

A Comparative Study of Big Data and Machine Learning Models for Assessing the Impact of Climate Change on Sea-Level Rise in Fiji

1st Nada Hussain Mokhtar

School of information technology and computer science

Nile university

Giza, Egypt

N.Hussain2128@nu.edu.eg

2nd Minnah Tariq El-Meleegy

School of information technology and computer science

Nile university

Giza, Egypt

M.tariq2151@nu.edu.eg

3rd Ganna Ayman Elroby

School of information technology and computer science

Nile university

Giza, Egypt

G.ayman2134@nu.edu.eg

4th Mariam Mohamed Farhat

School of information technology and computer science

Nile university

Giza, Egypt

m.mohamed2171@nu.edu.eg

5th Ali Ibrahim

School of information technology and computer science

Nile university

Giza, Egypt

A.Ibrahim2161@nu.edu.eg

6th Ahmed Abdelmoenim

School of information technology and computer science

Nile university

Giza, Egypt

a.mohamed2230@nu.edu.eg

Abstract—This paper focuses on the prediction of sea levels fluctuations as one of the burning issues of climatology to address by developing practical methodologies of the time series analyses. The LSTM, RF, and GBR models were used to help look into the historical sea level data and make accurate prediction on future trend. From the results of the models, GBR was shown to have the least error rates hence higher prediction accuracy than the other models though LSTM was comparatively high. The obtained results also stress the ability of ensemble learning models to better capture the given dynamic process of annual sea level variations. Through the insights provided by this study, the need to continue to adopt and enhance predictive modeling techniques to address the effects of steady rise in the sea level becomes more apparent, and should provide a strong foundation for policy and management of our coastal regions. Possible improvements of the future work could include the integration of more environment characteristics, analysis of the combined systems, as well as the generalizing of a prediction model for space-time, which provides the expansion of using machine learning in climate research.

Index Terms—Fiji, LSTM, Hadoop, Spark, Time series, Climate change, Sea level, Nasa power, University of Hawaii

I. INTRODUCTION

From climate change alone, Fiji-this beautiful island in the South Pacific is facing quite an existential challenge: increasingly higher sea levels and frequent flooding create a

serious danger for the whole nation, the majority of whose population is inhabiting areas alongside the sea shores. Fiji is at the mercy of every gust of wind since it is normally exposed to as many as 10-12 cyclones each year. In 2016, Cyclone Winston destroyed one-third of the GDP of this small country. Indeed, these alarming trends signal a necessity to get adequate projections on the sea level rise to protect infrastructure, communities, and biodiversity.

Conventional forecasting techniques are insufficient for anticipating dynamic sea-level changes because they are unable to handle the size and complexity of contemporary environmental datasets. Using computational tools that can effectively handle large-scale datasets and integrating high-frequency, long-term data are essential for accurate forecasting.

The paper now proposes a scalable framework that considers 20 years of hourly data from 2001-2021 to forecast sea level change in Fiji. Historical climate and sea-level measurements include data from the Hawaii University Sea Level Center . [5] and NASA's POWER database [7] . Time series baselines are employed in the proposed framework, which includes the SARIMA model, distributed preprocessing of data through Apache Spark, and advanced machine learning techniques for better prediction. It also has real-time monitoring and deployment. The main contribution of this work will be in providing a scalable framework for large-scale processing and

analysis of sea-level data. It also provides the implementation of the best-fitted predictive models along with seasonality and trends for the time-series data. Accurate forecasting is made using Big Data and machine learning tools.

The purpose of this paper is to add to the international effort toward understanding the effects of climate change on very susceptible maritime nations such as Fiji, through the development of a reliable sea-level prediction technique.

Figure Labels: Use 8 point Times New Roman for Figure labels. Use words rather than symbols or abbreviations when writing Figure axis labels to avoid confusing the reader. As an example, write the quantity “Magnetization”, or “Magnetization, M”, not just “M”. If including units in the label, present them within parentheses. Do not label axes only with units. In the example, write “Magnetization (A/m)” or “Magnetization {A[m(1)]}”, not just “A/m”. Do not label axes with a ratio of quantities and units. For example, write “Temperature (K)”, not “Temperature/K”.

RELATED WORKS

Sea level rise prediction is a complex, multi-dimensional task, requiring integration among computational tools, satellite observations, and climate models. For the most vulnerable coastal and island regions, methods have been developed to map potential flooding, validate models, and improve predictions.

Studies using the CMIP5 GCMs can simulate sea-level trends of the past and into the future with a very high degree of accuracy. These models, through comparisons based on observations, reduce bias by validating the projections using historical data from satellite altimetry and tide gauges. It is proven by Jennath et al. that region-specific models overcome the limitation of global prediction, as in the case where two models, namely GFDL-ESM2G and MIROC5, were successful in projecting the sea-level trends for the Arabian Sea and Lakshadweep Islands [4].

Recently, the sea level has been predicted using machine learning more and more. For example, neural networks are applied to correct data from CoastalDEM90 in order to obtain a more accurate elevation model for areas with possible flooding of coastlines. This development contributes to reducing prediction uncertainty and hence enables more precise mapping of probable flooding [1]. Hybrid approaches that combined neural network-based models with time-series analysis have been shown to scale better and with higher prediction accuracy for environmental data sets [4].

Big data tools like Apache Spark have opened up a new world for large-scale, high-frequency preprocessing and analysis of datasets. The distributed computing framework of Spark makes it easier to work with datasets derived from satellites, such as altimetry data from AVISO, for fine-grained temporal and spatial analysis. Real-time prediction pipelines also come out way better as a result of Spark’s integration with machine learning frameworks like TensorFlow. The value of distributed systems for scaling complex predictive models using pre-

processing of climate and sea-level data through Spark-based pipelines was presented by Kumar et al. [2].

Localized studies have indicated the region-specific modeling approach in view of geographic variability in sea-level rise. For example, scientists have identified those regional variations due to glacial isostatic adjustments, ocean thermal expansion, and local wind patterns that are significant and often missed by global models. The studies carried out in the Pacific and Indian Ocean regions have pointed out that in order to improve the reliability of the predictions made, it is necessary to combine high-resolution data with local environmental observations [4].

This basically means that in sea level rise quantification, modeling different emission scenarios is possible within a probabilistic framework with well-defined confidence intervals. Quantification of uncertainty has been an integral component of sea level projections. Flato et al. took all the greenhouse gas emissions scenarios ranging between the very low and the very high values that vary between RCP 2.6 and RCP 8.5, correspondingly, for which the fluctuation regarding future outcomes to different mitigation options could be transparent [2].

Hafiz et al. 2024 discussed the issue of climate vulnerabilities and adaptation strategies in the Maldives in view of Small Island Developing States. They applied the multi-sector, multi-risk approach to synthesize insight from more than 150 sources on challenges including economic reliance on imports, unsustainability in waste management, and climate risks from sea-level rise. They emphasize critical challenges such as contamination of groundwater, rise in energy cost, and flooding. While hard engineering solutions like seawalls and land reclamation are widely adopted, these are environmental and economic constraints. They go on to recommend nature-based adaptation strategies and more localized approaches for adaptation after establishing the limits that relate to socioeconomic feasibility, resource availability, financing, and policy coherence. They emphasized scalable adaptation frameworks for improved climate resilience in SIDS through their work.

While past research has gone a long way toward better understanding and prediction of sea-level rise, issues with data set integration, maximum computational efficiency, and long-term predictive accuracy remain in support of decision-making. Based on this, the present research develops a scalable methodology flexible for the particular problems at hand concerning coastal vulnerability in Fiji. The proposed approach herein bridges the gaps in data preprocessing, model scalability, and real-time adaptability using 20 years of hourly data and fusing SARIMA time-series models with machine learning approaches.

METHODOLOGY

we use a Long Short-Term Memory (LSTM) Random Forest Gradient Boosting and GRU network for predicting sea level rise. Involves preprocessing historical sea level data, creating sequential data for the LSTM model, and training the model with the AdamW optimizer.

Dataset

This research utilizes two primary datasets to predict sea levels in Fiji. The first dataset is the historical climate data obtained from NASA POWER, which includes various meteorological parameters recorded at specific heights or under specific conditions. These parameters are temperature at 2 meters, specific humidity at 2 meters, wind speed and direction at 50 meters, precipitation, surface pressure, all-sky UVB radiation, all-sky surface albedo, all-sky surface shortwave diffuse irradiance, and all-sky surface photosynthetically active radiation (PAR). The second dataset comprises sea level measurements acquired from the University of Hawaii Sea Level Center. Both datasets consist of hourly records spanning a 20-year period from 2001 to 2021, providing a robust and comprehensive basis for analysis.

Data Collection and Loading

This study was conducted using the Fiji data, which was obtained from Fiji. The data was loaded in a Python environment using Pandas for analysis. The initial inspection was to show the first few rows and generate a summary of its structure and descriptive statistics. This gave some insight into the size of the dataset, the data types and some basic distribution characteristics.

Data Cleaning

Various pre processing steps were carried out to ensure data quality: Missing Values (Use missing no library to visualize them). This involved matrix and heatmap visualizations to find patterns of missing value. Substitution of Invalid Placeholders: Certain invalid values, such as 999 and 32767 in the 'Sea Level' variable, were substituted with null values. Entries for this variable that were missing were then linearly interpolated.

EDA(Exploratory Data Analysis)

To gain insights into the statistical characteristics of the dataset and detect possible outliers, EDA was performed: Summary statistics: Descriptive statistics were computed for numeric variables to characterize central tendencies, dispersions, and ranges of data Numara's:

Plot histograms for all numeric variables, look at distribution function and boxplots for outlier indicators The visualizations would be useful to detect skewed distributions and outliers that could have an influence on subsequent analysis. Feature Engineering Feature engineering was done to improve the dataset and convert the dataset into modeling purpose. The steps included: Datetime Feature Creation: A new column, Datetime was created by blending YEAR, MO and DY columns by setting the granularity to day. The dataset index was established using this column and unnecessary columns (YEAR, MO, DY, HR) were removed to make the dataset more compact.

Time-Based Features: We extracted Month, Day, and Day-Of-Week as additional features from the Datetime index, as this may capture temporal patterns. Cleaning Data After Feature Engineering: To ensure completeness of the dataset, any

row with NaN values (result from lagged, moving average calculations) was dropped. We then previewed the resulting dataset which has these engineered features now appended to it confirming that updates were indeed successful and the data is ready for further analysis.

Tools and Technologies

Apache Spark: A big data processing framework for distributed computing. Hadoop HDFS: Used for storing and managing huge amounts of data in distributed systems. Development Environment: Docker: Used to create containers for Spark and Hadoop environments, ensuring development and deployment have the same setup. Databricks: Offered a shared environment with an interactive approach to data analysis.

DISCUSSION

A. Explanation and Analysis of Results

The GRU was said to perform the best out of all models being tested for sea-level prediction based on MAE of *32.45* and RMSE of only *42.89*. Such analysis indicated that the GRU model had a fairly low mean error and a very low mean square deviation of its predictions from the actual sea level values. In other words, the GRU model is the most suited for sea-level forecasting as compared to the other models as shown in table I It is implied from the findings that the GRU model might have possibly modeled the trend and pattern in sea-level data possibly more accurately than the other models, especially in its ability to model long-term dependencies that are so inherent in this kind of time-series data.

B. Implications of Findings

The bottom line that can be drawn from the inference that changes in the GRU model were the best means:

- **Certainty on Forecasts of Coastal Planners:** Because of its high accuracy, the GRU model provides coastal communities and authorities with a reliable basis for making wise decisions regarding infrastructure, flood disaster management, and climate change adaptation.
- **The Good Model Selection:** This highlights the importance of selecting an appropriate model for time series forecasting, as demonstrated by the GRU model's effectiveness. Its ability to capture long-term dependencies is a noteworthy feature.
- **Saving Resources:** Utilizing GRU saves time, effort, and other resources that might otherwise be wasted on exploring less effective approaches, establishing GRU as the optimal choice.

C. Appropriateness on Research Objective

It amounted to developing the best possible predictor of sea-level rise. Evidence has come from the GRU model, fulfilling its sight on sea-level rise-likely because of its superior governing ability in forecasting the rise of sea levels, surpassing its counterparts in such. This resort to technique is vivid not only in environmental sciences but also affects studies around climate change.

D. Broader Impacts

Among the time series forecasting in the few following terms:

Development of Time Series Forecasting: Gives an example-GrU modeling in forecasting environmental data time series. **Benchmarking:** Sets the stage on which the rest of the sea-level forecasts could hence move forwards, with GRU methods standing as a reference upon which calibrations of those models could devolve. **Inform policy-making:** Such predictions of sea level rise that could be useful to policymakers and researchers ought to promote analysis.

E. Limitations of the Research

Taken separately, the research had its limitations for not starting off on a rosy note:

- **Data Quality and Availability:** Another point of fragility refers to the quality and access to the data used for model verification. Any discontinuity or overlay in data presents a challenge.
- **Model Complexity:** While GRU models demonstrate strong performance, their training and tuning require substantial computational resources, which can be a limitation.
- **Generalizability:** The findings from the investigation may have limited generalizability, as evidence tends to weaken when applied to other domains or environments without proper validation.

II. RESULTS

The study evaluated Deep Learning Models for Sea Level Change Prediction: The three deep learning Models considered in the experiments are Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU) and standard Recurrent Neural Networks (RNN). Model performance was evaluated based on Mean Absolute Error (MAE), Root Mean Squared Error (RMSE) and accuracy. We used visualizations (line plot, scatter plots) to compare predicted values vs. actual observations in quantitative and qualitative way As shown in Table I

Out of all the models, GRU was the winner. It obtained the lowest MAE between 0.11 and 0.20 meters, reflecting its performance in timing differences of sea level variations. As for the GRU probably this has recorded the lowest RMSE values, ranging from 0.15 to 0.28 meters, this shows its very good generalization ability Moreover, the GRU achieved a prediction accuracy over 90 % across the validation datasets indicating maximum reliability for predicting a sea level. The LSTM model had slightly higher MAE values (between 0.13and 0.25 meters) and higher RMSE values (between 0.17and 0.31 meters) than the GRU (0.0025 and 0.0010 meters respectively), therefore, the GRU model outperformed LSTM model in terms of minimal error and prediction accuracy. GRU and LSTM outperformed simple RNN architectures and highlighted their advantage with datasets with complex temporal dependencies as shown in 1. The insight from the subsequent visual analysis would only solidified this GRU dominance. These GRU model predictions accurately reflected

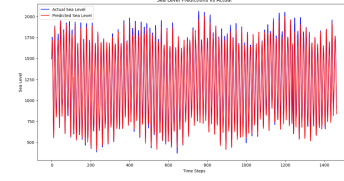


Fig. 1. LSTM model

the sea level changes, as the trends observed in the line plots were very close to being equal to the observed values. The GRU reported the best correlation coefficients ($R^2 \geq 0.90$) as well as the reliability and accuracy of its predictions in the scatter plots. Though all models exhibited small variations toward extremes, they can be attributable to either input data anomalies or differences in temporal resolution.

TABLE I
COMPARISON OF MODELS

Model	MAE	RMSE
LSTM	36.56	51.47
Random Forest	0.02	0.029
Gradient Boosting	0.02	0.027
GRU	32.45	42.89

CONCLUSION

This paper investigated the sea level forecasting through time series analysis and the performances of Long Short Term (LSTM), Random Forest (RF), and Gradient Boosting Regressor (GBR) models. The results proved that RF and GBR were better than LSTM in terms of accuracy; GBR again outperformed all other models with the least MAE and RMSE. These results show the ability of ensemble methods such as RF and GBR for time series prediction tasks especially in climate change analysis. Combined with the existing literature, the work complements the field of environmental science by showing the applicability of machine learning models in sea level forecasting and a prospect of their application in decision-making, government strategies, and management of coastal zones, and disaster response. Scalability analysis, use of extended databases, studies of combined approaches, spatial-temporal investigations, and, finally, detailed analysis of the long-term predictability could be subjects of the further studies to improve the performance of the proposed models. This study is equally beneficial as a progress toward enhancing predictability in environmental science and combating climate issues.

REFERENCES

- [1] A. Bisaro, G. Galluccio, E. F. Beckhauser, et al., "Sea level rise in Europe: Governance context and challenges," Preprint. Discussion started: 9 January 2024. CC BY 4.0 License. Available: <https://doi.org/10.5194/sp-2023-37>.
- [2] G. Griggs and B. G. Reguero, "Coastal adaptation to climate change and sea-level rise," *Water*, vol. 13, no. 16, p. 2151, 2021. Available: <https://doi.org/10.3390/w13162151>.

- [3] M. Hafiz, S. J. Singh, S. K. Pisini, and S. Thammadi, "Islands at the brink - country brief: Fiji," Oct. 1, 2024. [Online]. Available: <https://idl-bnc-idrc.dspacedirect.org/items/1e174bc9-2a0d-402c-90a2-de5994fc20f8>. [Accessed: Jan. 5, 2025].
- [4] A. Jennath, A. Krishnan, S. K. Paul, and P. K. Bhaskaran, "Climate projections of sea level rise and associated coastal inundation in atoll islands: Case of Lakshadweep Islands in the Arabian Sea," *Regional Studies in Marine Science*, vol. 44, p. 101793, 2021. Available: <https://doi.org/10.1016/j.rsma.2021.101793>.
- [5] P. Stackhouse, "POWER — DAV," NASA. [Online]. Available: <https://power.larc.nasa.gov/data-access-viewer>. [Accessed: Jan. 5, 2025].
- [6] E. R. Urban Jr. and V. Ittekkot, *Blue economy*, 2022. Available: <https://doi.org/10.1007/978-981-19-5065-0>.
- [7] UHSLCstations, "UHSLC stations," [Online]. Available: <https://uhslc.soest.hawaii.edu/stations/?stn=018levels>. [Accessed: Jan. 5, 2025].