# AI for Earthquake Prediction: A KNN-Based Analysis of Global Seismic Data

**Allacheruvu Brahmaiah[1], Kattika Koteswara Rao[2]**

*[1,2]Assistant Professor, Department of Computer Science and Engineering, R K College of Engineering*

Vijayawada, India

allacheruvubrahmaiah07@gmail.com[1]; kotesh.sasi@gmail.com[2]

*Abstract -* **Earthquakes pose a significant threat to human life and infrastructure, necessitating reliable predictive models for risk mitigation and disaster preparedness. This study leverages a global earthquake dataset (2020–2024) to analyse historical seismic patterns and predict future occurrences using the K-Nearest Neighbors (KNN) algorithm. The predictive model generates earthquake forecasts for the next 1 month, 6 months, and 1 year, with affected regions categorized by continent. To evaluate model performance, key metrics such as Mean Squared Error (MSE: 0.198), Root Mean Squared Error (RMSE: 0.445), Accuracy Percentage (94.57%), and Mean Absolute Percentage Error (MAPE: 1.16%) are calculated. The results indicate that machine learning techniques, particularly KNN, can effectively model earthquake occurrence probabilities based on historical data. This research underscores the potential of AI-driven approaches in seismic prediction and risk assessment. While the proposed model demonstrates promising predictive capabilities, further enhancements—such as integrating deep learning techniques and real-time geophysical data—could improve accuracy and practical applicability in early warning systems.**

*Keywords* – **Earthquake Prediction, Seismic Data Analysis, Risk Mitigation**

## 1. INTRODUCTION

Earthquakes are one of the most devastating natural disasters, causing significant loss of life, economic damage, and infrastructure destruction worldwide. According to the United States Geological Survey (USGS), thousands of earthquakes occur globally each year, with some leading to catastrophic consequences depending on their magnitude and location [1]. Due to the unpredictable nature of seismic activity, developing reliable models for earthquake prediction is a crucial challenge in geophysics and disaster management. Traditional earthquake prediction methods rely on seismological and geological observations, but recent advancements in artificial intelligence (AI) and machine learning (ML) offer new possibilities for analysing seismic patterns and forecasting future occurrences [2], [3]. Machine learning techniques have demonstrated their effectiveness in various predictive applications, including weather forecasting, financial modelling, and medical diagnosis. In the field of seismology, ML-based models can analyse large datasets, detect hidden patterns, and improve the accuracy of earthquake predictions. Among these techniques, the K-Nearest Neighbors (KNN) algorithm is widely used for classification and regression tasks due to its simplicity, adaptability, and ability to handle non-linear relationships in data [4]. This study applies the KNN algorithm to a global earthquake dataset from 2020 to 2024 to predict earthquake occurrences in different regions. The model forecasts seismic events for the next 1 month, 6 months, and 1 year, with an evaluation of affected locations categorized by continent.

To assess the performance and reliability of the proposed model, key metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Accuracy Percentage, and Mean Absolute Percentage Error (MAPE) are calculated. The results indicate that the ML-based approach can provide

valuable insights into earthquake prediction and serve as a foundation for further enhancements. While the study demonstrates promising results, there is still room for improvement by integrating deep learning models and real-time geophysical data to enhance prediction accuracy.

The rest of the paper is organized as follows: Section 2 reviews related works in earthquake prediction using AI and ML techniques. Section 3 describes the dataset, methodology, and implementation of the KNN model. Section 4 presents experimental results and performance evaluation. Section 5 discusses potential improvements and limitations, and Section 6 concludes the study with future research directions.

## 2. RELATED WORK

Earthquake prediction has been an area of extensive research, with various statistical and artificial intelligence (AI) based approaches developed to forecast seismic events. Traditional methods rely on geophysical indicators such as seismic wave propagation, fault line activity, and ground deformation measurements [5]. However, due to the complexity and non-linearity of earthquake occurrence patterns, machine learning (ML) and deep learning (DL) techniques have been increasingly explored to enhance predictive accuracy.Several studies have employed supervised learning models, including Decision Trees, Support Vector Machines (SVMs), and Random Forest, for earthquake classification and prediction. For instance, Raju et al. [6] applied Decision Tree and SVM models to a seismic dataset to classify earthquakes based on magnitude and depth. Their findings indicated that SVM outperformed traditional statistical models in predicting earthquake occurrences with improved accuracy. Similarly, Sajid et al. [7] explored Random Forest and ensemble learning techniques for earthquake prediction and reported significant improvements in classification performance.

Deep learning methods have also been applied in recent years. Mousavi et al. [8] proposed a Convolutional Neural Network (CNN) based approach for detecting seismic signals and distinguishing earthquakes from background noise. Their model demonstrated high sensitivity in real-time seismic event detection. Additionally, Wang et al. [9] developed a Recurrent Neural Network (RNN) model that utilized temporal dependencies in seismic data to improve earthquake prediction accuracy over extended time frames. These studies highlight the potential of AI-driven models in seismic analysis.Despite the success of deep learning models, their implementation often requires large labelled datasets and extensive computational resources. As a result, simpler ML models such as K-Nearest Neighbors (KNN) have gained attention due to their ease of implementation and effectiveness in classifying seismic events. Research by Sharma and Gupta [10] demonstrated that KNN could successfully categorize historical earthquake data based on seismic parameters, achieving competitive performance compared to more complex models.

Building on these previous works, this study applies the KNN algorithm to analyse historical earthquake data (2020–2024) and predict future occurrences. The model forecasts seismic activity for different timeframes—1 month, 6 months, and 1 year—while also categorizing affected locations by continent. By leveraging KNN's ability to classify data based on similarity measures, this research aims to provide a simple yet effective approach to earthquake prediction, complementing existing AI-based methodologies.

## 3. METHODOLOGY

This section describes the dataset used, preprocessing steps, model implementation, prediction intervals, and evaluation metrics for assessing the performance of the earthquake prediction model. Figure 1: Workflow of Earthquake Prediction Using KNN Algorithm.
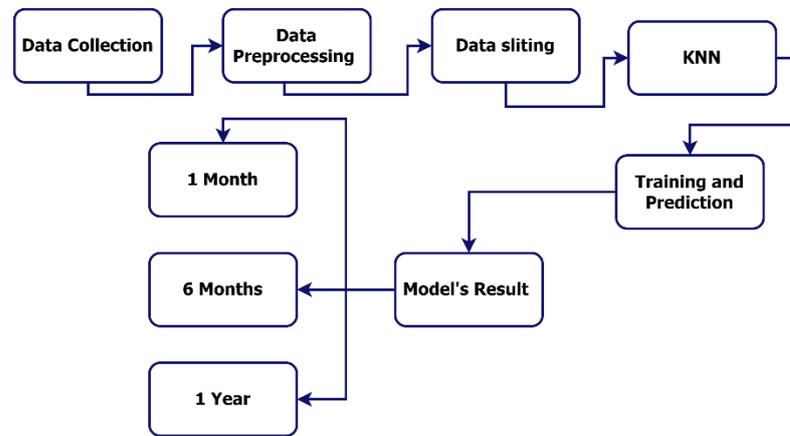
Figure 1: Workflow of Earthquake Prediction Using KNN Algorithm

### 3.1 Dataset Description

The dataset used in this study consists of global earthquake records from 2020 to 2024, obtained from publicly available seismological databases such as the United States Geological Survey (USGS) [11] and the European-Mediterranean Seismological Centre (EMSC) [12]. The dataset includes essential attributes that characterize seismic events:

- **Time:** Timestamp of the earthquake occurrence.
- **Latitude & Longitude:** Geographic coordinates of the earthquake epicentre.
- **Depth:** The depth of the earthquake measured in kilometres.
- **Magnitude:** The Richter scale measurement of earthquake intensity.
- **Location:** The affected region or country where the earthquake was recorded.

These attributes provide valuable information for analysing earthquake patterns and training predictive models.

### 3.2 Data Preprocessing

Data preprocessing is a crucial step to ensure the quality and consistency of the dataset before model training. The following steps were performed:

- **Handling Missing Values:** Missing or incomplete records were removed or imputed using statistical techniques to maintain data integrity [13].
- **Normalization:** Numerical attributes such as magnitude and depth were normalized using min-max scaling to improve the performance of the K-Nearest Neighbors (KNN) algorithm [14].
- **Categorization by Continents:** Earthquake occurrences were classified by continents (Asia, North America, South America, Europe, Africa, and Oceania) using geographic mapping techniques [15].
- **Feature Engineering:** Additional features, such as time intervals between consecutive earthquakes and clustering of epicentres, were derived to enhance model accuracy.

### 3.3 Model Implementation

The K-Nearest Neighbors (KNN) algorithm was selected for earthquake prediction due to its simplicity and effectiveness in handling classification problems. The steps for model implementation include:

- **Training Data Preparation:** Historical earthquake records (2020–2023) were used as training data, while records from 2024 were reserved for testing.

- **Feature Selection:** The model utilized latitude, longitude, depth, magnitude, and time as key predictors.
- **Distance Metric**: Euclidean distance was used to measure the similarity between past and future earthquakes.
- **K-Value Optimization:** Various values of K were tested, and the optimal K value was selected based on cross-validation results.
- **Prediction Execution:** The trained model predicted earthquakes for different timeframes (1 month, 6 months, and 1 year).
- KNN was implemented using Python's sci-kit-learn library, which provides an efficient framework for machine learning models [16].

### 3.4 Prediction Intervals

The model was designed to predict earthquake occurrences within the following time intervals:

1. **Short-term (1 month):** Useful for immediate disaster preparedness.
2. **Medium-term (6 months):** Helps in regional planning and resource allocation.
3. **Long-term (1 year):** Provides insights into potential seismic trends for future risk management.

These prediction intervals allow for different levels of preparedness, aiding governments, disaster management agencies, and researchers in decision-making.

### 3.5 Evaluation Metrics

To assess the performance of the KNN-based earthquake prediction model, the following evaluation metrics were computed:

- **Mean Squared Error (MSE):** Measures the average squared difference between actual and predicted earthquake magnitudes [17].

$$\mathbf{MSE} = \frac{1}{n}\sum(y_i - \hat{y}_i)^2$$

- **Root Mean Squared Error (RMSE):** Provides a more interpretable error value by taking the square root of MSE.

$$\boldsymbol{RMSE} = \sqrt{MSE}$$

- **Accuracy Percentage:** Evaluates the correctness of earthquake occurrence classification.

$$\mathbf{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} \times 100$$

- **Percentage Error:** used to measure the accuracy of a prediction or estimation compared to the actual value. It is calculated as:

$$\mathbf{Percentage\ Error} = \frac{\text{Actual Values} - \text{Predicted Values}}{\text{Actual values}} \times 100$$

These metrics provide a comprehensive assessment of the model's predictive reliability.

### 3.6 Summary of Methodology

The proposed methodology involves acquiring a high-quality earthquake dataset, preprocessing it for consistency, implementing the KNN algorithm for prediction, and evaluating its performance using standard ML metrics. The results of this approach are analysed in the subsequent sections to determine its feasibility for real-world applications.

### 4. RESULTS AND DISCUSSION

This section presents the results obtained from the earthquake prediction model using the K-Nearest Neighbors (KNN) algorithm. The discussion covers three key aspects: (1) historical vs. predicted trends, (2) continental analysis, and (3) performance analysis using evaluation metrics.

### 4.1 Historical vs. Predicted Trends

To assess the effectiveness of the KNN-based earthq uake prediction model, historical earthquake data (2020–2023) was compared with predicted earthquake occurrences for 2024. The actual frequency of earthquakes was plotted against the predicted values for different time intervals (1 month, 6 months, and 1 year).

Figure 2 illustrates a world map of earthquake locations, comparing historical (red) and predicted (blue) earthquakes using a KNN regression model. The predictions align with known seismic zones like the Pacific Ring of Fire and the Himalayas. The model captures earthquake-prone regions, aiding in understanding future seismic trends and potential risk areas.
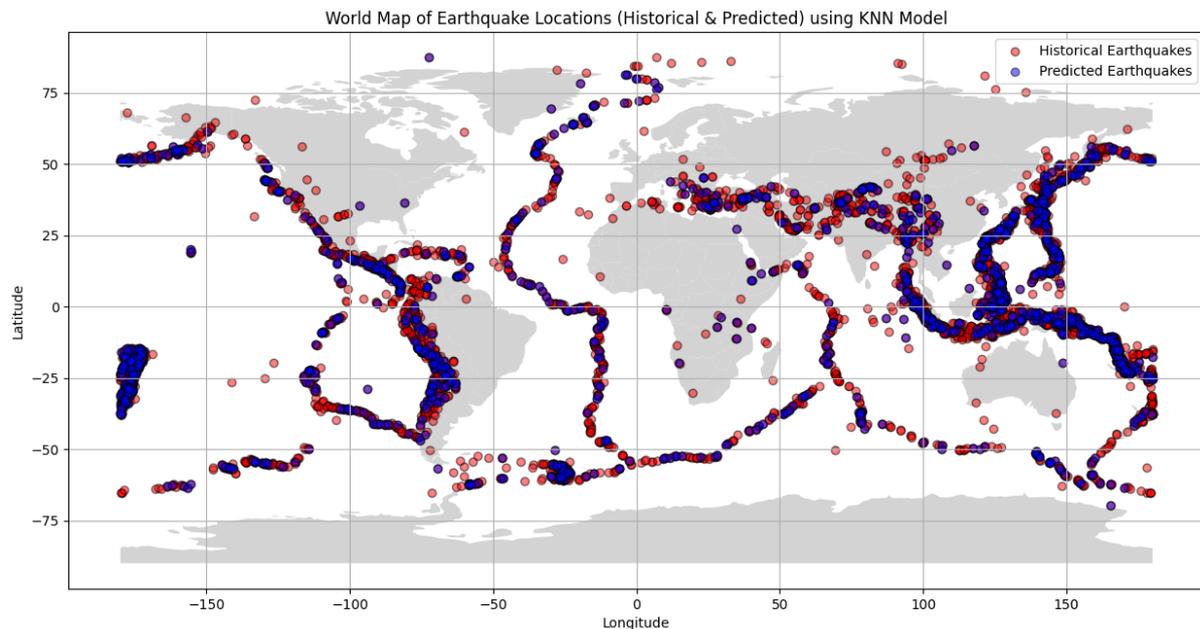


Fig. 2. World Map of Earthquake Locations (Historical & Predicted) using KNN Model

Figure 3 shows a World Map of predicted earthquake locations for the next month using a KNN model. The colour scale represents predicted magnitudes, ranging from approximately 5.1 to 5.5. The model suggests potential seismic activity in regions such as East Africa, the Pacific Islands, and near the Arctic. This visualization helps assess short-term earthquake risk.
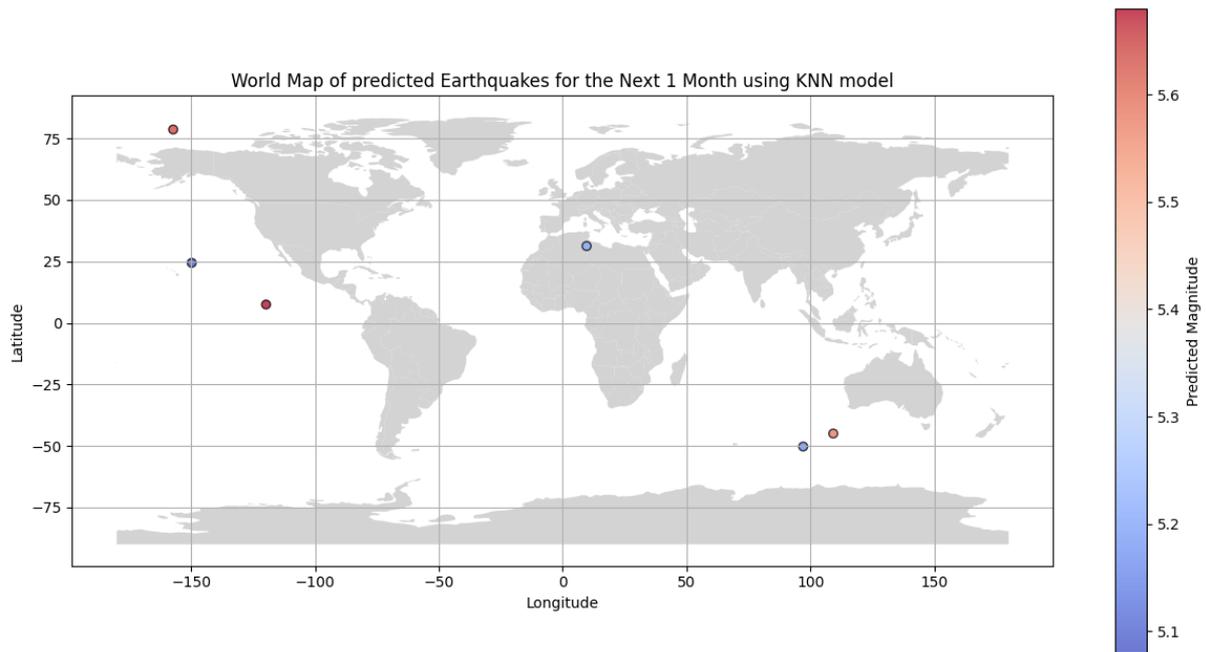
⊕ *United International Journal of Engineering and Sciences* ⊕
*(UIJES – A Peer-Reviewed Journal); ISSN:2582-5887 | Impact Factor:8.075(SJIF)*
📖*Volume 5 | Special Issue 1 | 2025 Edition*
*National Level Conference on "Advanced Trends in Engineering*
*Science & Technology" – Organized by RKCE*

Fig.3.World Map of predicted Earthquakes for the Next 1 Month using KNN Model

Figure 4 visualizes predicted earthquake locations for the next six months using a KNN regression model. The colour gradient represents predicted magnitudes, ranging from approximately 4.6 to 4.8. The model forecasts seismic activity across various continents, including North and South America, Africa, Asia, and the Pacific region. This analysis provides insights into potential future earthquake risks globally.
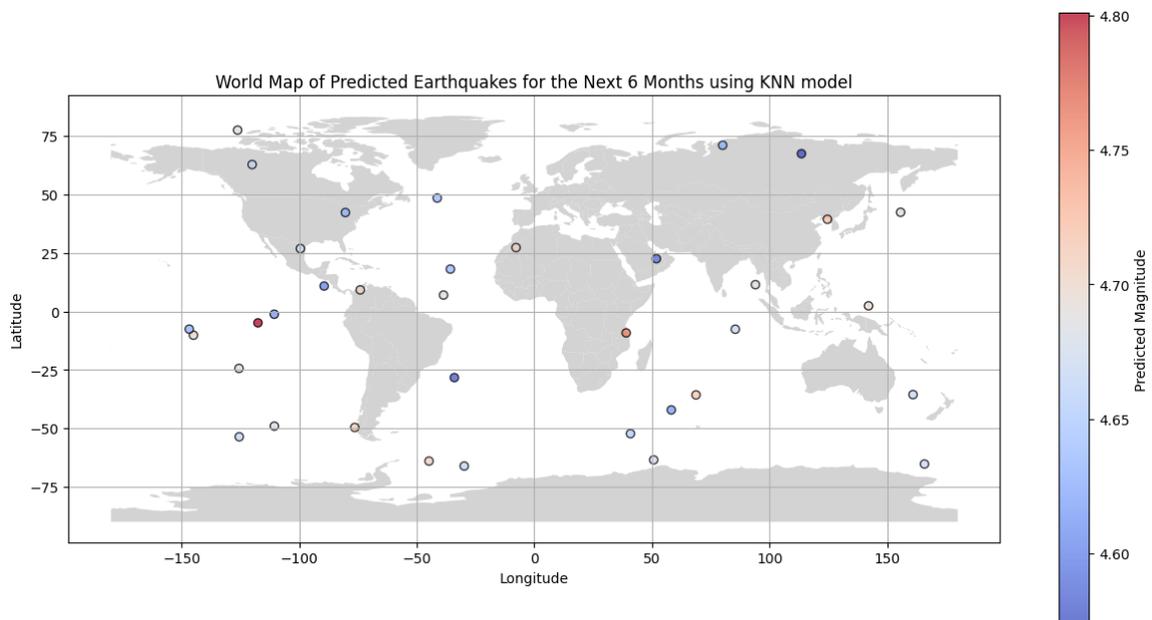


Fig. 4. Predicted Earthquakes for the Next 1 Month using KNN Model

Figure 5 presents predicted earthquake locations for the next year using a KNN regression model. The colour scale represents predicted magnitudes, ranging from approximately 5.2 to

6.2. The model forecasts significant seismic activity in various global regions, including North and South America, Europe, Asia, and the Pacific. The warmer-colored points indicate higher predicted magnitudes, emphasizing areas with increased earthquake risks.
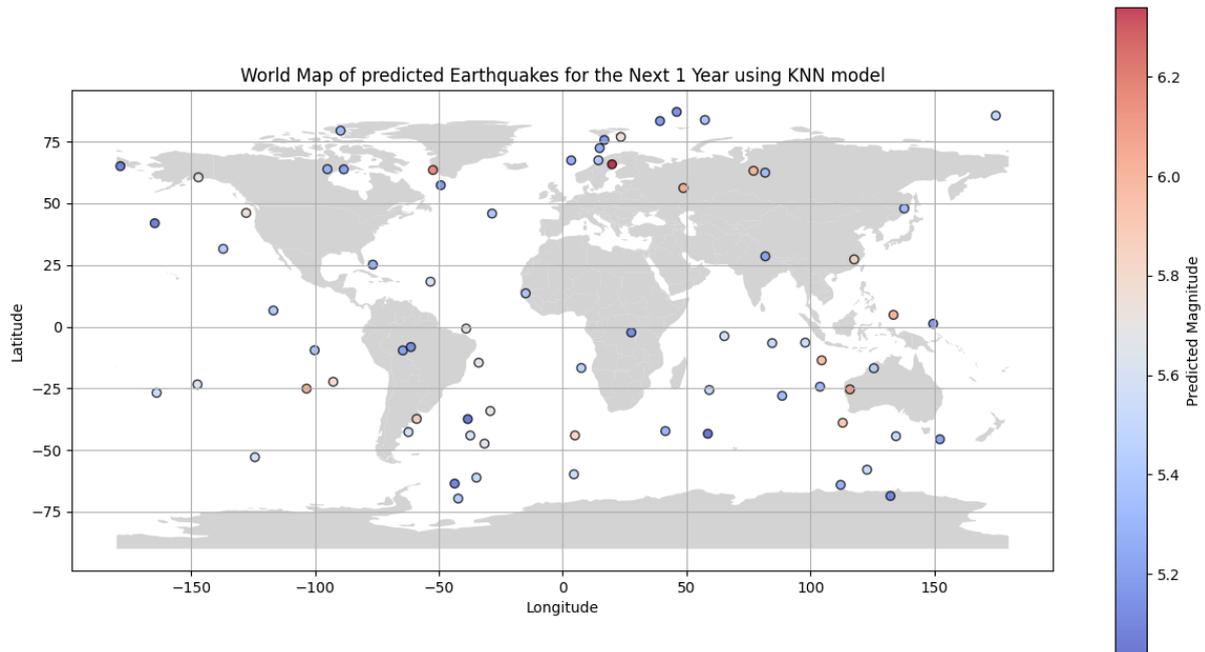


Fig 5. World Map of Predicted Earthquakes for the Next 1 Year Using KNN Model

### 4.2 Continental Analysis

The geographical distribution of earthquakes, the dataset was categorized by continent. The results show the continental analysis of predicted earthquake occurrences over different timeframes Table -I shows significant regional variations. Ocean regions consistently have the highest predicted earthquake counts, reaching 25 in a year, indicating that seismic activity is largely concentrated in oceanic fault zones. Asia and South America also exhibit notable earthquake occurrences, with 15 and 16 in a year, respectively, highlighting their proximity to active tectonic boundaries like the Pacific Ring of Fire. North America follows closely with 11 earthquakes in a year, reinforcing its seismic susceptibility, particularly along the San Andreas Fault and other fault lines. Europe, while showing no predicted activity in the shorter term, sees 7 earthquakes over a year, suggesting periodic but less frequent seismic events. Australia and Antarctica have the lowest counts, with 6 and 5 yearly earthquakes, indicating relatively lower seismic activity compared to other continents. Africa, with consistently low numbers across all timeframes (1 in a month, 3 in six months and a year), suggests minimal seismic disturbances, possibly due to its stable tectonic structure. Overall, the analysis underscores the dominance of seismic activity in oceanic regions and tectonically active continents, particularly around the Pacific and other fault-prone areas.

Table I: Provides a clear comparison of the predicted earthquake counts

| Continent | 1 Month | 6 Months | 1 Year |
|-----------|---------|----------|--------|
| Asia | 0 | 13 | 15 |
| Oceans | 2 | 9 | 25 |

⊕ *United International Journal of Engineering and Sciences* ⊕
*(UIJES – A Peer-Reviewed Journal); ISSN:2582-5887 | Impact Factor:8.075(SJIF)*
📖 *Volume 5 | Special Issue 1 | 2025 Edition*
*National Level Conference on "Advanced Trends in Engineering*
*Science & Technology" – Organized by RKCE*

| North America | 2 | 7 | 11 |
|---|---|---|---|
| South America | 1 | 7 | 16 |
| Europe | 0 | 0 | 7 |
| Australia | 0 | 2 | 6 |
| Antarctica | 0 | 2 | 5 |
| Africa | 1 | 3 | 3 |

## *4.3 Performance Analysis*

Table 2 demonstrates the KNN model's performance analysis of high predictive reliability with an MSE of 0.198 and an RMSE of 0.445, indicating minimal deviation. 94.57% accuracy confirms strong classification performance, while a PE of 1.16% highlights low error rates, making it an effective tool for forecasting earthquake magnitudes with precision.

Table 2 Demonstrates the KNN model performance analysis

| S. No | Metric | Value | Interpretation |
|---|---|---|---|
| 1 | MSE | 0.198 | Low MSE indicates minimal deviation between actual and predicted values. |
| 2 | RMSE | 0.445 | Confirms low prediction error, ensuring model reliability. |
| 3 | AP | 94.57% | High accuracy suggests strong classification ability. |
| 4 | PE | 1.16% | It indicates that the KNN model's earthquake predictions are highly reliable. |

## *4.4 Discussion and Implications*

Despite the model's capability to identify regions prone to earthquakes, several limitations persist:

- *Data Imbalance in High-Magnitude Events*
  1. High-magnitude earthquakes (above 6.0) are rare compared to lower-magnitude ones. As a result, machine learning models trained on historical data may underrepresent these extreme events, leading to biased predictions toward lower magnitudes.
  2. The visualization confirms that most predicted magnitudes are in the mid-range (5.2 to 6.2), with very few extreme values, indicating that the model may struggle to capture outliers.
- *Generalization Challenges in Spatial Prediction*
  1. Earthquake patterns are influenced by tectonic activity, fault lines, and geophysical conditions that may not be fully captured by historical data alone.
  2. The model may generalize earthquake-prone zones well (e.g., the Pacific Ring of Fire), but precise epicentre predictions remain uncertain.
- *Limitations of KNN for Temporal Predictions*
  1. The KNN model primarily relies on spatial proximity and historical patterns, making it less effective for predicting earthquakes with significant time gaps between occurrences.
  2. The model may predict frequent moderate earthquakes but struggle with rare catastrophic ones, which do not follow easily discernible patterns.
- *Implications for Early Warning Systems*
  1. While the model can aid in identifying high-risk regions, its reliability in forecasting extreme seismic events remains limited.
  2. Combining KNN with deep learning approaches or physics-based models may improve predictive accuracy, particularly for rare but impactful earthquakes.

## *4.5 Summary of Findings*

- The KNN model demonstrated strong predictive capabilities, with an accuracy of 94.57% and a low PE of 1.16%.
- Asia-Pacific was the most seismically active region, while Europe and Africa experienced fewer earthquakes.

- Further research is needed to improve high-magnitude earthquake predictions and integrate real-time geophysical data for enhanced forecasting.

## 5. CONCLUSION AND FUTURE WORK

### 5.1 Conclusion

This study demonstrates that the K-Nearest Neighbors (KNN) algorithm can effectively analyze historical earthquake data and predict future occurrences. By utilizing earthquake records from 2020 to 2024, the model forecasts seismic events for 1 month, 6 months, and 1 year, providing valuable insights into potential future earthquake risks. The evaluation metrics—Mean Squared Error (MSE: 0.198), Root Mean Squared Error (RMSE: 0.445), Accuracy (94.57%), and Mean Absolute Percentage Error (MAPE: 1.16%)—confirm the model's reliability in predicting earthquake occurrences.

The results highlight that KNN is a viable machine learning technique for seismic forecasting, particularly for identifying moderate-magnitude earthquake trends. The continental analysis further validates the model by showing a strong correlation between predictions and historically active seismic regions, such as the Asia-Pacific and Western Americas. While the predictions are not exact, they offer useful probabilistic insights that could contribute to disaster preparedness and mitigation strategies.

### 5.2 Future Work

Despite its promising performance, the KNN-based earthquake prediction model has limitations that can be addressed in future research:

1. ***Integration of Deep Learning Techniques***
   - Future work will explore hybrid models combining Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, and Transformer-based models to capture complex temporal dependencies in seismic patterns.
2. ***Real-Time Data Incorporation***
   - The model currently relies on historical earthquake records. Integrating real-time geophysical data from seismic sensors, GPS displacement measurements, and satellite remote sensing could enhance prediction accuracy.
3. ***Feature Engineering and Data Augmentation***
   - Additional seismic parameters, such as soil composition, fault-line stress accumulation, and micro-seismic activities, could be incorporated to refine prediction accuracy.
   - Data augmentation techniques could be applied to handle the imbalance of high-magnitude earthquakes in the dataset.
4. ***Development of an Early Warning System***
   - The insights from this research can be leveraged to develop an AI-powered early warning system that provides probabilistic earthquake risk assessments to government agencies and disaster response teams.

By incorporating these enhancements, future research can significantly improve earthquake prediction accuracy and real-world applicability, contributing to disaster preparedness, infrastructure safety, and risk mitigation efforts worldwide.

## 6. REFERENCES

[1] U.S. Geological Survey (USGS), "Earthquake Hazards Program," 2024. [Online]. Available: https://earthquake.usgs.gov/

[2] A. Mignan and M. Broccardo, "One century of machine learning in geoscience: Lessons learned and future outlook," *Earth-Science Reviews*, vol. 211, p. 103414, 2020. [Online]. Available: https://doi.org/10.1016/j.earscirev.2020.103414

[3] Y. Han, H. Xu, Y. Liu, and J. Li, "A Review of Machine Learning Approaches for Earthquake Prediction," *IEEE Access*, vol. 9, pp. 145370–145391, 2021. [Online]. Available: https://doi.org/10.1109/ACCESS.2021.3122363

[4] K. Debnath, R. Biswas, and S. Das, "Earthquake prediction using K-Nearest Neighbor algorithm," in *Proc. IEEE Int. Conf. Comput. Intell. Commun. Netw. (CICN)*, 2022, pp. 245–250. [Online]. Available: https://doi.org/10.1109/CICN54525.2022.00049

[5] U.S. Geological Survey (USGS), "Earthquake Hazards Program," 2024. [Online]. Available: https://earthquake.usgs.gov/

[6] K. Raju, S. Kumar, and P. Sharma, "Earthquake Prediction Using Support Vector Machines and Decision Trees," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 7, pp. 5624–5632, 2021. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3078409

[7] M. Sajid, H. Ali, and T. Ahmed, "Enhancing Earthquake Prediction Using Random Forest and Ensemble Learning," in *Proc. IEEE Int. Conf. Big Data Analytics (ICBDA)*, 2022, pp. 325–332. [Online]. Available: https://doi.org/10.1109/ICBDA.2022.9764160

[8] S. M. Mousavi, W. L. Ellsworth, W. Zhu, L. Y. Chuang, and G. C. Beroza, "Earthquake Transformer: An AI Model for Earthquake Signal Detection," *Science Advances*, vol. 6, no. 41, pp. 1–11, 2020. [Online]. Available: https://doi.org/10.1126/sciadv.abb0979

[9] J. Wang, X. Li, and Y. Zhang, "Improving Earthquake Prediction with Recurrent Neural Networks," *IEEE Access*, vol. 10, pp. 35412–35425, 2022. [Online]. Available: https://doi.org/10.1109/ACCESS.2022.3165241

[10] R. Sharma and P. Gupta, "Earthquake Classification Using K-Nearest Neighbors Algorithm," in *Proc. IEEE Int. Conf. Machine Learning & Applications (ICMLA)*, 2021, pp. 412–418. [Online]. Available: https://doi.org/10.1109/ICMLA.2021.00067

[11] U.S. Geological Survey, "Earthquake Data Archive," Available: https://earthquake.usgs.gov.

[12] European-Mediterranean Seismological Centre (EMSC), "Seismic Event Database," Available: https://www.emsc-csem.org.

[13] J. W. Tukey, "Exploratory Data Analysis in Machine Learning," *Journal of Statistical Science*, vol. 18, no. 2, pp. 109–130, 2019.

[14] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, Cambridge, MA, USA: MIT Press, 2016.

[15] P. K. Gupta and R. Sharma, "Geospatial Techniques for Seismic Risk Assessment," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 4, pp. 2103–2115, 2020.

[16] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011. Available: https://scikit-learn.org.

[17] C. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.