed. It helps to predict the temperature of a particular season which is beneficial in agriculture to be known well in advance for early identification and hence mitigation of diseases in crops.

The machine learning models usually consider three types of input features usually to be fed into the models. One consists of using the meteorological or environmental variables. Second category uses the historical temperature data and the final one which combines the former and later. Similarly the performance of the machine learning models hugely depends on the horizon of forecast time whether it is a short term prediction or long time. Another factor which affects the performance is the spatial scale. Global scale predictions are found to have smaller errors than the local scale predictions which are done for one particular station.

Most of the researches in the area of short term temperature prediction focuses more on Artificial Neural Network (ANN) strategies than Support Vector Machines (SVMs). Multiple-layer Perceptron Neural Network (MLPNN) and Radial Basis Function Neural Network (RBFNN) are the widely used architectures in ANN approaches. Due to the learning rate and lesser errors, the Levenberg-Marquardt algorithm is found to be more useful for optimizations[1]. The advancements in deep learning earned much momentum to Recurrent Neural Network (RNN) and LSTM in the fields of finance, computer vision and natural language processing. LSTMs are highly promising for long term multivariate time series forecasting. Along with making temperature predictions, this paper tries to analyze the performance of LSTM for short term univariate regression problem.

The paper also describes some of the techniques and methodologies regarding weather forecasts. Based on this comprehensive study we attempt to determine some of the accurate and precise techniques so that further studies can be made in these areas for better results. The following section

of the paper consists of related works and then relevant methods. The third section of the paper deals with the ARIMA and LSTM model for forecasting temperature of a particular station, which is followed by results and conclusion.

## II. RELATED WORK

This section discusses some of the existing machine learning models already implemented for weather forecasting in the agricultural sector.

Mohan et al. [2] have worked on finding an appropriate information model, which helps in accomplishing high precision and simplification for value forecast. This research introduces an effective technique, namely SOM- LDA (Self Organizing Map – Latent Dirichlet Allocation known as weighted- SOM) model is introduced for performing crop and weather prediction. The proposed data set contains more than 1000 records which contain the attributes such as crop type, soil type, rainfall rate, temperature, humidity, etc. Data mining techniques were employed to know which environment was suitable for particular types of yields and which weather is suitable for which crops. From the early years, these strategies aid in predicting rainfall, moisture, temperature, and wind speed. The crop's productivity increased dramatically as a result of this prediction. Weather forecasting is processed in two phases like dimension reduction and classification. In the dimension reduction phase, the reduced data is used for classification by employing DNN. Along with the WSOM approach, the DNN approach was utilized as a classifier to enhance the accuracy of the prediction rate. The performance of the proposed method was estimated based on precision, recall, sensitivity, specificity, and accuracy value. The proposed method achieved 80.09% accuracy. Finally, the experimental outcome showed that the proposed approach improved accuracy in weather and crop prediction by 7-23% when compared to the existing methods.

Lennard et al. [3] discussed a supervised scheme: SOM for surface rainfall analysis correlated with synoptic circulation. It was investigated for two types of stations in South Africa's distinct rainfall zones. This work assesses the ability of SOMs to match the synoptic movers of observed rainfall records, effectively downscaling synopses of large-scale data to a resolute response of the surface. The relationship between synoptic-scale circulation and as a result, rainfall response is investigated at two sites in South Africa during a 31-year period in several rainfall regimes. Daily rainfall is related to circulation archetypes at the two stations, and trends in the frequency of occurrence of each archetype and daily rainfall are explored using a bootstrapping methodology. A histogram of the slopes of the linear regressions is constructed using the linear regression of each bootstrapped replicate. If it falls above the 95th percentile the trend is taken into account significantly. Self-organizing maps were used to compress multivariate atmospheric data into generalized modes of synoptic-scale circulation across South Africa. They are a type of artificial neural network that reduces high-dimensional data space to produce a low-dimensional (typically two-dimensional), discretized representation of the input data and presents the multidimensional data distribution function's major characteristics as an array or map of data archetypes. To improve the accuracy of rainfall detection, certain features in SOM were necessary for rainfall prediction.

Weather prediction depends on various parameters like temperature, rainfall, humidity, wind speed, which vary from time to time and also with different geographical locations. This paper highlights the model which uses a Decision Tree to predict weather phenomena. The proposed model in [4] was implemented using the open-source data mining tool Rapid miner. The weather data set of the years 2013 and 2014 was collected. First, a model was created by training with a sample data set called the training set which contained all the attributes in the data set of 2013. Then the trained model was provided with a test set that contains all the attributes in the data set of 2014 except the rainfall. The decision trees were built by recursively splitting the values of attributes. These attributes were decided by the information gain (the one attribute having minimum entropy is selected), and accuracy of the decision tree. The validation operation is also used to check the performance. An accuracy of 100% rainfall prediction was shown on the training data but when done on scoring data 80.67% was achieved.

A toy model for global weather forecasting was demonstrated, and obstacles identified in [5]. It depicted the basic design decisions for a forecast system based on the use of neural networks. There are no dynamical equations of motion used in the toy model. In this paper they used the toy model to identify issues and suggest essential design choices for NN-based forecast systems that could lead to optimal results. A thorough grasp of how to improve the construction of network architectures and network training, as well as how to preserve conservation properties, by utilizing physical knowledge of the Earth system and the connectivity between degrees of freedom will be necessary. On one hand, it is expected that future computing hardware will allow NN models to make greater use of more observations and higher resolution. On the other hand, as it will be difficult to increase physical consistency inside networks, it is probable that stabilizing long-term integrations and representing intricate interactions between model features in extended simulations will be challenging.

Luan et al. [6] developed a system that rationalized the overall business of agricultural drought monitoring and forecasting, as well as irrigation amount forecasting, on the basis of enhanced accuracy. Based on IDL hybrid programming, an integrated business service system with strong cohesion but low coupling was designed within the four-layered architecture framework. Drought monitoring, forecasts, and irrigation water requirements are all integrated and closely linked under the direction of data. The system made extensive use of IDL's parallel computing technology. They established a system platform that combined drought monitoring and defense-related activities into a single entity and performed accurately as well as quickly. Due to the complexity of the agricultural drought and the disaster-inducing environment, related algorithms and models must be reviewed and enhanced on a regular basis. Meanwhile, given the multi-source nature of drought disasters, how to improve the system's functions and add features such as drought impact assessment so that it can better support the country's disaster prevention and reduction efforts will be a key research focus on the future.

Xuan et al. [7] discussed the networks for flood forecasting. The daily discharge and rainfall were employed as input data to the LST) neural network model for flood forecasting. It evaluates the model's effectiveness in flow rate

forecasting and the impact of input data features on the model's flood forecasting capability, two scenarios were studied. In projecting the maximum flow rate using separate data sets, the model was validated and assessed using metrics such as NSE value, RMSE value, and relative error values. Both scenarios performed admirably during the validation and testing phases. Furthermore, the simulation findings for all three forecast instances show no significant differences. However, when it comes to the second scenario, it showed somewhat better-than-expected outcomes due to the flood peak predicting factor. Despite the fact that the LSTM model effectively tackles sequential data problems, there are a few drawbacks to consider. LSTM models, in general, are data-driven models like physical-based models, fail to adequately simulate hydrological processes. As a result, these models should be integrated with meteorological models such as rainfall forecasting models to improve long-term forecast performance. The results of the research pointed to the possibility of using the LSTM model in the field of hydrology to design and manage real-time flood warning systems. This concept remains a possible choice for underdeveloped countries or huge river basins.

A study in [8] reviews deep learning models - ARIMA Model, ANN, SVM and Self Organizing Map (SOM). The proposed model predicts the precision of rainfall using the past meteorological data from weather radars as input. The required prediction is acquired by applying the techniques, Multi-Layer Perceptron (MLP) and Auto-Encoders using the non-linear features in the data. The meteorological data from weather radars [9] was taken and divided into training and testing datasets for implementation. The best Epoch value used from the training set is identified as it indicates the least Mean absolute error and was tested against with the testing set. The result from this process is combined to get the best result at the targeted value. This paper compares the existing methodologies and claims that in terms of RMSE and MSE the results exceed the remaining approaches. It also claims that due to the non-linear relationships in the datasets and their capability to learn from the past, ANN gives the best outcome in the approaches available.

A Survey by Naveen et al. [10] described the performance statistics of rainfall forecasting using various methods. The main objective of this paper is to develop a system to multi-scale forecast of weather and climate using Artificial Neural Networks and to understand the statistics of different models like Radical Basis Functions (RBF) using Neural Networks (NN), Genetic Algorithm (GA) and Hybrid Particle Swarm Optimization and Genetic Algorithm (HPSOGA). These techniques are implemented against the data collected from various weather forecasting stations on an hourly, monthly or yearly basis. The authors divided the data into two samples (480 and 80) and tested their accuracy using RMSE, CC and AARE techniques. They also compared Gaussian and Wavelet SVM approaches for water level forecasting and concluded that Wavelet SVM gave better results. This paper claims that as data is impacted by various sorts of irritations, using a more logical philosophy in gauging results in a more deterministic capacity with infrequent conditions instead of solitary deterministic capacity gives better results. The author [11] assessed five leading edge Numerical Weather Prediction (NWP) frameworks. As reference information utilized has gridded precipitation given by the Climate Prediction Centre (CPC), they have considered the figure lead-times up to 5 days. In

addition, to benchmark the expertise of these models, they considered precipitation gauges from one radar-based (Stage IV) and four satellite-based models for up-to 6 days.

The table I below shows the summary of the related works which presented weather-based forecasting models.

TABLE I. SUMMARY OF LITERATURE

Paper	Models	Input variables
[4]	Weather forecasting (2014)	Temperature, dew point, mean sea level pressure, mean station pressure, visibility, wind speed, sustained wind speed, precipitation
[3]	Rainfall response to synoptic states (2015)	Mean Sea-Level Pressure, moisture transport, mid-latitude cyclone
[6]	Forecasting irrigation amount and drought monitoring (2015)	Satellite images, evapotranspiration, soil texture interpolation, land coverage, albedo and vegetation coverage, temperature, relative humidity, sunlight, wind speed, and rainfall
[2]	Weather forecasting (2018)	Crop type, soil type, rainfall rate, temperature, humidity, etc.
[5]	simulation of dynamics of global atmosphere (2018)	Reanalysis data set, global snapshots of Z500 with 1860 grid points
[7]	Flood forecasting (2019)	Maximum daily precipitation, peak flood discharge values, daily flowrate
[9]	Rainfall forecast (2019)	Meteorological data
[11]	Rainfall prediction (2020)	Meteorological data

### III. PREDICTING TEMPERATURE USING ARIMA MODEL AND LSTM

This section encompasses the techniques used in this study for forecasting the temperature.

#### A. Auto Regressive Integrated Moving Average (ARIMA)

ARIMA model is used to predict, analyze, and forecast temperature. ARIMA (p, d, q) is a standard notation that replaces the parameters with integer values to signal that the ARIMA model is being utilized right away. The four steps of the ARIMA model are as follows:

- In the first stage, a series of responses is identified, which is then utilized to calculate time series and auto correlations using statement IDENTIFY.
- In this stage, the previously selected variables are estimated, as well as the parameters, using the statement ESTIMATE. This stage involves performing diagnostic checks on the variables and parameters that have been collected previously.
- Using the ARIMA model and the expression FORECAST, the predictive values of time series are forecasted, which are future values.

Time series forecasting can be deployed successfully by using ARIMA model. We use the model to forecast average temperature and analyze the trends. For that, the first step is

to analyze the stationarity of data which tells about seasonal patterns in data and the trends. The time series data then undergoes decomposition and then forecasting is done. The ARIMA method assumes stationarity of data.

Dickey Fuller test was performed to check for stationarity of data. The test gives the measures for p-value and critical value which helps to accept or reject hypothesis. The p-values obtained were less than the significance level of 5% and hence the null hypothesis is rejected. As the data we use in this experiment are only of particular months of July and August for five years, it shows stationarity in values. A stationary series has the characteristics that the mean and variance and autocorrelation show a constant pattern and they do not vary much over the time [12].

Analysing for autocorrelation is another important part of the model. It shows how a data is related with its past values in time. Trends in short and long terms of data can be identified from obtaining the autocorrelation values for different lags. Higher values depict strong correlation with the past data [13].

Forecasting of values is done using the ARIMA model. It helps to predict values with respect to the combination of values in the past. We used the *auto-arima* package available in the python library to get the prediction. In spite of predicting time series, the difference of the data from one timestamp ( $a_t$ ) to another ( $a_{t+1}$ ) is calculated in the method. That difference ( $Z_t$ ) forms the 'I' integrated part of the ARIMA model which specifies the middle argument (represented by d) of the model.

$$Z_t = a_{t+1} - a_t \quad (1)$$

The other arguments for ARIMA( $p, d, q$ ) model are the auto regression values represented by  $p$  and the count of moving average terms represented by  $q$  [14]. The data is divided into testing and training sets before feeding to the model. The regression is done for moving averages. The moving average is calculated using rolling mean function. The best ARIMA model for the data is selected based on the minimal AIC value and finally the RMSE value is obtained.

### B. Long Short-Term Memory (LSTM)

In sequence prediction challenges, LSTM networks are a type of recurrent neural network that can learn order dependence. The data are the boundary conditions of a model. The input data is used by the LSTM models to update a variety of values in the internal cell states. The LSTM model does not include mass or momentum conservation principles. As a result, topography data is not included in the input data. The LSTM models, on the other hand, learn these physical principles from the input data and observed data during the training process, and are designed to forecast values as accurately as feasible.

- The first stage in building an LSTM network is to identify information that isn't needed and will be left out of the cell.
- The decision and storage of information from the new input in the cell state, as well as updating the cell state
- The output values are based on the output cell stated in the last phase however they are filtered.

The importance of LSTM is that it can be used to solve sequence problems which are difficult and helps to conquer vanishing gradient problem[15]. LSTMs work on large memory blocks. We have used it in the study for regression problem. The values of temperature are taken for a particular season and predicts the values of temperature for next days. Values of input were converted into floating points which is recommended with models of neural networks. And as LSTMs are sensitive based on the scale of data, the data was normalized using MinMaxScaler class from the scikit-learn library in Python. For estimating the efficiency of the model on unseen data, we split the ordered data into training and testing sets. By creating a data set  $x$  where values are the average temperature on a particular day and  $y$  is the set with temperatures on next day.

The network design has one visible input layer and single hidden layer with four memory blocks and a single value predicting output layer. Sigmoid activation function was used. Predictions were generated for both data sets. Inverse transforms were used to get back the data values in same units of inputs before calculating the RMSE values. The regression problem was solved using the LSTM for predicting the daily temperature.

## IV. RESULTS

Data required for the study was obtained from the UCI Machine Learning Repository. The data set consisted of maximum and minimum temperatures for the month of July and August for consecutive five years [16]. Data was pre-processed and average temperature was calculated for the values.

The time series plot of data is visualized as below in figure 1. It also shows the rolling mean and rolling standard deviation in the graph. The values of temperature are in degree Celsius and are plotted along the Y axis. X axis shows the number of days used to analyse the series.

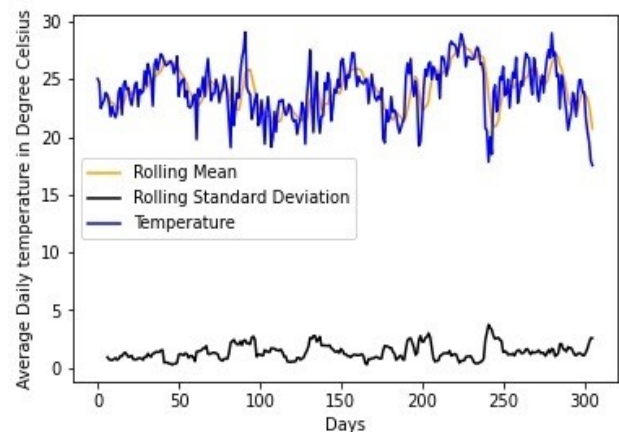


Fig. 1. Rolling Mean and Rolling Standard Deviation

The data from dickey fuller test showed it to be stationary data. The decomposition of data using the *seasonal\_decompose* function shows a stationary pattern and irregular trend which are plotted in the following figure 2. The autocorrelation function value was obtained as 0.77.

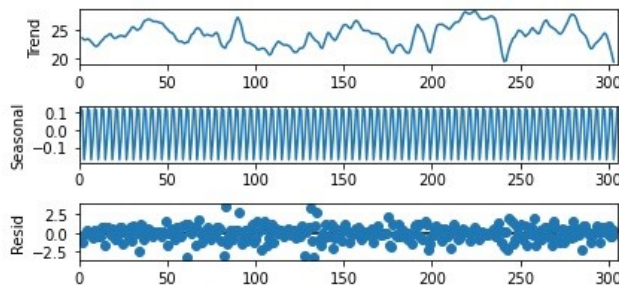


Fig. 2. Trend and seasonality analysis

The ARIMA model was trained with data and tested. The *auto\_arima* model imported from *pmdarima* library of Python also selects the best ARIMA model for the data set. Stepwise search is implemented to minimize AIC. The results were obtained as the table II below.

TABLE II. RESULTS OF AIC VALUES

Model	AIC value	Time
ARIMA(2,0,2)	900.312	0.68 sec
ARIMA(0,0,0)	1089.201	0.01 sec
ARIMA(1,0,0)	902.005	0.14 sec
ARIMA(0,0,1)	964.910	0.07 sec
ARIMA(0,0,0)	2248.164	0.01 sec
ARIMA(1,0,2)	899.345	0.29 sec
ARIMA(0,0,2)	940.244	0.10 sec
ARIMA(1,0,1)	902.492	0.22 sec
ARIMA(1,0,3)	900.511	0.53 sec
ARIMA(0,0,3)	925.053	0.15 sec
ARIMA(2,0,1)	902.015	0.50 sec
ARIMA(2,0,3)	901.759	0.67 sec
ARIMA(1,0,2)	912.606	0.11 sec

As illustrated by the table, the model ARIMA (1,0,2) was selected as the best model with minimal AIC value. It denotes that for autoregression the lag value is set to 1, 0 tells that it is not required to do anything more to make it stationary, and it considers a moving average window of size 1. Finally, the optimal parameters were used to fit the model to sort out the coefficients of regression. And the forecast was made. The prediction is plotted as in figure3 below.

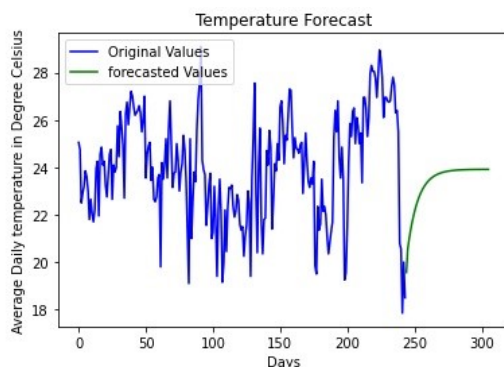


Fig. 3. Average Temperature values predicted using ARIMA model

The predicted outcomes in the study are solely dependent on the actual collected data at the participating gauge station. The RMSE value was obtained as 2.55 for the prediction in ARIMA Model. However, LSTM shows better prediction with an RMSE value of 1.50 for seen data and 1.51 in degree Celsius for the unseen data in this univariate short term time series prediction.

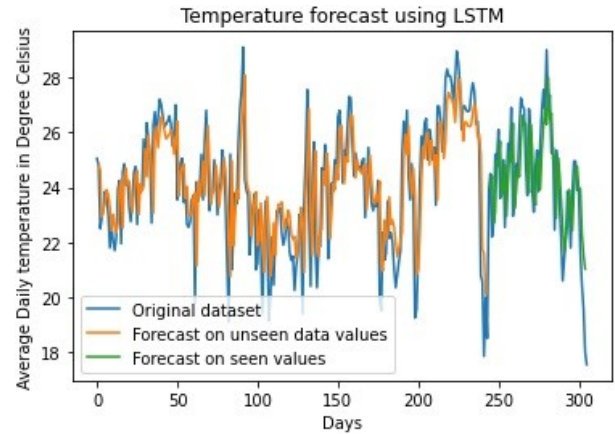


Fig. 4. Average Temperature predicted using LSTM

The LSTM depicted a much better prediction of seasonal temperature than the ARIMA model. LSTM works well for difficult sequential patterns of data. The forecast and accuracy can be improved on larger sets of training data in the LSTM model.

## CONCLUSION

Weather prediction has always been a challenging area for research because of its dynamic nature. Within the broad area of research, there were several methods that were used to predict weather to the utmost accuracy. More continuous data can be included to enhance the result on the ARIMA model. This paper has shown ARIMA model and LSTM for predicting average daily temperature at a station for a particular season. The future scope of this study is to forecast more parameters in relation to the environmental parameters as well which could be customized for particular crops. The lack of promised performance [17] of deep learning methods on univariate time series prediction is left open for research.

## REFERENCES

- [1] Cifuentes, J., Marulanda, G., Bello, A. and Reneses, J., 2020. Air temperature forecasting using machine learning techniques: a review. *Energies*, 13(16), p.4215.
- [2] Mohan, P., & Patil, K. K. (2018). "Deep learning-based weighted SOM to forecast weather and crop prediction for agriculture application". *Int. J. Intell. Eng. Syst*, 11, 167-176.
- [3] C. Lennard and G. Hegerl, "Relating changes in synoptic circulation to the surface rainfall response using self-organizing maps", *Climate Dynamics*, Vol.44, No.3-4, pp.861-879, 2015.
- [4] A. Geetha and G. M. Nasira, "Data mining for meteorological applications: Decision trees for modeling rainfall prediction", 2014 IEEE International Conference on Computational Intelligence and Computing Research, 2014, pp. 1-4, DOI: 10.1109/ICCIC.2014.7238481

- [5] Düben, Peter & Bauer, Peter. (2018). "Challenges and design choices for global weather and climate models based on machine learning". *Geoscientific Model Development*. 11. 3999-4009. 10.5194/gmd-11-3999-2018.
- [6] Q. Luan, X. Fang, C. Ye and Y. Liu, "An integrated service system for agricultural drought monitoring and forecasting and irrigation amount forecasting", 2015 23rd International Conference on Geoinformatics, 2015, pp. 1-7, doi:10.1109/GEOINFORMATICS.2015.7378617.
- [7] Le X-H, Ho HV, Lee G, Jung S. "Application of Long Short-Term Memory (LSTM) Neural Network for Flood Forecasting". *Water*. 2019; 11(7):1387. <https://doi.org/10.3390/w11071387>
- [8] Z. Basha, N. Bhavana, P. Bhavya and S. V, "Rainfall Prediction using Machine Learning & Deep Learning Techniques", 2020 International Conference on Electronics and Sustainable Communication Systems (ICESC), 2020, pp. 92-97, doi:10.1109/ICESC48915.2020.9155896.
- [9] Christodoulou, Christodoulos & Michaelides, Silas & Gabella, Marco & Pattichis, C., (2004). Prediction of rainfall rate based on weather radar measurements. 2. 1393 - 1396 vol.2. 10.1109/IJCNN.2004.1380153.
- [10] L. Naveev and H.S. Mohan, "Atmospheric Weather Prediction Using various machine learning Techniques: A Survey", 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC), 2019, pp. 422-428, doi: 10.1109/ICCMC.2019.8819643.
- [11] Beda Luitel, Gabriele Villarini and Gabriel A. Vecchi, "Verification of the skill of numerical weather prediction models in forecasting rainfall from U.S. landfalling tropical cyclones", *Journal of Hydrology*, 2016.
- [12] <https://www.itl.nist.gov/div898/handbook/pmc/section4/pmc442.html> (Accessed on 14<sup>th</sup> July 2021).
- [13] <https://towardsdatascience.com/mastering-time-series-analysis-in-python-8219047a0351>, (Accessed on 14<sup>th</sup> July 2021).
- [14] <https://people.duke.edu/~mrau/411arim.htm>, (Accessed on 14<sup>th</sup> July 2021).
- [15] Available online: <https://machinelearningmastery.com/time-series-prediction-lstm-recurrent-neural-networks-python-keras/>.
- [16] D.Cho, C. Yoo, J.Im, & D. H. Cha, "Comparative assessment of various machine learning - based bias correction methods for numerical weather prediction model forecasts of extreme air temperatures in urban areas", In *Earth and Space Science*, 7(4), e2019EA000740, 2020.
- [17] Available online: <https://www.datasciencecentral.com/profiles/blogs/arima-sarima-vs-lstm-with-ensemble-learning-insights-for-time-ser> (Accessed on 31<sup>st</sup> August 2021)