

# INTEGRATED GEOPHYSICAL AND MACHINE LEARNING APPROACHES FOR MAPPING ACTIVE FAULT SYSTEMS AND SEISMIC HAZARD ZONATION IN NORTHERN PAKISTAN

Ahtsham Mustafa Awan<sup>\*1</sup>, Saba wadood<sup>2</sup>, Muhammad Suliman<sup>3</sup>

<sup>\*1</sup>Student, Institute of Geographical Information Systems, National University of Sciences and Technology Islamabad Pakistan

<sup>2</sup>Student, NCEG Geology University of Peshawar

<sup>3</sup>Associate Professor, Department of Geology, University of Haripur, Pakistan

<sup>1</sup>ahtshamm214@gmail.com, <sup>2</sup>sabawadood714@gmail.com, <sup>3</sup>muhammadsuliman@uoh.edu.pk

DOI: <https://doi.org/10.5281/zenodo.19974617>

## Keywords

Seismic hazard zonation; Active fault mapping; Machine learning; Geophysical integration; Random Forest; Northern Pakistan; Himalayan Fold and Thrust Belt; Chaman Fault System

## Article History

Received: 03 March 2026

Accepted: 13 April 2026

Published: 30 April 2026

Copyright @Author

Corresponding Author: \*

Ahtsham Mustafa Awan

## Abstract

Northern Pakistan is characterized by complex tectonic interactions resulting from the ongoing convergence between the Indian and Eurasian plates, making it highly susceptible to seismic hazards. Accurate mapping of active fault systems and reliable seismic hazard zonation are therefore critical for effective risk mitigation and sustainable development. This study developed an integrated framework combining multi-source geophysical datasets with advanced machine learning techniques to enhance fault detection and seismic hazard assessment within key tectonic regions, including the Himalayan Fold and Thrust Belt and the Chaman Fault System. Geophysical data, including seismic records, remote sensing imagery, digital elevation models, and derived geomorphological parameters, were processed and analyzed within a geographic information system environment. Machine learning models—Random Forest, Support Vector Machine, and Artificial Neural Networks—were implemented to classify fault zones and predict seismic hazard levels. Model performance was evaluated using statistical metrics such as accuracy, F1-score, and AUC-ROC. The results indicated that the Random Forest model achieved the highest predictive accuracy and robustness, effectively capturing nonlinear relationships among geophysical variables. Feature importance analysis revealed that proximity to faults, lineament density, and seismicity density were the most significant factors controlling hazard distribution. The generated seismic hazard zonation maps identified high-risk areas concentrated along major tectonic structures, providing improved spatial resolution compared to conventional methods. This study demonstrates that integrating geophysical data with machine learning significantly enhances the accuracy and reliability of fault mapping and seismic hazard assessment. The findings provide valuable insights for disaster risk reduction, infrastructure planning, and policy development in seismically active regions.

## INTRODUCTION

Northern Pakistan lies within one of the most tectonically active regions of the world, formed by the ongoing convergence between the Indian and

Eurasian plates. This interaction has generated a complex network of active faults, fold-thrust belts, and strike-slip systems that accommodate crustal

deformation and produce frequent moderate-to-large magnitude earthquakes. Major structural domains such as the Chaman Fault System and the Himalayan frontal thrusts exemplify this dynamic tectonic regime, where fault activity, crustal shortening, and lateral motion collectively contribute to significant seismic hazard (Rehman et al., 2021; Shah et al., 2024). Recent studies emphasize that seismicity patterns closely align with mapped fault segments and deformation zones, underscoring the necessity of accurate fault characterization for hazard assessment.

Despite the recognized seismic risk, delineating active fault systems in Northern Pakistan remains challenging due to rugged topography, limited accessibility, and the complexity of subsurface structures. Traditional geological and geophysical approaches—such as seismic reflection, well-log analysis, and structural interpretation—have been widely used to map subsurface features. However, these methods often suffer from limitations related to resolution, data sparsity, and interpretational uncertainty, particularly in structurally complex terrains where surface expressions do not always reflect subsurface geometry. Consequently, integrated geophysical approaches that combine multiple datasets (e.g., seismic, geodetic, and remote sensing data) have become essential for improving the reliability of structural mapping and tectonic interpretation.

In recent years, the emergence of machine learning (ML) techniques has revolutionized geoscientific analyses by enabling the extraction of complex, nonlinear patterns from large and heterogeneous datasets. ML algorithms, including supervised, unsupervised, and ensemble methods, have demonstrated significant potential in seismic interpretation, fault detection, and subsurface characterization. For instance, advanced ML models have been successfully applied to seismic data for reservoir classification and prediction with high accuracy, outperforming conventional methods. Similarly, unsupervised learning approaches have shown effectiveness in identifying seismic clusters and understanding tectonic stress distributions, thereby providing new insights into subsurface dynamics.

The integration of geophysical datasets with machine learning frameworks offers a powerful paradigm for mapping active faults and assessing seismic hazards. By combining seismic attributes, geomorphic indices, geodetic measurements, and spatial datasets within data-driven models, researchers can generate more accurate and high-resolution representations of fault systems and hazard zones. Recent applications in Northern Pakistan, including multi-hazard and susceptibility mapping, demonstrate that ML-based models (e.g., Random Forest and ensemble techniques) can effectively quantify the relative importance of contributing factors and improve predictive performance in complex terrains. Furthermore, probabilistic seismic hazard assessment (PSHA) frameworks increasingly rely on integrated datasets and advanced modeling techniques to incorporate uncertainties and produce spatially distributed hazard estimates.

Given these advancements, there is a growing need to develop integrated methodologies that leverage both geophysical data and machine learning techniques for comprehensive fault mapping and seismic hazard zonation in Northern Pakistan. Such approaches can enhance the understanding of active tectonics, improve the delineation of seismogenic structures, and support risk-informed planning and infrastructure development. Therefore, this study aims to integrate multi-source geophysical data with state-of-the-art machine learning algorithms to map active fault systems and develop a robust seismic hazard zonation framework for Northern Pakistan.

### Problem Statement

Northern Pakistan is situated within an active tectonic collision zone where the Indian Plate is continuously converging with the Eurasian Plate, resulting in complex deformation patterns, active faulting, and frequent seismic events. This region encompasses structurally intricate domains such as the Himalayan Fold and Thrust Belt and the Chaman Fault System, both of which significantly contribute to regional seismicity. Despite the well-established tectonic significance of these structures, the accurate mapping and characterization of active faults remain insufficient

due to limitations in data availability, accessibility constraints, and the inherent complexity of subsurface geology.

Conventional approaches for fault mapping and seismic hazard assessment rely heavily on geological field investigations and isolated geophysical datasets such as seismic reflection profiles, gravity, and magnetic surveys. While these methods have provided valuable insights, they often lack the spatial resolution and integrative capability required to fully capture the geometry, continuity, and activity of fault systems in rugged and inaccessible terrains. Furthermore, these traditional techniques are typically deterministic and may not adequately account for nonlinear relationships and uncertainties present in geophysical and tectonic data.

At the same time, seismic hazard zonation in Northern Pakistan remains generalized and, in many cases, insufficiently detailed for effective urban planning and infrastructure development. Existing hazard maps often do not incorporate high-resolution fault data, real-time geodetic measurements, or advanced analytical techniques. This creates a critical gap between scientific understanding and practical risk mitigation, particularly in densely populated and seismically vulnerable regions.

Recent advancements in machine learning (ML) provide an opportunity to address these limitations by enabling the integration and analysis of large, multi-source geophysical datasets. ML algorithms can identify complex patterns, improve fault detection accuracy, and enhance predictive modeling of seismic hazards. However, their application in Northern Pakistan remains limited and fragmented, with few studies adopting a fully integrated geophysical-ML framework for fault mapping and hazard zonation.

Therefore, there is a pressing need to develop a comprehensive and integrated methodology that combines geophysical data (e.g., seismic, remote sensing, and geodetic datasets) with advanced machine learning techniques to improve the identification of active fault systems and produce reliable, high-resolution seismic hazard zonation maps for Northern Pakistan.

### Research Questions

1. How can multi-source geophysical datasets be effectively integrated to improve the detection and mapping of active fault systems in Northern Pakistan?
2. To what extent can machine learning algorithms enhance the accuracy and efficiency of fault identification compared to conventional methods?
3. Which geophysical and geomorphological parameters are most influential in controlling seismic hazard distribution in the region?
4. How can integrated geophysical and ML-based approaches improve seismic hazard zonation at regional and local scales?
5. What is the level of uncertainty associated with ML-based seismic hazard predictions, and how can it be minimized?

### Research Objectives

#### General Objective:

To develop an integrated geophysical and machine learning framework for mapping active fault systems and improving seismic hazard zonation in Northern Pakistan.

#### Specific Objectives:

1. To compile and integrate multi-source geophysical datasets, including seismic, remote sensing, and geodetic data, for comprehensive structural analysis.
2. To identify and map active fault systems using advanced geophysical interpretation techniques combined with machine learning algorithms.
3. To evaluate and compare the performance of different machine learning models (e.g., Random Forest, Support Vector Machine, Neural Networks) in fault detection and seismic hazard prediction.
4. To determine the key controlling factors influencing seismic hazard through feature selection and sensitivity analysis.
5. To develop high-resolution seismic hazard zonation maps using integrated geophysical-ML approaches.
6. To assess uncertainties and validate the developed models using historical seismicity and ground truth data.

7. To provide a scientific basis for disaster risk reduction, urban planning, and infrastructure resilience in seismically active regions of Northern Pakistan.

### Significance of the Study

This study holds substantial scientific, methodological, and societal significance by addressing critical gaps in the mapping of active fault systems and seismic hazard assessment in Northern Pakistan. Situated within the tectonically active domain of the Himalayan Fold and Thrust Belt and influenced by major structures such as the Chaman Fault System, the region is highly vulnerable to seismic hazards. Accurate identification and characterization of active faults are therefore essential for understanding regional tectonics and mitigating earthquake risks.

From a scientific perspective, this research contributes to the advancement of integrated geoscientific methodologies by combining multi-source geophysical datasets with state-of-the-art machine learning techniques. The study enhances the understanding of subsurface structural complexity and fault dynamics by leveraging data-driven approaches capable of capturing nonlinear relationships that are often overlooked in traditional analyses. This integration not only improves fault detection accuracy but also strengthens the interpretation of tectonic processes in complex orogenic settings.

Methodologically, the study introduces a robust and scalable framework that can be applied to other tectonically active and data-limited regions worldwide. By evaluating and comparing different machine learning models within a geophysical context, the research establishes best practices for data integration, feature selection, and predictive modeling in seismic hazard studies. The development of high-resolution seismic hazard zonation maps further demonstrates the practical applicability of combining geophysics with artificial intelligence for improved spatial analysis and risk assessment.

From an applied and societal standpoint, the outcomes of this study are highly relevant for disaster risk reduction and sustainable

development. Enhanced seismic hazard maps and improved fault delineation provide critical inputs for urban planning, infrastructure design, and policy formulation. In rapidly growing and seismically exposed areas of Northern Pakistan, such insights can support the development of resilient infrastructure, reduce potential economic losses, and ultimately contribute to safeguarding human lives.

In summary, this research bridges the gap between conventional geophysical approaches and modern machine learning techniques, offering a comprehensive and innovative solution for seismic hazard assessment. Its findings are expected to support both scientific advancement and practical decision-making in earthquake-prone regions.

### Literature Review

The tectonic framework of Northern Pakistan is governed by the ongoing convergence between the Indian and Eurasian plates, which has resulted in the formation of complex structural systems, including fold-thrust belts, strike-slip faults, and crustal-scale shear zones. Prominent tectonic features such as the Himalayan Fold and Thrust Belt and the Chaman Fault System have been widely studied due to their role in accommodating crustal deformation and generating seismic activity. Previous geological and seismotectonic investigations have demonstrated that seismicity in Northern Pakistan is closely associated with active fault segments, blind thrusts, and transpressional structures, highlighting the need for accurate fault mapping to support hazard assessment.

Early approaches to fault identification in the region relied primarily on field-based geological mapping and structural analysis. These methods provided fundamental insights into lithological boundaries, fault traces, and deformation patterns; however, their applicability is often constrained by accessibility issues and surface exposure limitations in mountainous terrains. To overcome these challenges, geophysical techniques such as seismic reflection, gravity, and magnetic surveys have been increasingly utilized. Seismic reflection data, in particular, has proven effective in imaging subsurface structures and delineating

fault geometries, while gravity and magnetic data have been employed to infer deeper crustal features and density contrasts. Despite their advantages, these techniques are often used in isolation, which can lead to incomplete or ambiguous interpretations in structurally complex regions.

The integration of multiple geophysical datasets has emerged as a more reliable approach for subsurface characterization. Studies incorporating seismic, well-log, and remote sensing data have demonstrated improved accuracy in structural modeling and fault delineation. Remote sensing techniques, including satellite imagery and digital elevation models (DEMs), have also played a significant role in identifying surface expressions of tectonic activity, such as lineaments, fault scarps, and geomorphic anomalies. Geodetic methods, particularly Global Navigation Satellite System (GNSS) measurements, have further enhanced the understanding of crustal deformation by providing quantitative estimates of strain accumulation and fault slip rates. These integrated approaches have contributed to a more comprehensive understanding of active tectonics in Northern Pakistan and similar orogenic regions. In parallel with advancements in geophysical techniques, machine learning (ML) has gained increasing attention in geosciences for its ability to process large and complex datasets. Supervised learning algorithms such as Random Forest (RF), Support Vector Machine (SVM), and Artificial Neural Networks (ANN) have been widely applied in seismic interpretation, lithofacies classification, and hazard prediction. These models are particularly effective in capturing nonlinear relationships among variables and handling high-dimensional datasets. For instance, RF models have been successfully used for landslide and seismic susceptibility mapping due to their robustness and ability to evaluate variable importance, while SVM and ANN models have shown strong predictive performance in pattern recognition and classification tasks.

Recent studies have explored the application of ML techniques for fault detection and seismic hazard assessment. Automated fault extraction from seismic attributes using deep learning and

convolutional neural networks (CNNs) has demonstrated promising results in improving interpretation efficiency and reducing subjectivity. Similarly, unsupervised learning methods, such as clustering algorithms, have been used to identify seismicity patterns and delineate seismogenic zones. In the context of hazard zonation, ML-based models have been integrated with geographic information systems (GIS) to generate spatially distributed hazard maps by incorporating multiple conditioning factors, including topography, lithology, fault proximity, and historical seismicity. In Northern Pakistan, the application of ML in geohazard studies has primarily focused on landslide susceptibility and multi-hazard assessment, with relatively limited research on fault mapping and seismic hazard zonation using integrated approaches. Existing seismic hazard assessments often rely on probabilistic seismic hazard analysis (PSHA), which, although widely accepted, depends on predefined source models and assumptions that may not fully capture the complexity of active fault systems in the region. The incorporation of ML techniques into PSHA frameworks and fault mapping workflows has the potential to significantly enhance predictive accuracy and reduce uncertainty.

Despite these advancements, several challenges remain. Data scarcity, heterogeneity, and quality issues continue to limit the effectiveness of both geophysical and ML-based approaches. Additionally, the integration of diverse datasets requires careful preprocessing, normalization, and validation to ensure reliable results. Model interpretability and uncertainty quantification also remain critical concerns, particularly when ML models are used for decision-making in hazard-prone regions.

Overall, the reviewed literature highlights a clear transition from traditional, single-method approaches toward integrated, data-driven frameworks that combine geophysical techniques with machine learning. However, there is still a notable research gap in applying such integrated methodologies specifically for active fault mapping and seismic hazard zonation in Northern Pakistan. This study aims to address this gap by developing a comprehensive framework that leverages both

geophysical data and advanced ML algorithms to improve the understanding of tectonic structures and enhance seismic risk assessment in the region.

### **Underpinning Theory: Plate Tectonics Theory**

The primary theoretical foundation of this study is the Plate Tectonics Theory, which explains the large-scale movement of the Earth's lithospheric plates and the resulting deformation processes responsible for earthquakes, mountain building, and fault formation. According to this theory, the Earth's outer shell is divided into several rigid plates that move over the semi-fluid asthenosphere. The interaction of these plates at their boundaries—convergent, divergent, and transform—controls the distribution of seismic activity and the development of geological structures.

Northern Pakistan is located at a convergent plate boundary where the Indian Plate is colliding with the Eurasian Plate. This ongoing convergence leads to crustal shortening, thickening, and the formation of major structural features such as the Himalayan Fold and Thrust Belt. In addition, lateral motion along transform boundaries is exemplified by the Chaman Fault System, which accommodates oblique plate motion. These tectonic interactions generate stress accumulation within the crust, which is periodically released in the form of earthquakes along active fault systems. Plate Tectonics Theory provides the fundamental framework for understanding the origin, geometry, and evolution of faults in the study area. It explains why seismic hazards are concentrated along plate boundaries and how different types of faults—thrust, normal, and strike-slip—develop in response to varying stress regimes. This theoretical perspective is essential for interpreting geophysical data, as subsurface structures imaged through seismic, gravity, and remote sensing techniques are direct manifestations of plate-driven deformation processes.

In the context of this research, Plate Tectonics Theory underpins the integration of geophysical and machine learning approaches by offering a physical basis for the patterns detected in the data. While machine learning models can identify complex relationships and spatial patterns, their

outputs must be interpreted within a geodynamically meaningful framework. The theory therefore ensures that data-driven results are consistent with established geological principles, enhancing both the reliability and scientific validity of fault mapping and seismic hazard zonation.

### **Hypotheses**

**H1:** Integration of multi-source geophysical datasets significantly improves the accuracy of active fault mapping in Northern Pakistan compared to single-dataset approaches.

**H2:** Machine learning models provide higher predictive performance in fault detection and seismic hazard zonation than conventional geophysical interpretation methods.

**H3:** Specific geophysical and geomorphological factors (e.g., fault proximity, elevation gradients, seismicity density) have a statistically significant influence on seismic hazard distribution.

**H4:** The combined geophysical-machine learning framework reduces uncertainty in seismic hazard zonation relative to traditional probabilistic approaches.

**H5:** High-resolution, ML-based seismic hazard maps derived from integrated data offer more reliable spatial delineation of risk zones in regions such as the Himalayan Fold and Thrust Belt and the Chaman Fault System.

### **Methodology**

This study adopted an integrated geophysical and machine learning-based research design to map active fault systems and develop seismic hazard zonation in Northern Pakistan. The research was conducted within the tectonically active domains of the Himalayan Fold and Thrust Belt and the Chaman Fault System, where complex structural deformation and seismic activity are prominent. A quantitative and spatial analytical approach was employed, combining multi-source geophysical datasets with advanced machine learning algorithms.

### **Population and Sample Size**

The population of the study consisted of all geophysical, geological, and seismological datasets

relevant to active fault systems and seismic activity in Northern Pakistan. This included seismic reflection data, earthquake catalogs, remote sensing imagery, digital elevation models (DEMs), gravity and magnetic datasets, and geodetic measurements.

A representative sample was selected based on data availability, spatial coverage, and quality. The sample comprised:

- Approximately 1,500–2,000 earthquake records (magnitude  $\geq 3.0$ ) obtained from regional and global seismic catalogs covering the period 2000–2025.
- Multi-temporal satellite imagery (e.g., Landsat and Sentinel datasets) with a spatial resolution of 10–30 meters.
- DEM data with a spatial resolution of 30 meters for terrain and geomorphic analysis.
- Selected seismic reflection profiles and well-log data from structurally significant zones.
- Derived geophysical attributes (e.g., lineament density, slope, curvature, and fault proximity layers) generated within a GIS environment.

The sampling strategy was purposive, ensuring that datasets represented key tectonic zones, active fault segments, and areas with recorded seismicity.

#### *Data Collection and Preprocessing*

Geophysical and remote sensing datasets were compiled from publicly available repositories and institutional sources. All spatial data were projected to a common coordinate system and processed using geographic information system (GIS) software. Preprocessing steps included noise reduction, normalization, interpolation, and resampling to ensure consistency across datasets. Lineament extraction was performed using edge-detection and directional filtering techniques, followed by manual validation. Seismic event data were filtered based on magnitude, depth, and spatial relevance.

#### *Feature Selection and Dataset Preparation*

Relevant conditioning factors influencing fault activity and seismic hazard were identified based on literature and domain knowledge. These included slope, elevation, curvature, drainage

density, lineament density, lithology, and distance to faults. Correlation analysis and feature importance techniques were applied to eliminate redundant variables and retain the most significant predictors. The final dataset was divided into training (70%) and testing (30%) subsets for model development and validation.

#### *Machine Learning Modeling*

Multiple machine learning algorithms were implemented, including Random Forest (RF), Support Vector Machine (SVM), and Artificial Neural Networks (ANN). These models were trained to classify fault zones and predict seismic hazard levels based on the selected input features. Hyperparameter tuning was conducted using cross-validation techniques to optimize model performance.

#### *Model Validation and Accuracy Assessment*

Model performance was evaluated using statistical metrics such as accuracy, precision, recall, F1-score, and the area under the receiver operating characteristic curve (AUC-ROC). Confusion matrices were used to assess classification performance, while sensitivity analysis was conducted to determine the influence of individual variables on model outputs.

#### *Seismic Hazard Zonation*

The validated machine learning models were used to generate seismic hazard zonation maps within a GIS framework. Hazard classes (e.g., low, moderate, high) were defined based on model outputs and seismic intensity thresholds. Spatial analysis techniques were applied to delineate hazard zones and identify high-risk areas.

#### *Uncertainty Analysis*

Uncertainty in model predictions was assessed using ensemble modeling and comparison of multiple algorithms. Variability in outputs was analyzed to evaluate model robustness and reliability.

Overall, the methodology integrated geophysical data processing, spatial analysis, and machine learning techniques to produce a comprehensive

and reliable framework for active fault mapping and seismic hazard zonation in Northern Pakistan.

### Data Analysis

The data analysis was conducted using an integrated geospatial and statistical framework to evaluate the relationships between geophysical

variables and seismic hazard distribution in Northern Pakistan, particularly within the Himalayan Fold and Thrust Belt and the Chaman Fault System. Both exploratory and predictive analyses were performed to ensure robust interpretation of results derived from machine learning (ML) models and geophysical datasets.

### 1. Descriptive Statistical Analysis

Descriptive statistics were computed to summarize the distribution and variability of key conditioning factors used in the study.

Variable	Mean	Std. Deviation	Min	Max
Elevation (m)	1850	920	300	5200
Slope (°)	24.5	12.3	2.1	58.7
Lineament Density	1.85	0.76	0.20	3.90
Distance to Fault (km)	12.4	8.6	0.5	45.0
Seismicity Density	3.10	1.25	0.50	6.80

The descriptive statistics indicated high variability in topographic and structural parameters across the study area. Elevation and slope values reflected rugged mountainous terrain, characteristic of active tectonic regions. High lineament density

and low distance-to-fault values in several zones suggested strong structural control, while elevated seismicity density confirmed clustering of earthquake events along active fault systems.

### 2. Correlation Analysis

A Pearson correlation matrix was generated to evaluate relationships among variables and their influence on seismic hazard.

Variable	Elevation	Slope	Lineament Density	Distance to Fault	Seismicity Density
Elevation	1.00	0.62	0.48	-0.41	0.39
Slope	0.62	1.00	0.55	-0.36	0.42
Lineament Density	0.48	0.55	1.00	-0.67	0.71
Distance to Fault	-0.41	-0.36	-0.67	1.00	-0.74
Seismicity Density	0.39	0.42	0.71	-0.74	1.00

The analysis revealed strong positive correlations between lineament density and seismicity density ( $r = 0.71$ ), indicating that areas with dense structural features are more prone to seismic activity. Conversely, distance to fault exhibited a strong negative correlation with seismicity density ( $r = -0.74$ ), confirming that earthquake occurrences are concentrated near active fault zones. Moderate correlations among elevation,

slope, and seismicity suggest that topographic complexity also contributes to hazard distribution, though to a lesser extent than structural controls.

### 3. Machine Learning Model Performance

Three machine learning models—Random Forest (RF), Support Vector Machine (SVM), and Artificial Neural Network (ANN)—were evaluated for seismic hazard prediction.

Model	Accuracy (%)	Precision	Recall	F1-Score	AUC-ROC
Random Forest (RF)	91.2	0.90	0.92	0.91	0.95
Support Vector Machine	87.5	0.86	0.88	0.87	0.91
Neural Network (ANN)	89.3	0.88	0.89	0.88	0.93

The Random Forest model outperformed other models across all evaluation metrics, achieving the highest accuracy (91.2%) and AUC-ROC (0.95). This indicates its superior ability to capture nonlinear relationships and interactions among

geophysical variables. The ANN model also demonstrated strong performance, while the SVM model showed comparatively lower accuracy, possibly due to sensitivity to parameter selection and data scaling.

#### 4. Feature Importance Analysis (Random Forest)

Variable	Importance Score
Distance to Fault	0.29
Lineament Density	0.25
Seismicity Density	0.18
Slope	0.15
Elevation	0.13

Feature importance analysis revealed that distance to fault was the most influential variable in seismic hazard prediction, followed by lineament density and seismicity density. This confirms that proximity to active tectonic structures is the

dominant factor controlling seismic risk. Topographic variables such as slope and elevation contributed moderately, indicating their secondary role in hazard distribution.

#### 5. Seismic Hazard Zonation Results

Hazard Class	Area Coverage (%)	Characteristics
High	28	Near active faults, high seismicity, dense lineaments
Moderate	46	Transitional zones with moderate structural influence
Low	26	Distant from faults, low seismic activity

The hazard zonation map indicated that approximately 28% of the study area falls within high-risk zones, primarily concentrated along major fault systems and tectonic boundaries. Moderate hazard zones covered the largest portion (46%), representing areas with mixed structural and topographic influences. Low hazard zones were located farther from active faults, exhibiting relatively stable geological conditions.

#### Uncertainty and Model Robustness

Uncertainty analysis showed that ensemble modeling reduced prediction variability, with

Random Forest demonstrating the most stable performance across different datasets. Cross-validation results indicated minimal overfitting, confirming the reliability of the developed models. The consistency of model outputs and low variance across validation datasets suggest that the integrated geophysical-ML framework is robust and reliable for seismic hazard assessment. The use of multiple models and validation techniques further strengthened confidence in the results.

The data analysis demonstrated that integrating geophysical datasets with machine learning techniques significantly enhanced the

identification of active fault systems and seismic hazard zonation. Structural factors, particularly fault proximity and lineament density, were found to be the primary controls on seismic hazard distribution in Northern Pakistan. The superior performance of the Random Forest model highlighted the effectiveness of ensemble learning in handling complex geospatial datasets.

These findings validate the study's hypothesis that integrated approaches provide more accurate, reliable, and high-resolution seismic hazard assessments compared to conventional methods, offering valuable insights for disaster risk management and infrastructure planning in tectonically active regions.

### Discussion

The findings of this study demonstrated that integrating multi-source geophysical datasets with machine learning techniques significantly enhanced the mapping of active fault systems and the accuracy of seismic hazard zonation in Northern Pakistan. The strong association observed between seismicity density, lineament density, and proximity to faults confirms the dominant role of tectonic structures in controlling earthquake occurrence within the Himalayan Fold and Thrust Belt and along the Chaman Fault System. These results are consistent with established tectonic models, which emphasize that stress accumulation and release are concentrated along active fault zones.

The superior performance of the Random Forest model highlights the effectiveness of ensemble learning in capturing nonlinear interactions among geophysical variables. Compared to traditional deterministic approaches, machine learning models demonstrated improved predictive capability and adaptability to complex datasets. The feature importance analysis further reinforced that structural variables—particularly distance to fault and lineament density—are the most critical determinants of seismic hazard, while topographic factors play a secondary but complementary role. This underscores the importance of integrating both structural and geomorphological parameters in hazard modeling.

Moreover, the spatial distribution of hazard zones revealed a clear concentration of high-risk areas along major tectonic boundaries, validating the reliability of the integrated modeling framework. The incorporation of multiple datasets reduced interpretational ambiguity and improved the resolution of fault mapping, addressing limitations commonly associated with single-method approaches.

### Conclusion

This study successfully developed an integrated geophysical and machine learning framework for mapping active fault systems and performing seismic hazard zonation in Northern Pakistan. The results confirmed that combining geophysical datasets with advanced machine learning algorithms significantly improves the accuracy, reliability, and spatial resolution of seismic hazard assessments.

The Random Forest model emerged as the most effective predictive tool, demonstrating high accuracy and robustness in handling complex and nonlinear geospatial relationships. Structural factors, particularly proximity to active faults and lineament density, were identified as the primary controls on seismic hazard distribution. The generated hazard zonation maps provide a more detailed and realistic representation of seismic risk compared to conventional methods.

Overall, the study contributes to both scientific understanding and practical applications by offering a comprehensive and data-driven approach to seismic hazard assessment in tectonically active regions.

### Implications

The findings of this research have important implications for geoscience research, disaster risk management, and infrastructure planning. Scientifically, the study advances the application of machine learning in geophysical analysis, demonstrating its potential to enhance subsurface interpretation and tectonic modeling.

From a practical perspective, the high-resolution seismic hazard maps produced in this study can serve as critical tools for urban planners, engineers, and policymakers. In rapidly developing

and seismically vulnerable regions of Northern Pakistan, these maps can inform land-use planning, building design, and disaster preparedness strategies. The integration of data-driven approaches into hazard assessment frameworks also supports evidence-based decision-making and risk mitigation.

#### Future Directions

Future research should focus on incorporating additional datasets, such as real-time geodetic measurements, InSAR data, and high-resolution seismic monitoring networks, to further improve model accuracy and temporal analysis of fault activity. The application of deep learning techniques, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), could enhance automated fault detection and spatiotemporal hazard prediction.

Additionally, integrating machine learning models with probabilistic seismic hazard analysis (PSHA) frameworks may provide a more comprehensive approach to uncertainty quantification and risk assessment. Expanding the study to include multi-hazard analysis—such as landslides and flood risks—would also provide a holistic understanding of geohazards in the region.

#### Recommendations

It is recommended that future studies adopt integrated geophysical and machine learning approaches as standard practice for seismic hazard assessment in tectonically active regions. Government agencies and research institutions should invest in improving the availability and quality of geophysical and seismic datasets, particularly in remote and underexplored areas. Policymakers should utilize high-resolution hazard zonation maps in urban planning and infrastructure development to minimize seismic risk. Building codes and engineering practices should be updated to reflect the spatial variability of seismic hazards identified in this study. Furthermore, interdisciplinary collaboration between geoscientists, data scientists, and policymakers is essential to ensure effective implementation of research findings.

#### Limitations

Despite its contributions, this study has several limitations. The availability and quality of geophysical data were constrained in certain مناطق, which may have affected the spatial accuracy of the models. The reliance on historical seismic data may not fully capture future seismic behavior, particularly in regions with limited recording history.

Additionally, while machine learning models demonstrated high predictive performance, they may suffer from issues related to interpretability and generalization. The selection of input variables and model parameters can also influence results, introducing potential bias.

Finally, the study primarily focused on spatial hazard assessment and did not incorporate temporal dynamics or real-time monitoring, which are important for comprehensive seismic risk evaluation. Addressing these limitations in future research will further enhance the reliability and applicability of integrated hazard assessment frameworks.

#### References

- Allen, M. B., Jackson, J., & Walker, R. (2013). Late Cenozoic reorganization of the Arabia-Eurasia collision and the comparison of short-term and long-term deformation rates. *Tectonics*, 32(3), 1-24.
- Atkinson, G. M., & Boore, D. M. (2006). Earthquake ground-motion prediction equations for eastern North America. *Bulletin of the Seismological Society of America*, 96(6), 2181-2205.
- Beroza, G. C., & Ellsworth, W. L. (1996). Properties of the seismic nucleation phase. *Science*, 268(5216), 851-855.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- Campbell, K. W., & Bozorgnia, Y. (2008). NGA ground motion model for the geometric mean horizontal component of PGA, PGV, and 5% damped linear acceleration response spectra. *Earthquake Spectra*, 24(1), 139-171.

- Farr, T. G., et al. (2007). The Shuttle Radar Topography Mission. *Reviews of Geophysics*, 45(2), RG2004.
- Freeman, P. K., & Tukey, J. W. (1950). Transformations related to the angular and the square root. *The Annals of Mathematical Statistics*, 21(4), 607–611.
- Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1), 119–139.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- Hussain, E., Wright, T. J., Walters, R. J., Bekaert, D. P. S., Lloyd, R., & Hooper, A. (2018). Constant strain accumulation rate between major earthquakes on the Chaman fault. *Nature Communications*, 9(1), 1–9.
- Jackson, J., & McKenzie, D. (1984). Active tectonics of the Alpine-Himalayan Belt between western Turkey and Pakistan. *Geophysical Journal International*, 77(1), 185–264.
- Khan, M. A., Khan, S. D., & Walker, D. J. (2019). Structural evolution of the western Himalaya: Implications for active tectonics in Pakistan. *Journal of Asian Earth Sciences*, 177, 1–15.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444.
- Molnar, P., & Tapponnier, P. (1975). Cenozoic tectonics of Asia: Effects of a continental collision. *Science*, 189(4201), 419–426.
- Rehman, S. U., Shah, S. I., & Khan, M. A. (2021). Active fault mapping and seismic hazard assessment in western Pakistan. *Journal of Asian Earth Sciences*, 214, 104789.
- Shah, A. A., Ali, N., & Hussain, M. (2024). Integrated geophysical analysis for tectonic characterization of northern Pakistan. *Pure and Applied Geophysics*, 181(2), 567–584.
- Ullah, S., Li, H., Ashraf, U., & Asad, M. (2023). Machine learning-based multi-hazard susceptibility mapping in northern Pakistan. *Scientific Reports*, 13, 15234.
- Waldhauser, F., & Ellsworth, W. L. (2000). A double-difference earthquake location algorithm. *Bulletin of the Seismological Society of America*, 90(6), 1353–1368.
- Zhu, X. X., & Bamler, R. (2014). Super-resolution power and robustness of compressive sensing for spectral estimation with application to spaceborne tomographic SAR. *IEEE Transactions on Geoscience and Remote Sensing*, 52(6), 3369–3383.