**Title page:**

# Comparative analysis of Auto Regressive Integrated Moving Average and Long Short-Term Memory in Prediction of Water Inundation Frequency with improved accuracy.

Kovi Sai Ganesh[1], S. Kalaiarasi [2]

Kovi Sai Ganesh [1]
Research Scholar,
Department of Computer Science and Engineering,
Saveetha School of Engineering,
Saveetha Institute of Medical and Technical Sciences,
Saveetha University, Chennai, Tamil Nadu, India, Pincode: 602105.
192124028.sse@saveetha.com


S. Kalaiarasi [2]
Project Guide, Corresponding Author,
Department of Computer Science and Engineering,
Saveetha School of Engineering,
Saveetha Institute of Medical and Technical Sciences,
Saveetha University, Chennai, Tamil Nadu, India, Pincode: 602105.
kalaiarasis.sse@saveetha.com

## ABSTRACT

**Aim**: The project aims to evaluate the accuracy of Long Short-Term Memory (LSTM) and ARIMA(Autoregressive Integrated Moving Average) in predicting the frequency of water inundations, which can be one of the most dangerous natural disasters that can seriously harm infrastructure and property. **Materials and Methods:** Two groups of algorithms are proposed. Long Short-Term Memory is compared with Autoregressive Integrated Moving Average. Both algorithms work in predicting the frequency of Water Inundation for accuracy. Accuracy is analyzed for Water Inundation frequency. LSTM's ability to handle long-term dependencies and ARIMA model can be useful for identifying long-term trends and seasonal patterns in water inundation frequency. The algorithm uses the properties of training data to create a model, which it then uses to estimate the value of new data. The Long Short-Term Memory and Autoregressive Integrated Moving Average of sample size (N=20) are used to recognize water inundation frequency. The significance value of the data set was predicted using SPSS with a G-power value above 80%. **Results:** Long Short-Term Memory has obtained an accuracy of 93.0095% which is comparatively higher than Autoregressive Integrated Moving Average with an accuracy of 87.5310%. There is a significant difference between the two groups with a significance value of 0.174(p<0.05).**Conclusion:** In Conclusion, these findings demonstrate the LSTM's greater predictive power in predicting the frequency of water inundation.For applications requiring the capture of long-term dependencies in the data, LSTM performs better than ARIMA.

**Keywords** : Autoregressive Integrated Moving Average, Deep Learning, Floods,Long Short-Term Memory(LSTM),Machine Learning, Neural Networks,Recurrent Neural Network, Water Inundation Frequency.

## INTRODUCTION

The temporary or permanent covering of land by water, known as "water inundation," is a major global concern. The frequency and intensity of floods are increasing due to a complex interaction between land-use alterations, urbanization, and climate change(Zhao et al. 2023). These events affect both inland and coastal regions, with the former suffering from flash floods and riverine overflows, and the latter from storm surges and tsunamis. Precisely anticipating the regularity of these flooding incidents has become essential for anticipating disasters, allocating resources wisely, and protecting susceptible populations(Hu et al. 2020).Now let's go into the world of deep learning, where techniques such as and Long Short-Term Memory (LSTM) have great potential for deciphering the complex temporal relationships seen in water inundation dynamics(Üneş et al. 2019). These advanced models have the ability to handle sequential data.

A thorough statistical analysis of flood occurrences highlights how serious the problem is. Figure 1 (Cho et al. 2022) provides a graphic depiction of the magnitude of the harm that occurs as a result of natural disasters. It emphasizes the increasing effects that have been shown throughout two noteworthy decades, 1980 to 1999 and 2000 to 2019. An

extensive scan of the Scopus database during our investigation of the research environment turned up an astounding 2,801 publications on Science Direct and more than 17,700 papers on Google Scholar over the previous five years. This increase in academic interest highlights the complexity of the problem at issue as well as the growing emphasis on water inundation prediction.Floods and other forms of flooding provide a complex threat that goes beyond the local ecosystem(Talbot et al. 2018).Besides being vulnerable to direct bodily harm, there is also a risk to public safety, disruptions to agricultural output, and possible effects on the integrity of the socioeconomic system as a whole(Engeland et al. 2018). This threat is so extensive that creative and practical remedies are required.The rising frequency and unpredictability of water inundation occurrences highlight the critical need for precise forecast models in order to minimize potential damages(Najibi and Devineni 2018).

This project's main goal is to carry out an exhaustive and efficient assessment of the effectiveness of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) models for water inundation frequency forecasting. Using cutting edge deep learning (DL) methods, our main goal is to increase these models' accuracy and reliability so that they may significantly improve early warning systems and catastrophe preparedness efforts(J. Chen et al. 2019). Our research aims to investigate the complex dynamics surrounding water inundation, taking into account important elements including land-use patterns, urbanization, and climate change, going beyond just improving prediction skills.

## MATERIALS AND METHODS

The investigation was conducted in the Saveetha Institute of Medical and Technical Sciences' Data Analytics lab, which has an extremely flexible setup that makes for thorough analyses and accurate findings. Ten volunteers in all were split into two groups for the study: Group 2 used the ARIMA (Autoregressive Integrated Moving Average) approach, whereas Group 1 used the LSTM (Long Short-Term Memory) method. Using G-Power, the study was able to maintain an 80 percent statistical power at the 0.8 (beta) power level and the 0.05 (alpha) significance level. Furthermore, a 95 percent confidence interval was maintained for the methodical calculation and assessment of the two groups' differences.

The present investigation's estimation of the frequency of flooding caused by water was based on the "Kerala Floods" dataset(Devakumar 2019), which was stored in CSV format. In order to improve frequency determination accuracy, the dataset included variables such as year, month, and annual rainfall. For training and testing, the dataset was thoroughly cleaned, with missing values removed, normalization applied, and averages or medians substituted for null values. Next, the preprocessed dataset with all of its features was added to the LSTM algorithm.

A particular dataset with values spanning a variety of data was used in a CSV file format to apply the strategies described in this study. Using the presented algorithms, a thorough analysis of all the records in the dataset was conducted, enabling a thorough inspection and comparison. An Intel dual-core processor with 8 GB of RAM was the hardware used for this endeavor. The Jupyter Notebook, Python, and MySQL database were used in the software configuration to create a strong foundation for running the algorithm and performing the comparative analysis.

**Long Short-Term Memory**

A particular kind of neural network called LSTM (Long Short-Term Memory) was created to get over the drawbacks of conventional RNNs in terms of identifying long-term dependencies. It includes memory cells and gating mechanisms that let it to retain information over extended sequences, making it well-suited for activities whereby situational and timing dependence are critical(Y.-C. Chen et al. 2021). When working with non-linear relationships in time-series data, LSTM is very good at collecting complex patterns.

**Autoregressive Integrated Moving Average**

An effective time series forecasting model that combines moving average, differencing, and autoregressive components is called ARIMA (AutoRegressive Integrated Moving Average). The moving average component models the association between an observation and residual errors, the autoregressive component records the relationship between the current and previous observations, and the differencing component changes the data to attain stationarity(Musarat et al. 2021). Because of its adaptability to handle a variety of temporal patterns, ARIMA is frequently used to forecast future values in time series data. Its parameters, which are written as (p, d, q), stand for the autoregressive, differencing, and moving average components, in that order. This allows for customisation according to the features of the particular dataset(Supatmi, Huo, and Sumitra 2019).

## Statistical Analysis

The latest version of IBM SPSS, version 26, was used for the analysis, which concentrated on accuracy values acquired from a sample size of 20 utilizing the ARIMA (Autoregressive Integrated Moving Average) and LSTM (Long Short-Term Memory) algorithms. Groups 1 and 2 were given the LSTM and ARIMA models, respectively, and the models were trained using the features of the dataset. The accuracy values for each epoch were then subjected to T-test computations, allowing for a comparison of the two algorithms' performances.

## RESULTS

The accuracy of the initial data table for both ARIMA (Autoregressive Integrated Moving Average) and Long Short-Term Memory (LSTM) is displayed in Table 3. The accuracy numbers are determined using the Long Short-Term Memory (LSTM) and the ARIMA (Autoregressive

Integrated Moving Average), both of which have sample sizes of 20. The accuracy value of the ARIMA (Autoregressive Integrated Moving Average) technique is 87.5530%, whereas the Long Short-Term Memory (LSTM) algorithm has a mean accuracy value of 93.0095%. Table 4, which displays these accuracy numbers. This suggests that the recommended algorithm is more accurate than the others. This study used the Long Short-Term Memory (LSTM) approach instead of the Autoregressive Integrated Moving Average (ARIMA).

Table 5 displays the "F" value as 2.010 and the "Sig" value as 0.164 with equal variances assumed. The 95% confidence intervals for the difference for equal variances assumed and not assumed are 7.18793 and 7.19368 for the higher case and 3.72507 and 3.71932 for the lower case (table 5). Plotting mean accuracy on the Y-axis and groups on the X-axis results in the bar graph displayed in Figure 2. The graph clearly shows that LSTM (Long Short-Term Memory) is more accurate than ARIMA (Autoregressive Integrated Moving Average).

## DISCUSSION

According to the research, LSTM (Long Short-Term Memory) appears to perform better than Regions ARIMA (Autoregressive Integrated Moving Average), with a significant value of 0.164 (Two-tailed, $p > 0.05$). While the mean accuracy of Regions using the ARIMA (Autoregressive Integrated Moving Average) classifier is 87.5530%, the mean accuracy analysis using the LSTM (Long Short-Term Memory) technique is 93.0095%. This suggests that when compared to ARIMA, LSTM performs better.

Both Long Short-Term Memory (LSTM) and ARIMA (Autoregressive Integrated Moving Average) are effective techniques, each having certain benefits and drawbacks. Applications such as time series forecasting and natural language processing profit immensely from LSTM's ability to accurately predict long-term dependencies in sequential data(Subha and Saudia 2023). Its gating system allows selective memory, which enhances its ability to recall complex temporal connections. Nevertheless, LSTM can be computationally expensive and challenging to train. However, because of its computing efficiency and ability to accommodate non-linear interactions, ARIMAis more suited for modeling stationary time series data(Azad et al. 2022). Its relatively simple training method makes it interesting for smaller or simpler datasets. But it's crucial to consider ARIMA's tendency toward overfitting and its limited ability to capture perpetual dependencies.

The size, diversity, and feature selection of the dataset, the ability to be general and unproven computational choices are among the potential limits of the study. To get over these limitations and progress the area, additional research should focus on expanding datasets that forecast increased frequency of flooding from various causes. Promising approaches for boosting prediction accuracy and the understanding of Water Inundation Frequency include focusing ensemble models, interpretability, and real-world applications in environmental science and crisis management(Dayal et al. 2019). Further investigation into alternative artificial intelligence

techniques, DL approaches, parameterization efficiency, and actual validating may yield enhanced models for forecasting.

## CONCLUSION

The Long Short-Term Memory (LSTM) algorithm's mean accuracy for forecasting the frequency of flooding was 93.0005%, whereas the ARIMA (Autoregressive Integrated Moving Average) algorithm's mean accuracy for the regions under study was 87.5530%. The effectiveness of the LSTM and ARIMA models in forecasting the frequency of water inundations was examined in this study. Both models fared well, according to the data, with LSTM outperforming ARIMA in terms of accuracy. Predicting the frequency of water inundations, which is impacted by past hydrological trends, was made easier by LSTM's capacity to identify long-term relationships in time series data. Although ARIMA performed well enough, it was not as good at capturing intricate temporal correlations. Because LSTM can handle sequential data better than other models, the investigation found that it is a better option for forecasting the frequency of water inundations.

# DECLARATION

## Conflict of Interests

This paper contains no declared conflicts of interest. To maintain our commitment to academic integrity and prevent any inadvertent participation with concerns related to academic dishonesty, we carefully checked that our work is original.

## Acknowledgment

## Authors Contribution

Author KSG actively participated in the data synthesis, analysis, and gathering. On the other hand, Author SKA made a significant contribution to the research proposal, verified the data, and offered critical criticism throughout the paper review procedure.

## Funding

# REFERENCES

Azad, Abdus Samad, Rajalingam Sokkalingam, Hanita Daud, Sajal Kumar Adhikary, Hifsa Khurshid, Siti Nur Athirah Mazlan, and Muhammad Babar Ali Rabbani. 2022. "Water Level Prediction through Hybrid SARIMA and ANN Models Based on Time Series Analysis: Red Hills Reservoir Case Study." *Sustainability: Science Practice and Policy* 14 (3): 1843.

Chen, Junfei, Qian Li, Huimin Wang, and Menghua Deng. 2019. "A Machine Learning Ensemble Approach Based on Random Forest and Radial Basis Function Neural Network for Risk Evaluation of Regional Flood Disaster: A Case Study of the Yangtze River Delta, China." *International Journal of Environmental Research and Public Health* 17 (1): 49.

Chen, Yue-Chao, Jia-Jia Gao, Zhao-Hui Bin, Jia-Zhong Qian, Rui-Liang Pei, and Hua Zhu. 2021. "Application Study of IFAS and LSTM Models on Runoff Simulation and Flood Prediction in the Tokachi River Basin." *Journal of Hydroinformatics* 23 (5): 1098–1111.

Cho, Minwoo, Changsu Kim, Kwanyoung Jung, and Hoekyung Jung. 2022. "Water Level Prediction Model Applying a Long Short-Term Memory (LSTM)–Gated Recurrent Unit (GRU) Method for Flood Prediction." *WATER* 14 (14): 2221.

Dayal, Deen, S. Swain, A. K. Gautam, S. S. Palmate, A. Pandey, and S. K. Mishra. 2019. "Development of ARIMA Model for Monthly Rainfall Forecasting over an Indian River Basin," May, 264–71.

Devakumar, K. P. 2019. "Kerala Floods 2018." https://www.kaggle.com/imdevskp/kerala-floods-2018.

Engeland, Kolbjørn, Donna Wilson, Péter Borsányi, Lars Roald, and Erik Holmqvist. 2018. "Use of Historical Data in Flood Frequency Analysis: A Case Study for Four Catchments in Norway." *Hydrology Research* 49 (2): 466–86.

Hu, Lanxin, Efthymios I. Nikolopoulos, Francesco Marra, and Emmanouil N. Anagnostou. 2020. "Sensitivity of Flood Frequency Analysis to Data Record, Statistical Model, and Parameter Estimation Methods: An Evaluation over the Contiguous United States." *Journal of Flood Risk Management* 13 (1): e12580.

Musarat, Muhammad Ali, Wesam Salah Alaloul, Muhammad Babar Ali Rabbani, Mujahid Ali, Muhammad Altaf, Roman Fediuk, Nikolai Vatin, et al. 2021. "Kabul River Flow Prediction Using Automated ARIMA Forecasting: A Machine Learning Approach." *Sustainability: Science Practice and Policy* 13 (19): 10720.

Najibi, Nasser, and Naresh Devineni. 2018. "Recent Trends in the Frequency and Duration of Global Floods." *Earth System Dynamics* 9 (2): 757–83.

Subha, J., and S. Saudia. 2023. "Robust Flood Prediction Approaches Using Exponential Smoothing and ARIMA Models." *Artificial Intelligence and Sustainable Computing*, 457–70.

Supatmi, Sri, Rongtao Huo, and Irfan Dwiguna Sumitra. 2019. "Implementation of Multiplicative Seasonal ARIMA Modeling and Flood Prediction Based on Long-Term Time Series Data in Indonesia." *Artificial Intelligence and Security*, 38–50.

Talbot, Ceara J., Elena M. Bennett, Kelsie Cassell, Daniel M. Hanes, Elizabeth C. Minor, Hans Paerl, Peter A. Raymond, et al. 2018. "The Impact of Flooding on Aquatic Ecosystem Services." *Biogeochemistry* 141 (3): 439–61.

Üneş, Fatih, Mustafa Demirci, Bestami Taşar, Yunus Ziya Kaya, and Hakan Varçin. 2019. "Estimating Dam Reservoir Level Fluctuations Using Data-Driven Techniques." *Polish*

*Journal of Environmental Studies* 28 (5): 3451–62.

Zhao, Chenchen, Chengshuai Liu, Wenzhong Li, Yehai Tang, Fan Yang, Yingying Xu, Liyu Quan, and Caihong Hu. 2023. "Simulation of Urban Flood Process Based on a Hybrid LSTM-SWMM Model." *Water Resources Management* 37 (13): 5171–87.

## TABLES AND FIGURES

**Table 1.** Representing LSTM (Long Short-Term Memory) algorithm  pseudocode**.**

| |
|---|
| **Input:** Kerala floods data |
| **Output:**  Prediction of Water Inundation Frequency |
| **Step 1: Data Collection and Preprocessing**<br><br>Input: Raw data from Kerala floods dataset.<br><br>Output: Processed data specific to Kerala floods.<br><br>Description: Collect and preprocess the Kerala floods dataset, handling missing values, and converting the data into a format suitable for time-series analysis. This may involve cleaning satellite imagery data, integrating weather records, and ensuring the data covers the relevant time periods of the floods. |
| **Step 2:  Feature Selection:**<br><br>Input: Processed Kerala floods data.<br><br>Output: Subset of relevant features for LSTM modeling.<br><br>Description: Identify features crucial for predicting flooding in Kerala. This may include rainfall data, water levels, land cover information, and historical flood occurrences. |
| **Step 3:Data Splitting:**<br><br>Input: Selected features from Kerala floods data.<br><br>Output: Training and testing datasets.<br><br>Description: Split the processed Kerala floods data into training and testing sets, considering the temporal aspect to ensure that the model generalizes well to future flood occurrences. |
| **Step 4:  Model Selection:**<br><br>Input: LSTM<br><br>Output: LSTM Model for Kerala floods.<br><br>Description: Choose the LSTM model for its capability to capture sequential dependencies in time-series data, which is relevant for predicting the frequency of floods in Kerala over time. |

**Step 5: Model Training:**

Input: Training data and selected features from Kerala floods.

Output: Trained LSTM model for Kerala floods.

Description: Train the LSTM model using the selected features and historical data specific to Kerala floods. Adjust LSTM parameters to effectively capture temporal patterns in the dataset.

**Step 6:Model Evaluation:**

Input: Trained LSTM model and testing data for Kerala floods.

Output: LSTM Model performance metrics for Kerala floods.

Description: Evaluate the LSTM model's performance on the testing dataset from the Kerala floods, considering factors such as prediction accuracy, sensitivity to flood events, and robustness in capturing temporal patterns.

**Step 7: Interpretability:**

Input: Trained LSTM model for Kerala floods.

Output: Insights into LSTM model behavior for Kerala floods.

Description: Explore the interpretability of the LSTM model specifically in the context of Kerala floods. Understand how the model utilizes historical information to predict flooding events in the region.

**Step 8: Deployment:**

Input: Trained LSTM model and selected features for Kerala floods.

Output: Deployed LSTM model for predicting floods in Kerala.

Description: Deploy the trained LSTM model for predicting floods in Kerala in a real-world setting. This may involve integrating the model into a system that can provide timely flood predictions for decision-makers in the region.

**Table 2**. Representing ARIMA(Autoregressive Integrated Moving Average) algorithm pseudocode.

| |
|---|
| **Input:** Kerala floods data |
| **Output:** Prediction of Water Inundation Frequency |
| **Step 1: Data Collection and Preprocessing**<br><br>Input: Raw data from Kerala floods dataset.<br><br>Output: Processed data specific to Kerala floods.<br><br>Description: Collect and preprocess the Kerala floods dataset, handling missing values, and converting the data into a format suitable for time-series analysis. This may involve cleaning satellite imagery data, integrating weather records, and ensuring the data covers the relevant time periods of the floods. |
| **Step 2:  Feature Selection:**<br><br>Input: Processed Kerala floods data.<br><br>Output: Subset of relevant features for LSTM modeling.<br><br>Description: Identify features crucial for predicting flooding in Kerala. This may include rainfall data, water levels, land cover information, and historical flood occurrences. |
| **Step 3:Data Splitting:**<br><br>Input: Selected features from Kerala floods data.<br><br>Output: Training and testing datasets.<br><br>Description: Split the processed Kerala floods data into training and testing sets, considering the temporal aspect to ensure that the model generalizes well to future flood occurrences. |
| **Step 4:  Model Selection:**<br><br>Input: ARIMA<br><br>Output: ARIMA Model.<br><br>Description: Choose the ARIMA model for its effectiveness in modeling time-series data with trends and seasonality. ARIMA is particularly suitable for forecasting tasks with historical patterns. |

**Step 5: Model Training:**

Input: Training data and selected features.

Output: Trained ARIMA model.

Description: Train the ARIMA model using the selected features and historical data. Adjust ARIMA parameters and orders (p, d, q) for auto-regression, differencing, and moving average to achieve improved accuracy.

**Step 6:Model Evaluation:**

Input: Trained ARIMA model and testing data.

Output: Performance metrics for the ARIMA Model.

Description: Evaluate the ARIMA model's performance on the testing dataset, considering metrics like Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). Assess its accuracy in predicting water inundation frequency.

**Step 7: Interpretability:**

Input: Trained ARIMA model.

Output: Insights into ARIMA model behavior.

Description: Explore the interpretability of the ARIMA model to understand how it leverages historical information and selected features to make predictions regarding water inundation frequency.

**Step 8: Deployment:**

Input: Trained ARIMA model and selected features.

Output: Deployed ARIMA model for predicting water inundation frequency.

Description: Deploy the trained ARIMA model in a real-world setting to provide efficient predictions of water inundation frequency. Ensure smooth integration for practical applications.

**Table 3.** Representing the Accuracy(93.005%) of LSTM & Accuracy(87.553%) of ARIMA Algorithms. The ARIMA Algorithm is less accurate than the LSTM Algorithm

| SAMPLE NO | LSTM(%) | ARIMA(%) |
|---|---|---|
| 1 | 96.81 | 92.31 |
| 2 | 96.25 | 91.97 |
| 3 | 95.89 | 90.65 |
| 4 | 95.74 | 90.21 |
| 5 | 94.78 | 90.11 |
| 6 | 94.68 | 89.98 |
| 7 | 94.45 | 89.12 |
| 8 | 94.32 | 88.91 |
| 9 | 92.56 | 88.37 |
| 10 | 92.55 | 88.26 |
| 11 | 92.41 | 88.14 |
| 12 | 91.37 | 87.91 |
| 13 | 92.31 | 86.37 |
| 14 | 91.63 | 86.17 |
| 15 | 91.49 | 86.14 |
| 16 | 91.34 | 85.52 |
| 17 | 91.14 | 84.34 |
| 18 | 90.85 | 82.22 |
| 19 | 89.67 | 82.19 |
| 20 | 88.95 | 82.17 |
| AVERAGE | 93.0095 | 87.553 |

**Table 4.**

## Group Statistics

Group Statistics Results-LSTM has a mean accuracy (93.0095%), St. Deviation (2.23460), whereas ARIMA has a mean accuracy (87.553%), St. Deviation (3.10432).

|  | Group | N | Mean | Std. Deviation | Std. Error Mean |
|---|---|---|---|---|---|
| Accuracy | LSTM | 20 | 93.0095 | 2.23460 | 0.49967 |
|  | ARIMA | 20 | 87.5530 | 3.10432 | 0.69415 |

**Table 5. Independent Samples Test**

The Equal variances assumed and the Equal variances not assumed are provided for the independent sample t-test for equality of means.

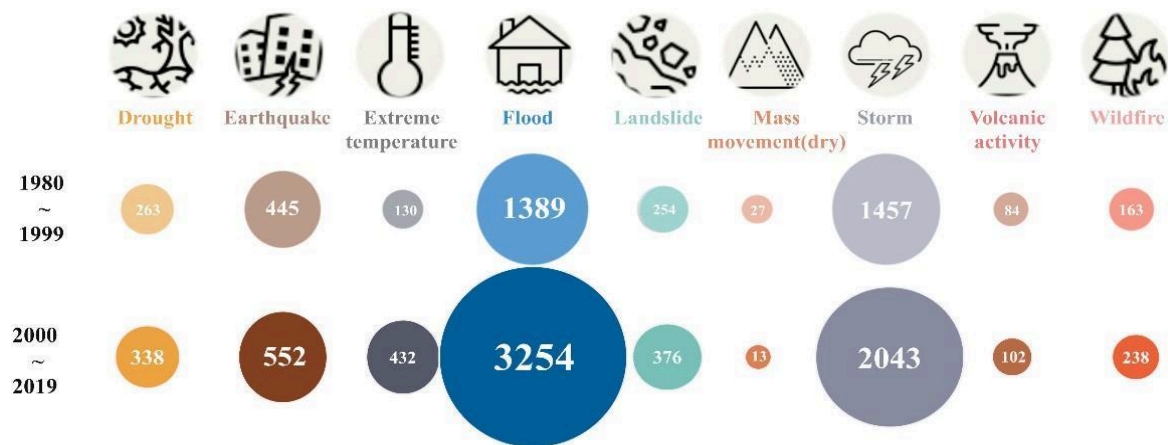|  |  | Levene's Test for Equality of Variances | | t-test for Equality of Means | | | | | | 95% Confidence Interval of the Difference | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | F | Sig. | t | df | Sig. (2-tailed) | Mean Difference | Std. Error Difference | Lower | Upper |
| Accuracy | Equal variances assumed | 2.010 | 0.164 | 6.380 | 38 | 0.000 | 5.45650 | 0.85529 | 3.72507 | 7.18793 |
|  | Equal variances not assumed |  |  | 6.380 | 34.523 | 0.000 | 5.45650 | 0.85529 | 3.71932 | 7.19368 |

**Figure 1.** Compares the categories of natural hazards events from 1980 to 1999 to 2000 to 2019 [5].
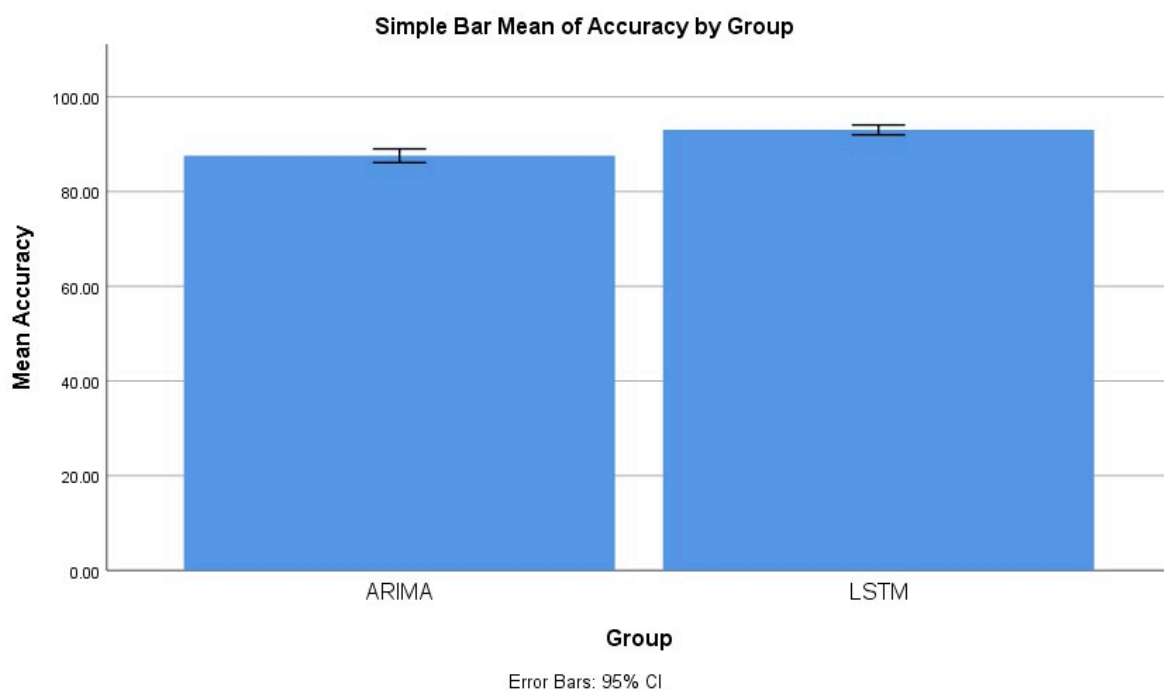


**Figure 2.** Shows a mean accuracy comparison between the ARIMA and LSTM algorithms. Compared to the ARIMA algorithm's mean accuracy of 87.553, the LSTM's mean accuracy of 93.0095 is superior. X-axis: ARIMA (Autoregressive Integrated Moving Average) vs. LSTM (Long Short-Term Memory). Y-axis: Mean accuracy.