

# Gold Mineral Prospectivity Mapping Using the Ensemble Model Approach, Case Study of Northwestern Thanh Hoa Province, Vietnam

Truong, X. Q.,<sup>1\*</sup> Tran, T. L.,<sup>2</sup> Dang, T. C.,<sup>3</sup> Tran, V. A.<sup>4</sup> and Truong, X. L.<sup>5</sup>

<sup>1</sup>School of Interdisciplinary Sciences and Arts, Vietnam National University, Hanoi, Faculty of Architecture, Urban Design and Sustainable Sciences, Building G7, Vietnam National University, Hanoi, 144 Xuan Thuy, Cau Giay, Hanoi, Vietnam, E-mail: txquang@vnu.edu.vn\*

<sup>2</sup>Graduate School of Science, Department of Geosciences, Geoinformatics lab, 558-8585 Osaka, Sumiyoshi Ward, Sugimoto 3 Chome-3-138, Japan, E-mail: sp22871@st.omu.ac.jp

<sup>3</sup>Faculty of Information Technology, Hanoi University of Natural Resources and Environment, 41A Phu Dien Street, Phu Dien Ward, Hanoi, Vietnam, E-mail: dtchien@hunre.edu.vn

<sup>4</sup>Hanoi University of Mining and Geology, Department of Photogrammetry and Remote Sensing, 18 Vien, Duc Thang, Bac Tu Liem, 100000 Hanoi, Viet Nam, E-mail: tranvananh@humg.edu.vn

<sup>5</sup>Hanoi University of Mining and Geology, Faculty of Information Technology, 18 Pho Vien, Duc Thang, Bac Tu Liem, 100000 Hanoi, Vietnam, E-mail: txluan@gmail.com

\*Corresponding Author

DOI: <https://doi.org/10.52939/ijg.v21i9.4437>

## Abstract

*In Vietnam, gold is considered a significant mineral resource. Although gold mining activities in Thanh Hoa province contribute economically, assessing and forecasting their spatial distribution continues to pose difficulties. Mineral prospectivity mapping (MPM) is crucial for investigating, surveying, planning, and managing natural resource exploitation, including gold deposits. Recently, the machine learning models have yielded compelling results. Although many individual machine learning models have been successfully applied, challenges in MPM remain due to data limitations and the complex nonlinear relationships between existing factors and deposits. To address the above challenges, this paper presents the first application of machine learning in general, and ensemble models in particular, for building mineral prospectivity maps (MPM) in the study area. An ensemble model integrating Random Forest (RF), Support Vector Machine (SVM), and XGBoost with ten selected conditioning factors was employed to enhance predictive accuracy. The study investigates the potential of stacking ensemble learning methods for MPM using a dataset of 438 points, consisting of 219 gold placer sampling sites and 219 non-deposit sites, divided into 70% for training and 30% for testing. The prediction model results using the Receiver Operating Characteristic (ROC) curve, with Area Under the Curve (AUC) values for RF, SVM, and XGBoost and Ensemble at 0.83, 0.87, 0.81 and 0.93. Compared with three single methods, stacking ensemble had the highest AUC. The result provides a statistical approach for constructing mineral prospectivity map (12.5-meter resolution) at the regional scale using geological, geophysical, and remote sensing data.*

**Keywords:** Ensemble Model, Gold Mineral Potential Mapping, MPM, Northwestern Thanh Hoa Province

## 1. Introduction

Mineral Prospectivity Mapping (MPM) plays a vital role in guiding exploration and optimising resource exploitation. Recently, machine learning methods have been widely employed to improve MPM models, with promising outcomes in applications such as gold deposit prediction and planning. These models rely on datasets that can be processed, analysed, and visualised through computer-based and GIS techniques [1].

Various machine learning methods have been employed to generate MPM. These include Logistic Regression and the weights-of-evidence method, as suggested by [2][3] and [4]; artificial neural networks (ANN) [5]; and support vector machines [1]. Another popular machine learning algorithm, Random Forest, has been successfully applied in several mineral exploration projects [5][6] and [7].

Predictive models built with XGBoost showed slightly better performance than those using Random Forest in generating optimal MPM results [8]. Recent methodological advancements in MPM have increasingly focused on deep learning algorithms, such as Convolutional Neural Networks [9][10][11] and [12]. A deep learning model has become a powerful tool for mineral exploration targeting in recent years. However, it requires accurate data. For example, geological image features are small and irregular, the image similarity is high, and the degree of influence of geological prospecting factors from different data sources on ore mineralisation varies [13].

In general, the types of data required for MPM mainly include (a) geological data, (b) geophysical data, (c) geochemical data, and (d) remote sensing data [10] and [14]. In addition to these data types, specific details and comprehensive information about known mineral or point deposits can enhance the accuracy of mineral deposit predictions. Numerous studies have addressed the issue of data imbalance resulting from the limited number of positive samples in geological datasets [13][14] and [15]. The finiteness and unbalance of geological exploration data often lead to large model errors or strong overfitting characteristics in machine learning models.

Recently, ensemble learning techniques that combine multiple base learners to improve predictive accuracy have been applied to classification problems, demonstrating superior performance over individual. Due to their effectiveness, these methods have gained rapid popularity. According to estimates, approximately 7,160 articles referenced “ensemble learning” in 2021 [16]. They have been successfully applied in various domains, such as landslide susceptibility mapping [17], soil organic matter content estimation [18], and mineral potential mapping models [19] and [20]. The study area is located in a remote and sparsely populated region, characterized by complex topography, dense vegetation, and limited accessibility. The population consists predominantly of ethnic minority communities, and the deeply concealed mineralisation is associated with all existing geological formations, particularly those related to fault activities. Geological and mineral data remain limited and fragmented; to date, only small-scale geological maps (1:200,000 and 1:50,000) have been produced. The geological structure of the region is highly complex [21] and [22]. Given the limited availability of geological data and the extensive area under consideration, Vietnamese authorities have

identified the need to prioritise specific regions for further exploration. Accordingly, this study applies machine learning models to predict mineral prospectivity, providing a timely and necessary approach to support more targeted and efficient resource assessment.

Due to the scarcity of geological and mineral data, as well as the lack of research applying machine learning to mineral resource assessment, progress in this field remains limited. To address the above challenges, we incorporated additional geological drilling data collected in 2021 and 2022, and this study presents the first application of machine learning in general, and ensemble models in particular, for building mineral prospectivity maps (MPM) in the study area. The proposed ensemble framework integrates three machine learning algorithms Random Forest (RF), Support Vector Machine (SVM), and Extreme Gradient Boosting (XGBoost) to develop an accurate MPM map that helps reduce computation time in forecasting and survey costs. The results demonstrate that applying this ensemble approach yields better predictive performance compared to the single-method techniques commonly employed in mineral resource assessment.

## 2. Study Area and Geological Setting

The study area, located in northwestern Thanh Hoa Province, spans about 653 km<sup>2</sup> with complex terrain and northeast–southwest trending ridges (Figure 1). Geologically, it lies within the Thanh Hoa structural zone, characterized by a wedge-shaped anticline bounded by the Mesozoic Son La and Sam Nua tectonic blocks. Most of the area belongs to the Paleozoic intracontinental tectonic belt of Northwestern Vietnam [21] and [22]. Of eight identified formations, seven are mineralized, including Ham Rong (€3-O1hr), Nam Pia (D1np), Ban Pap (D2bp), Toc Tat (D3tt), Bac Son (C-Pbs), Yen Duyet (P3yd), and Co Noi (T1cn) (Figure 2).

The Cam Thuy igneous complex is related with volcanic activities, relatively well developed in the area is of truly volcanic facies, composed of truly extrusive facies, consisting of porphyritic basalt, basalt diabase and explosive facies, consisting of basaltic tuff, basaltic breccia, basaltic tuff cobble and sub-volcanic formations in the form of veins, strata, stocks composed of gabbro, gabbro diabase and diabase porphyry related to gold mineralisation [23][24] and [25]. These magmatic rocks exhibit a close spatial and genetic relationship with gold mineralization.

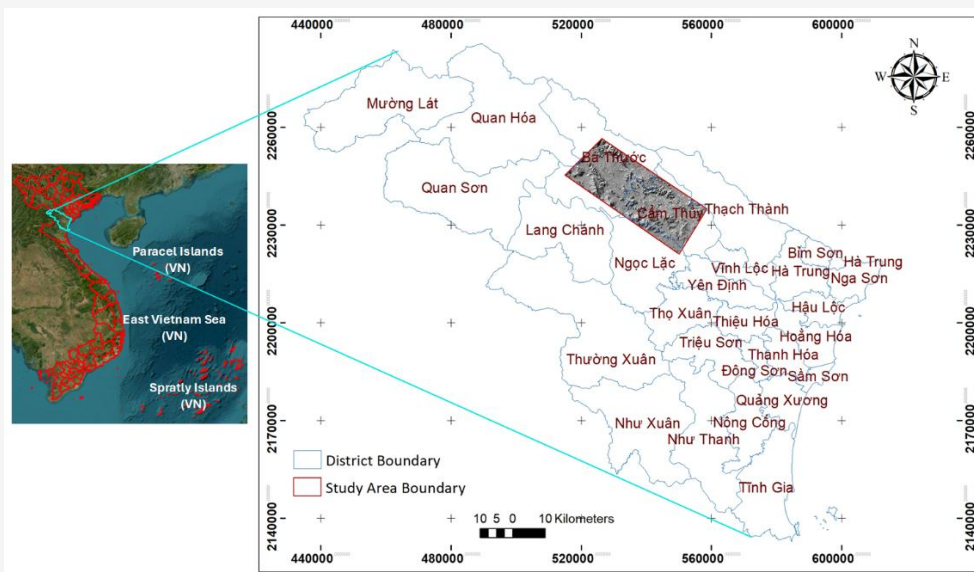


Figure 1: Northwestern Thanh Hoa province, Vietnam

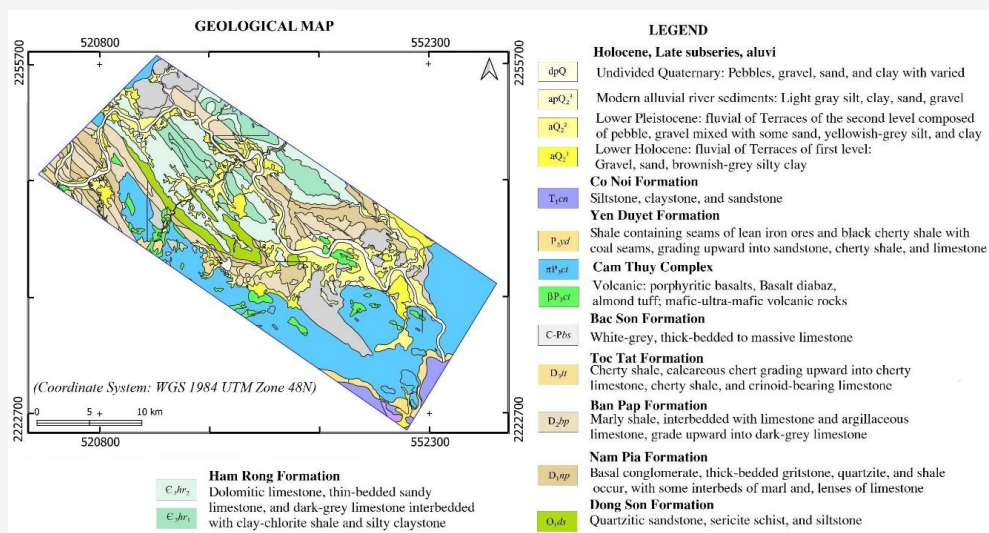


Figure 2: Geology of the study area (source: Department of Geology and Minerals of Vietnam)

The Cam Thuy magmatic complex comprises truly extrusive phases, such as basalt porphyry and basalt diabase; explosive facies, including basaltic tuff, basaltic breccia, and basaltic tuff conglomerate; as well as sub-volcanic formations occurring as veins, layers, and small stocks, composed of gabbro, gabbro-diabase, and diabase porphyry. These magmatic rocks are spatially and genetically associated with gold mineralization. Faulting activities in the study area are intense, multidirectional and complicated. The Northwest-Southeast fault system is dominant. Numerous plumose fractures and minor faults, filled with hydrothermal solutions, are observed along it.

The fault system controls the mineralisation processes. The Northeast-Southwest fault system is of smaller extent, mainly consisting of reverse faults with lateral displacements, creating well-developed cataclastic zones and shear fractures, which play the role as conductors for hydrothermal solutions and places for precipitation and accumulation of minerals. The sub-longitudinal and sub-latitudinal fault systems account for a small proportion, with not so deep dissection but with high intensity, causing great displacements. Gold ore pockets and lodes of significant value have formed at the intersections of this fault system and the aforementioned ones.

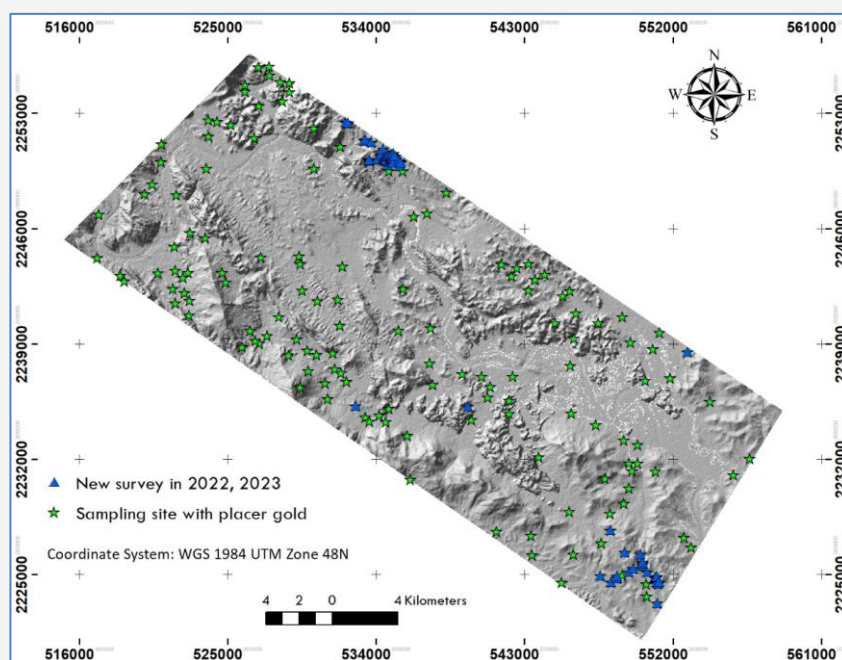
The main types Au-mineralisation in the area are summarised as follows: (1) Gold - quartz - low-sulfide mineralisation, distributed in the fissures and fissure zones of basalt and diabatic rocks and low-grade thermally metamorphosed sedimentary rocks; (2) Gold mineralisation in propylitic alteration zones, where propylite is distributed along fracture and fault systems, areal propylite coincides with explosion centers overlying basaltic tuff breccias; (3) Gold – antimony mineralisation: distributed within carbonate sedimentary units in the form of pockets and veins, with small-scale deposits; (4) Gold and polymetallic mineralisation, consists of stockwork veins and veinlets within limestone and calcareous shale; and (5) Carlin-like gold mineralisation, distributed in the carbonate sedimentary rocks.

Thus, it can be concluded that Geological periods, as well as magmatic, sedimentary, and metamorphic rocks, as well as structural features particularly minor structures are interrelated and play important but varying roles in the gold mineralisation process. However, gold mineralisation is mainly concentrated in zones of rock destruction (cataclastic zones), where porphyritic rocks occur in the form of sulfide-altered metamorphic rocks, along with gold-bearing hydrothermal sulfide veins and vein networks. Therefore, even small-scale tectonic activities play a significant role in ore formation. A total of 219 samples were collected and analyzed using various techniques. Among these, 142 samples contained placer gold. To enhance the dataset, an

additional 77 samples collected in 2022 and 2023 were incorporated and subsequently subjected to laboratory analyses, including chemical composition determination, fire assay, and atomic absorption spectroscopy (AAS) (Figure 3).

### 3. Data Used

In this study, the data layers generated from geological, geophysical, and remote sensing data related to gold mineralization were prepared to build ten evidential layers for MPM mapping. After processing with QGIS software, rasterized, converted to the 12.5-meter resolution and reclassified evidence layers, fault buffer, lineament buffer, lithology, magma, density of sampling sites with placer gold, placer aureole, bouguer gravity anomaly, magnetic anomaly, and mineral potential zones. In this context, geological age reflects the period associated with gold mineralization. Lithology determines the potential of rocks to host and accumulate gold, particularly in fractured zones. Magmatic activity supplies heat and metals, forming hydrothermal fluids enriched in gold that migrate through rocks or fault zones. Faults and lineaments provide pathways and act as structural traps conducive to the accumulation of gold-bearing hydrothermal fluids. Magnetic anomaly and Bouguer gravity anomaly data aid in identifying concealed structures such as faults, magmatic bodies, or metamorphic zones critical controls on gold deposition.



**Figure 3:** Sampling sites with gold placers distributed in the study area

**Table 1:** Data used

Map layers	Description
<b>Source:</b> Department of Geology and Minerals of Vietnam.	
1. Fault buffer (m)	The final map, generated from geological and petrographic-structural maps (1:50.000)
2. Magma	Generated from geological and petrographic-structural maps (1:50,000)
3. Lithology	Generate from the geological and petrographic-structural map (1:50.000).
4. Placer aureole	Generate from the geological and petrographic-structural map (1:50.000).
5. Bouguer gravity anomaly map (mGal)	Classification method: natural breaks with a 12.5-meter resolution
6. Magnetic anomaly map (nT)	Classification method: natural breaks with 12.5-meter resolution
7. Density of the sampling site with placer gold	The map was created using IDW (Inverse Distance Weighting)
8. Geological age	Generated from the geological map (1:50,000)
<b>Source:</b> Sentinel-1 and Sentinel-2, resolution (10 m) Date: 19-Nov2023 ( <a href="https://browser.dataspace.copernicus.eu/">https://browser.dataspace.copernicus.eu/</a> )	
9. Lineament buffer (m)	Generate from geological, petrographic-structural maps (1:50.000) and Remote sensing interpreter in both optical and SAR images (12.5 m).
10. Mineral potential zones	Implemented based on the Sentinel-2 image (RS); classification method: natural breaks with 12.5 meter resolution.

The distribution of placer aureoles and the density of sampling sites containing placer gold are closely related to areas with favorable conditions for gold accumulation. Finally, the mineral prospectivity layer was delineated using alteration indices derived from Sentinel-2 imagery. All ten map layers were converted to raster format with a spatial resolution of 12.5 meters (Table 1) and (Figure 4).

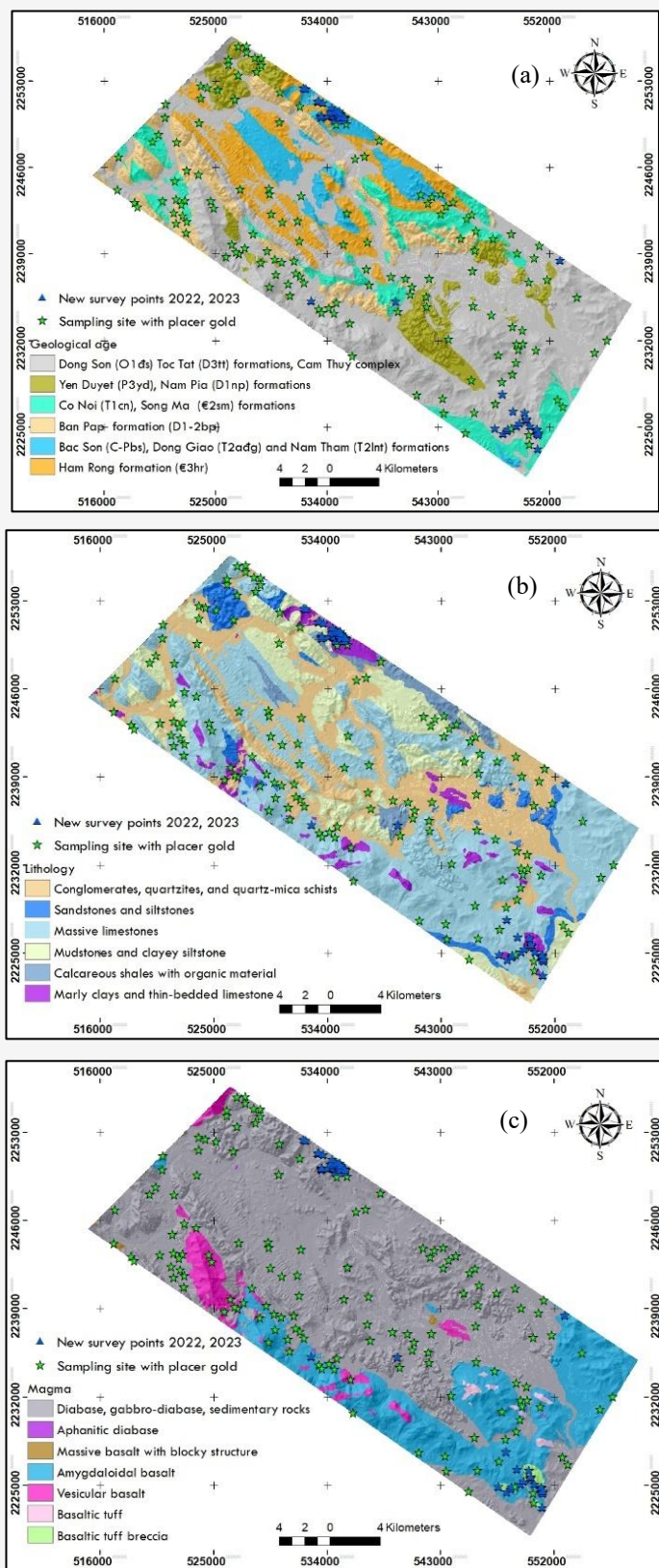
The selection of the 219 non-deposit samples is based on geological features that are unfavourable for gold mineralisation. The collected samples originated from intact, undeformed rocks with no evidence of small-scale fractures in various orientations, which significantly limits the potential for developing hydrothermal ore-bearing veins. These samples also comprise unmetamorphosed rock types and lithologies from volcanic craters. The identified lithologies are competent rocks resistant to brittle deformation and dynamic metamorphism, characterised by low joint density, such as quartzitic limestone, quartzitic sandstone, massive limestone, and white-grey cherty shale. Rock types that are not conducive to contact metamorphism or metasomatic alteration at medium to low temperatures, such as quartzite, dolomite, massive limestone, and dense basalt; Regions situated at significant distances from eruption centres, volcanic craters, and major tectonic faults trending northwest–southeast.

No uniform guideline exists for splitting the data [26]. In this study, the dataset was partitioned into a training subset comprising 70% of the data and a testing subset consisting of the remaining 30%. The split was stratified according to class, ensuring a

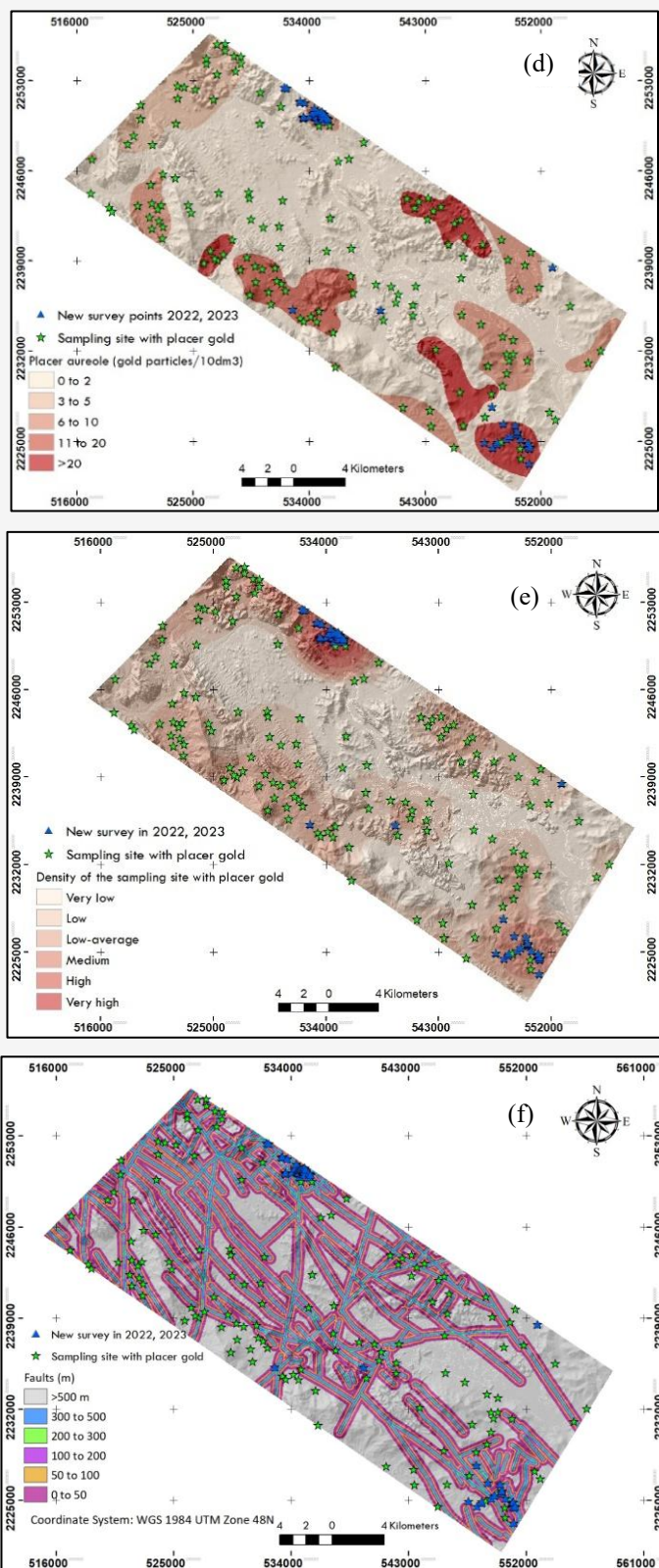
balanced representation of points with gold placer and non-deposit points in the training and testing sets. A random state was specified to ensure the reproducibility of the models. The dataset includes 438 points with 219 sampling sites with gold placer and 219 with non-deposit. Pixel values from the ten raster layers were extracted at the locations of all 438 points in the dataset. Finally, the extracted features and their corresponding labels (gold placer = 1, non-gold = 0) were used to construct the training and testing datasets (Figure 3).

#### 4. Methodology

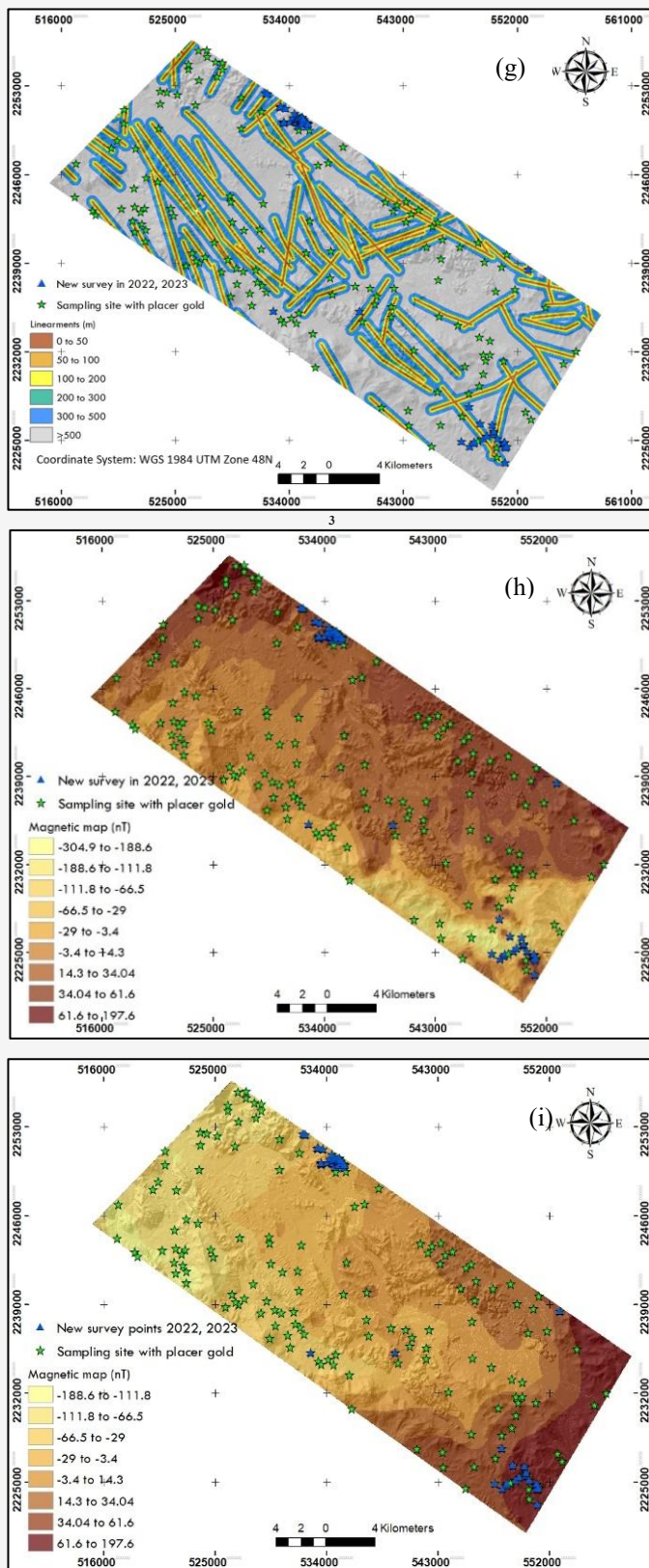
The Random Forest (RF) model is a supervised ensemble learning algorithm that constructs multiple decision trees using random subsets of both data samples and features, forming a "decision forest" Each tree contributes a vote for the predicted class, and the majority vote determines the final prediction [7] and [27]. It is widely used in mineral prospectivity mapping (MPM) studies because it can effectively handle non-linear relationships and integrate input data from multiple sources such as: remote sensing, geophysics, and geology [28]. This built-in randomness enhances model robustness and helps prevent overfitting through the bagging mechanism. RF is well known for its high accuracy in binary classification tasks and is computationally efficient, often requiring shorter training times than other supervised methods [29] and [30]. Additionally, RF can handle low-quality data and missing values effectively, making it a versatile tool for various predictive modelling tasks [27].



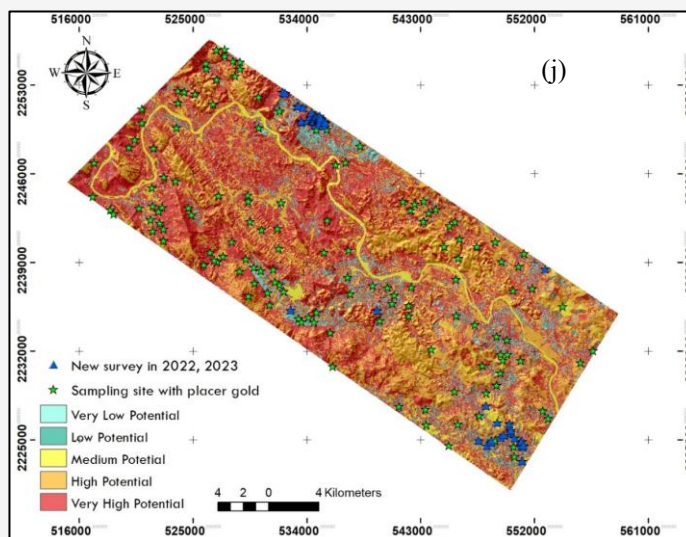
**Figure 4:** Ten map factors of the MPM using machine learning models  
(a) Geological age, (b) Lithology, (c) Magma (Continue next page)



**Figure 4:** Ten map factors of the MPM using machine learning models:  
 (d) Placer aureole, (e) Density of the sampling site with placer gold,  
 (f) Fault buffer (Continue next page)



**Figure 4:** Ten map factors of the MPM using machine learning models:  
 (g) Lineament density, (h) Magnetic anomaly map,  
 (i) Bouguer gravity anomaly map (Continue next page)



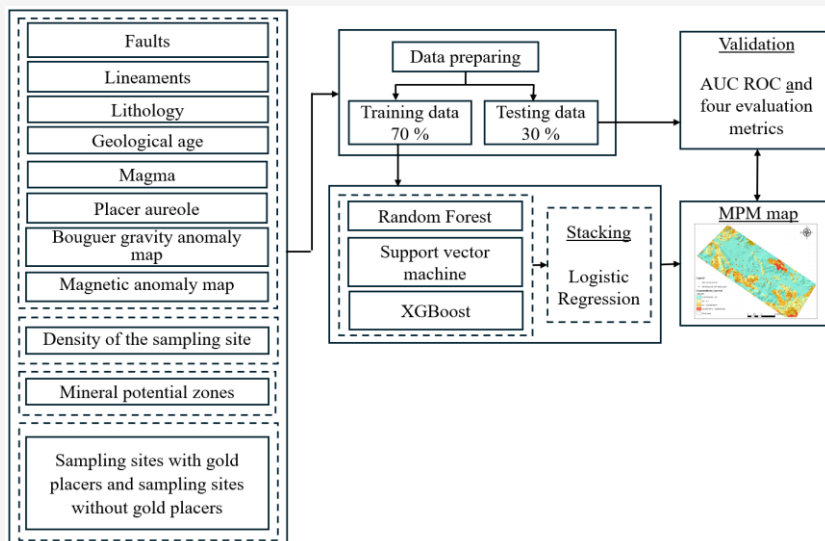
**Figure 4:** Ten map factors of the MPM using machine learning models:  
(j) Mineral potential zones (Continue from previous page)

The SVM algorithm is a well-established supervised machine learning method. It learns from training data and predicts outcomes by separating data into categories using hyperplanes, where the nearest points from each class are maximised. Support vectors are the points that lie closest to the hyperplane and influence its position. The SVM often generalises well and avoids overfitting. SVM can also be extended to handle nonlinear problems using kernel functions [31] and [29]. SVM also demonstrates robustness in handling datasets characterized by a limited sample size and collected from multiple sources. Although it does not directly provide the importance of variables, it can model complex boundaries between mineralized and non-mineralized areas through the use of kernel functions [28].

Gradient Boosting constructs an ensemble of weak learners, where each subsequent model sequentially corrects the errors of its predecessors [30]. This approach requires selecting an appropriate differentiable loss function and is valued for its flexibility to accommodate various loss functions without designing new algorithms [31]. An optimized implementation, known as XGBoost, enhances this framework by incorporating regularization to mitigate overfitting and by significantly improving computational efficiency through row and column sampling, parallel computing, and optimized data structures [32]. Instead of repeatedly sorting feature values during tree construction, it stores pre-sorted blocks of data to reduce both time and memory costs [33]. These characteristics make the method particularly suitable for mineral potential mapping, where datasets often combine heterogeneous variables from multiple

sources and require efficient handling of large spatial datasets.

Ensemble models are machine learning techniques that combine multiple algorithms to improve predictive performance compared with a single model [17]. This approach mitigates variance and bias, enhancing accuracy by leveraging the complementary strengths of individual learners, thereby producing more robust and generalisable outcomes. Common strategies include Bagging, Boosting, Voting, and Stacking. In homogeneous settings, such as Bagging, multiple identical models are trained on different random subsets of data and their predictions are aggregated, with Random Forest being a well-known example. Boosting, on the other hand, improves accuracy by sequentially training models, where each model corrects the errors of its predecessor, as seen in AdaBoost and XGBoost. In contrast, a heterogeneous can be implemented using Voting and Stacking. Voting combines several models and makes decisions based on majority voting (hard voting) or probability averaging (soft voting), and is commonly used in classification problems. Stacking takes the results predicted by the base-classifiers as the input attributes, and the meta-learner merges the different predictions into the final prediction [18] and [34]. The objective of this paper is to integrate models within a heterogeneous ensemble, a collection of machine learning algorithms with distinct structures in order to leverage their complementary strengths. Among various ensemble approaches, this study applies a stacking framework that combines three algorithms widely recognized for their robust performance in statistical analysis: Random Forest (RF), XGBoost, and Support Vector Machine (SVM).



**Figure 5:** The training process of the stacking ensemble learning method framework

Figure 5 illustrates the main steps for implementing mineral prospectivity mapping in the study area, including:

(1) The dataset for the ensemble model was pre-processed using QGIS software and converted to a 12.5-m resolution (Table 1).

(2) The processed dataset comprises 219 sampling sites containing gold placers and 219 without gold placers, with a 70%–30% split for model training and testing, respectively.

(3) The modelling phase employs RF, SVM, XGBoost, and their ensemble model for creating MPM. To enhance predictive performance, a meta-model that integrates the strengths of these algorithms is utilized, with logistic regression chosen as the meta-learner, given its widespread application in ensemble construction and final prediction generation [19] and [20].

(4) Evaluating models using AUC curves and four statistical metrics.

### 5. Model Evaluation Metrics

The Area Under the Receiver Operating Characteristic (ROC) curve (AUC) is a widely used metric for assessing model quality, with values ranging from 0.5 to 1.0. Higher AUC values indicate stronger predictive performance [35]. In addition, four fundamental statistical measures accuracy, precision, recall and F1-score, are commonly applied to evaluate the classification performance of the MPM.

These metrics are derived from the confusion matrix, which summarizes the model's predictions by showing the number of correct and incorrect classifications across all classes. The ROC curve and its corresponding AUC value quantify the model's ability to discriminate between classes. The combined use of these evaluation methods offers a comprehensive assessment of the algorithm's effectiveness, particularly in terms of overall accuracy and performance on imbalanced datasets [14]. The statistical measures are defined in Equations 1 to 4:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad \text{Equation 1}$$

$$Recall = \frac{TP}{TP + FN} \quad \text{Equation 2}$$

$$Precision = \frac{TP}{TP + FP} \quad \text{Equation 3}$$

$$F1\ score = \frac{2Precision \times Recall}{Precision + Recall} \quad \text{Equation 4}$$

A true positive (*TP*) occurs when a deposit sample is correctly identified as “deposit,” whereas a false negative (*FN*) refers to a deposit sample that is mistakenly classified as “non-deposit.” Conversely, a true negative (*TN*) is a non-deposit sample accurately recognized as “non-deposit,” while a false positive (*FP*) denotes a non-deposit sample incorrectly labeled as “deposit.”

Using these definitions, the ROC curve is constructed by plotting the true positive rate (TPR) along the y-axis against the false positive rate (FPR) along the x-axis. The closer the ROC curve lies to the upper-left corner, the stronger the classification capability of the model. Accuracy, as presented in Equation (Equation 1), is a widely used and straightforward metric. However, in cases where the data distribution between sampling sites with and without gold placers is imbalanced, accuracy may be insufficient and potentially misleading in evaluating the performance of an ensemble model. Recall (Equation 2), also known as sensitivity, quantifies the proportion of actual positive instances that are correctly identified by the model. Precision (Equation 3) measures the proportion of predicted positive instances that are truly positive. The F1-score (Equation 4) is utilized to represent the balance between precision and recall, which are inherently in trade-off, by calculating their harmonic mean.

## 6. Results and Discussion

The MPM outcome, including model training and evaluation, was implemented in Python using the scikit-learn package [35], which provides extensive resources for machine learning. Within the ensemble learning framework, parameter and hyperparameter tuning was systematically performed via GridSearchCV with k-fold cross-validation [36], enabling an exhaustive search for the optimal configuration. The key parameters of each model are summarised in Table 2. The stacking ensemble model was structured in two layers: outputs from the three base classifiers (RF, SVM, and XGBoost) served as input features for a higher-level predictive model. To enhance stability and generalizability, a k-

fold cross-validation strategy was applied during training. Logistic Regression (LR) was adopted as the meta-learner to generate the final prediction. The performance of the stacking ensemble was evaluated against the individual classifiers using the independent testing dataset. The final MPM was subsequently reclassified into four potential levels: low, moderate, high, and very high. To determine the optimal threshold for zones with very high potential, the researchers tested several threshold values, assessed their performance, and selected the one that yielded the highest F1 score. The MPM was also constructed with four classes: low potential (0.0–0.4), moderate (0.4–0.6), high potential (0.6–0.8), and very high potential (0.8–1). The predictive capacity was evaluated using the testing dataset, which provides an independent measure of model performance.

The results of the ensemble model and its base classifiers were evaluated using the statistical metrics presented in Table 3. The combination of RF–SVM–XGBoost outperformed the individual classifiers. This integrated model provided the highest accuracy (0.86) and achieved an AUC of 0.93, compared with 0.83 for RF, 0.86 for SVM, and 0.81 for XGBoost (Figure 6). The validation results demonstrate that MPM can reduce the uncertainty associated with multiple variables in mineral potential modelling. The ensemble classified 2% of the study area as very high potential, while high, moderate, and low potential zones accounted for 15%, 31%, and 52% of the total area, respectively (Table 3). The Random Forest model achieved an overall accuracy of 0.80, classifying 11% of the study area as belonging to the very high-potential zone.

**Table 2:** The optimised parameters of the models

Methods	Parameters
Random Forest	n_estimators = 200; criterion = 'gini'; max_depth = None; max_features = 'sqrt'
XGBoost	gamma: 0.3, learning_rate: 0.01, max_depth: 3, n_estimators: 150
SVM	C=1, gamma='scale', kernel='sigmoid', tol=0.0001

**Table 3:** Performance of the machine learning models for MPM

Metrics	Random Forest	SVM	XGBoost	Stacking
TP (%)	42.50	45.00	41.00	45.8
TN (%)	37.00	38.00	38.00	40.1
FP (%)	8.50	5.20	9.50	4.2
FN (%)	12.00	11.80	11.50	9.9
Accuracy	0.80	0.83	0.79	0.86
Precision	0.85	0.90	0.82	0.91
Recall	0.78	0.80	0.78	0.82
F1-Score	0.81	0.85	0.80	0.86

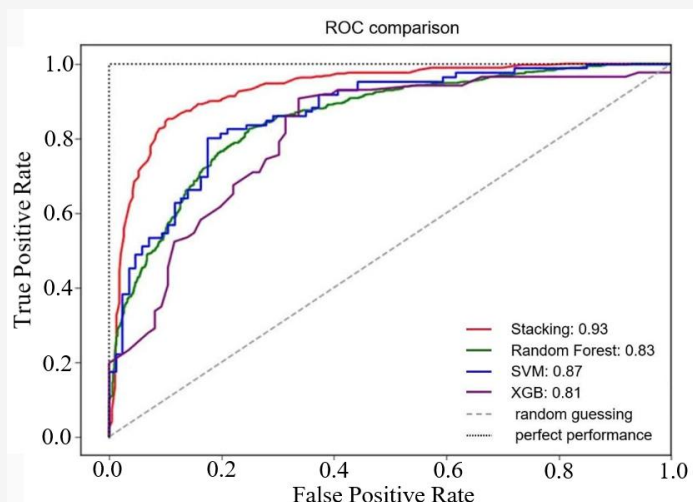


Figure 6: AUC curves of the MPM models

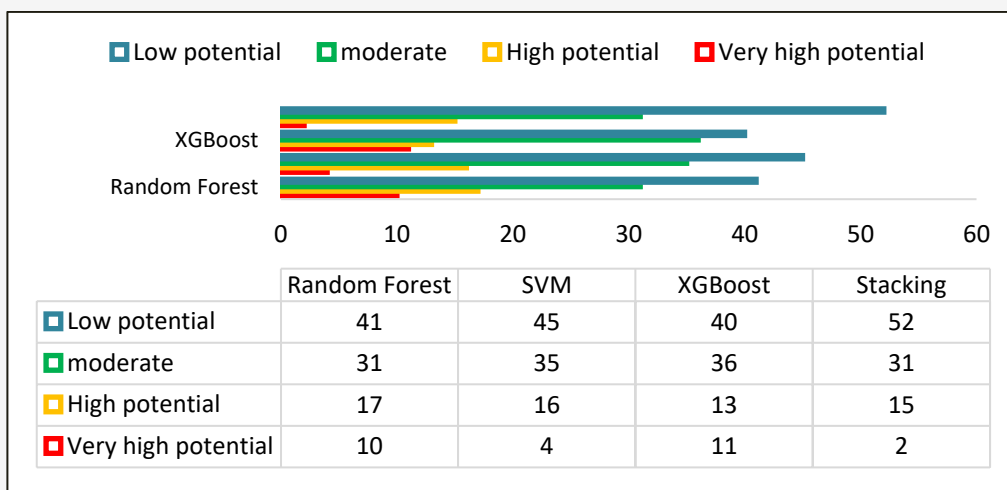


Figure 7: The total area of the predicted classes expressed as percentages

The Support Vector Machine (SVM) reached a slightly higher accuracy of 0.83, but identified only 4% of the area as very high potential. The XGBoost model, with an accuracy of 0.79, delineated 10% of the study area as having very high gold potential (Table 3) and (Figure 7).

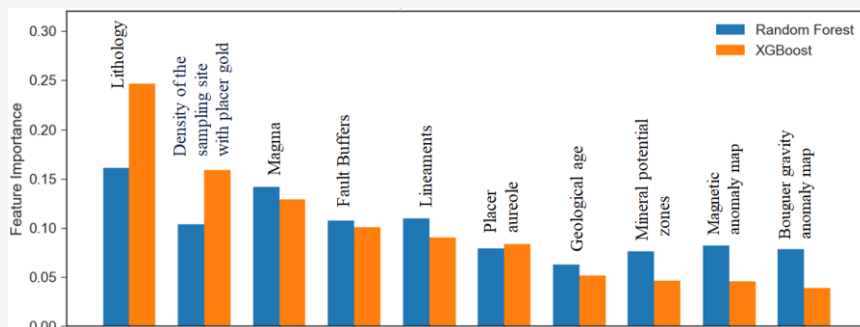
In ensemble learning, the ability to assess the importance of input parameters (feature importance) depends on the type of algorithm employed. Among them, Random Forest and XGBoost are capable of interpreting the importance of the model's input features [37] and [38]. XGBoost assigns greater importance to lithology and density of the sampling site with placer gold, while Random Forest distributes importance more evenly. For intermediate features such as Magma, Fault Buffers, and Lineaments, both models show quite similar evaluations. In contrast, for less influential features

like Placer aureole, Geological age, Mineral potential zones, Magnetic anomaly map, and Bouguer gravity anomaly map, Random Forest maintains stable values (0.07–0.09), whereas XGBoost significantly reduces their importance. This indicates that XGBoost emphasizes a few dominant features, whereas Random Forest provides a more balanced distribution across variables (Figure 8).

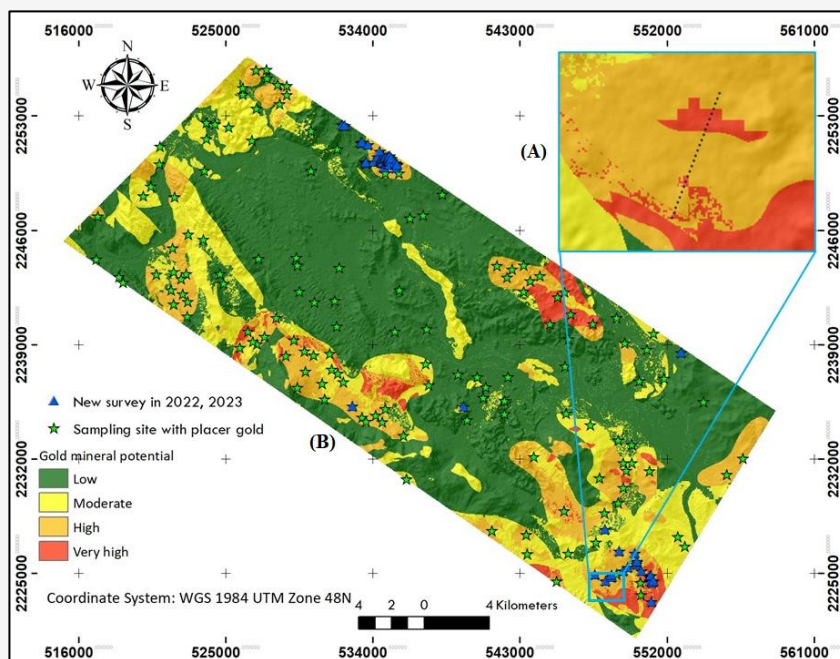
In addition to assessing the reliability of the ensemble model using the specified algorithmic method as above, the predictive mapping results were validated against actual conditions in the study area. Specifically, most of the predicted prospective gold mineralisation zones coincided or nearly coincided with ancient volcanic vents, thermal contact metamorphic zones, and fault intersections, generally extending in a northwest-southeast direction.

These prospective areas also overlap with or are located near known ore occurrences. To further verify the results, the authors randomly selected several points predicted to be prospective but that had not yet been surveyed or recorded by geologists, and conducted telluric geophysical measurements. Specifically, 25 measurement points were carried out in the Doi Vo area, Cam Tam commune, Cam Thuy district, Thanh Hoa province. The results showed that these locations had low resistivity values, fluctuating around  $\leq 200$  ohm·m, which are associated with hydrothermal alteration zones, including sericitisation, dolomitisation, and pyritisation, containing gold mineralisation veins. This outcome provides additional evidence that the predictive map is reliable. These prediction results provide valuable information for future gold mineral exploration, especially for locations predicted to have high or very

high gold mineral potential that geologists have not investigated or documented (Figure 9(a)). Machine learning models in general and ensemble models in particular have been applied for the first time in Vietnam through this study. The research area is characterized by complex terrain, limited accessibility, and a lack of detailed geological and mineral resource data. However, it is considered to hold significant potential for gold mineralization. Therefore, the outcomes of this study are expected to have important implications for future mineral exploration in Vietnam. Such a meta-learning approach can provide accurate guidance for future gold mineralisation surveys, particularly in highly prospective areas identified on the prediction map but not yet recognised by previous geologists (Figure 9(b)).



**Figure 8:** Feature importances of the Random Forest and XGBoost models



**Figure 9:** (a) Spatial distribution of 25 measurement points and the resulting MPM map  
(b) Gold mineral potential map for the study area based on Random Forest

To optimize time and costs, priority should be given to areas classified as having very high potential for initial exploration activities, followed by those with high and very high potential. Areas with moderate potential should be considered for long-term exploration plans, while those assessed as low potential are not recommended for investment. Given the relatively small exploration area, drilling and telluric (geophysical) methods are the primary approaches proposed for further investigation.

## 7. Conclusion

The stacking ensemble model integrating RF, SVM, and XGBoost surpassed the individual classifiers, demonstrating its effectiveness in mineral prospectivity mapping. This approach reduced uncertainties in mineral potential modelling and provided a more robust framework for identifying prospective zones. Validation with geological evidence and field geophysical surveys confirmed the reliability of the predictive map. The study represents the first application of ensemble machine learning models for mineral prospectivity mapping in Vietnam, offering new insights into the gold potential of Thanh Hoa Province. Overall, the findings provide a valuable methodological basis to guide and optimize future mineral exploration in complex geological settings. Future work may extend this framework by incorporating deep learning architectures and higher-resolution remote sensing datasets to further refine prediction accuracy. In addition, conducting geochemical surveys and data collection will be necessary to support the development of more comprehensive and sustainable exploration strategies.

## Acknowledgments

This work is funded by the Ministry of Science and Technology of Vietnam (MOST), with project code number ĐTĐL.CN-85/21. The authors thank the anonymous reviewers for improving the paper.

## References

- [1] Zuo, R. and Carranza, E. J. M., (2011). Support Vector Machine: A Tool for Mapping Mineral Prospectivity. *Computers and Geosciences*, Vol. 37(12), 1967-1975. <https://doi.org/10.1016/j.cageo.2010.09.014>.
- [2] King, G. and Zeng, L., (2017)., Logistic Regression in Rare Events Data. *Political Analysis*, Vol. 9(2), 137-163. <https://doi.org/10.1093/oxfordjournals.pan.a004868>.
- [3] Xiong, Y. and Zuo, R., (2018)., GIS-based Rare Events Logistic Regression for Mineral Prospectivity Mapping. *Computers and Geosciences*, Vol. 111, 18-25. <https://doi.org/10.1016/j.cageo.2017.10.005>.
- [4] Zhao, J., Sui, Y., Zhang, Z. and Zhou, M., (2023). Application of Logistic Regression and Weights of Evidence Methods for Mapping Volcanic-Type Uranium Prospectivity. *Minerals*, Vol. 13(5). <https://doi.org/10.3390/min13050608>.
- [5] Sun, T., Chen, F., Zhong, L. X., Liu, W. M. and Wang, Y., (2019). GIS-based Mineral Prospectivity Mapping Using Machine Learning Methods: A Case Study from Tongling Ore District, Eastern China. *Ore Geology Reviews*, Vol. 109, 26–49. <https://doi.org/10.1016/j.oregeorev.2019.04.003>.
- [6] Chen, M. and Xiao, F., (2023). Projection Pursuit Random Forest for Mineral Prospectivity Mapping. *Mathematical Geosciences*, Vol. 55, 963–987. <https://doi.org/10.1007/s11004-023-10070-0>.
- [7] Zhang, Z., Zuo, R. and Xiong, Y. A., (2016). Comparative Study of Fuzzy Weights of Evidence and Random Forests for Mapping Mineral Prospectivity for Skarn-Type Fe Deposits in the Southwestern Fujian Metallogenic Belt, China. *Science China Earth Sciences*, Vol. 59, 556–572. <https://doi.org/10.1007/s11430-015-5178-3>.
- [8] Parsa, M., (2021). A Data Augmentation Approach to XGboost-based Mineral Potential Mapping: An Example of Carbonate-Hosted ZnPb Mineral Systems of Western Iran, *Journal of Geochemical Exploration*, Vol. 228. <https://doi.org/10.1016/j.gexplo.2021.106811>.
- [9] Li, Y. S., Peng, C., Ran, X. J., Xue, L. F. and Chai, S. L., (2021). Soil Geochemical Prospecting Prediction Method Based on Deep Convolutional Neural Networks-Taking Daqiao Gold Mine in Gansu Province, China as an Example. *China Geology*, Vol. 4(3), 1-14. <https://doi.org/10.31035/cg2021044>.
- [10] Li, Q., Chen, G. and Luo, L., (2023). Mineral Prospectivity Mapping Using Attention-Based Convolutional Neural Network, *Ore Geology Reviews*, Vol. 156. <https://doi.org/10.1016/j.oregeorev.2023.105381>.
- [11] Xiong, Y. and Zuo, R., (2020). Recognizing Multivariate Geochemical Anomalies For Mineral Exploration By Combining Deep Learning And One-Class Support Vector Machine. *Computers and Geosciences*, Vol. 140. <https://doi.org/10.1016/j.cageo.2020.104484>.

- [12] Wang, Z. and Zuo, R., (2022). Mineral Prospectivity Mapping Using a Joint Singularity-Based Weighting Method and Long Short-Term Memory Network. *Computers and Geosciences*, Vol. 158. <https://doi.org/10.1016/j.cageo.2021.104974>.
- [13] Gao, L., Wang, K., Zhang, X. and Wang, C., (2023). Intelligent Identification and Prediction Mineral Resources Deposit Based on Deep Learning. *Sustainability*, Vol. 15(13). <https://doi.org/10.3390/su151310269>.
- [14] Sun, K., Chen, Y., Geng, G., Lu, Z., Zhang, W., Song, Z., Guan, J., Zhao, Y. and Zhang, Z., (2024). A Review of Mineral Prospectivity Mapping Using Deep Learning. *Minerals*, Vol. 14. <https://doi.org/10.3390/min14101021>.
- [15] Chen, Y. and Shayilan, A., (2022). Dictionary Learning for Multivariate Geochemical Anomaly Detection for Mineral Exploration Targeting. *Journal of Geochemical Exploration*, Vol. 235. <https://doi.org/10.1016/j.gexplo.2022.106958>.
- [16] Mohammed, A. and Kora, R., (2023). A Comprehensive Review on Ensemble Deep Learning: Opportunities and Challenges. *Journal of King Saud University - Computer and Information Sciences*, Vol. 35(2), 757-774. <https://doi.org/10.1016/j.jksuci.2023.01.014>.
- [17] Lee, S. M. and Lee, S. J., (2024). Landslide Susceptibility Assessment of South Korea Using Stacking Ensemble Machine Learning. *Geoenvironmental Disasters*, Vol. 11(7). <https://doi.org/10.1186/s40677-024-00271-y>.
- [18] Wu, M., Dou, S., Lin, N., Jiang, R. and Zhu, B., (2023). Estimation and Mapping of Soil Organic Matter Content Using a Stacking Ensemble Learning Model Based on Hyperspectral Images. *Remote Sensing*, Vol. 15. <https://doi.org/10.3390/rs15194713>.
- [19] Zhang, Y., Li, M., Han, S., Ren, Q. and Shi, J., (2019). Intelligent Identification for Rock-Mineral Microscopic Images Using Ensemble Machine Learning Algorithms. *Sensors*, Vol. 19. <https://doi.org/10.3390/s19183914>.
- [20] Wang, K., Zheng, X., Wang, G., Liu, D. and Cui, N., (2020). A Multi-Model Ensemble Approach for Gold Mineral Prospectivity Mapping: A Case Study on the Beishan Region, Western China. *Minerals*, Vol. 10(12). <https://doi.org/10.3390/min10121126>.
- [21] Dovjikov, A. E., (1965). *Geology of Northern Vietnam* (in Vietnamese) (Ed.). Publishing House for Science and Technology, Ha Noi, Vietnam.
- [22] Tran, V. T., (2009). *Geology and Resources of Vietnam* (in Vietnamese). 1<sup>st</sup> Edition by Science and Technology Publishing House, Ha Noi, Vietnam.
- [23] Tran, V. T., and K. Vu. (2011). *Geology and Earth Resources of Vietnam*; General Department of Geology and Minerals of Vietnam. Publishing House for Science and Technology. Hanoi.
- [24] Tran, T. H., Le, T. X., Zaw, K. and Manaka, T., (2015). Structural Controls on Gold Mineralization in the Southeastern Truong Son Fold-Thrust Belt and Its Significance in Regional Metallogeny. *Proceedings of the PACRIM 2015 Congress*, Hong Kong, China, 521–532.
- [25] Hua, X., Zenga, G., Zhang, Z., Lid, W., Liua, W., Gong, Y. and Yao, S., (2018). Gold Mineralization Associated with Emeishan Basaltic Rocks: Mineralogical, Geochemical, and Isotopic Evidences from the Lianhuashan Ore Field, Southwestern Guizhou Province, China. *Ore Geology Reviews*, Vol. 95, 604-619. <https://doi.org/10.1016/j.oregeorev.2018.03.016>.
- [26] Kuhn, M. and Johnson, K., (2019). Feature Engineering and Selection: A Practical Approach for Predictive Models (1<sup>st</sup>Ed.). *Chapman and Hall/CRC.*, New York. <https://doi.org/10.1201/9781315108230>.
- [27] Breiman, L., (2001). Random Forests. *Machine Learning*, Vol. 45, 5–32. <https://doi.org/10.1023/A:1010933404324>.
- [28] Lachaud, A., Adam, M. and Mišković, I., (2023). Comparative Study of RandomForest and Support Vector Machine Algorithms in Mineral Prospectivity Mapping with Limited Training Data. *Minerals*, Vol. 13. <https://doi.org/10.3390/min13081073>.
- [29] Nguyen, D., Chou, T., Hoang, T., and Chen, M. (2023). Flood Susceptibility Mapping Using Machine Learning Algorithms: A Case Study in Huong Khe District, Ha Tinh Province, Vietnam. *International Journal of Geoinformatics*, Vol. 19(7), 1–15. <https://doi.org/10.52939/ijg.v19i7.2739>.
- [30] Friedman, J., Hastie, T. and Tibshirani, R., (2000). Additive Logistic Regression: A Statistical View of Boosting (with Discussion and a Rejoinder by the Authors). *The Annals of Statistics*, Vol. 28(2), 337-407. <https://doi.org/10.1214/aos/1016218223>.

- [31] Tran, V. A., Khuc, T. D., Truong, X. Q., Nguyen, A. B. and Phi, T. T., (2024). Application of potential Machine Learning Models in Landslide Susceptibility Assessment: A Case Study of Van Yen District, Yen Bai Province, Vietnam, *Quaternary Science Advances*, Vol. 14. <https://doi.org/10.1016/j.qsa.2024.100181>.
- [32] Chen, T. and Guestrin, C., (2016). XGBoost: A Scalable Tree Boosting System. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16). *Association for Computing Machinery*, New York, NY, USA, 785–794. <https://doi.org/10.1145/2939672.2939785>.
- [33] Zhang, W., He, Y., Wang, L., Liu, S. and Meng, X., (2023). Landslide Susceptibility Mapping Using Random Forest and Extreme Gradient Boosting: A Case Study of Fengjie, Chongqing. *Geological Journal*, Vol. 58(6), 2372–2387. <https://doi.org/10.1002/gj.4683>.
- [34] Tewari, S. and Dwivedi, U. D., (2020). A Comparative Study of Heterogeneous Ensemble Methods for the Identification of Geological Lithofacies. *Journal of Petroleum Exploration and Production Technology*, Vol. 10, 1849–1868. <https://doi.org/10.1007/s13202-020-00839-y>.
- [35] Fawcett, T., (2006). An Introduction to ROC Analysis. *Pattern Recognition Letters*, Vol. 27(8), 861–874. <https://doi.org/10.1016/j.patrec.2005.10.010>.
- [36] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Müller, A., Nothman, J., Louppe, G., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, É., (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, Vol. 12, 2825–2830.
- [37] Jain, R., Tripathi, N. K., Pant, M., Anutariya, C. and Silpasuwanchai, C., (2024). Investigating Gender and Age Variability in diabetes prediction: A Multi-Model Ensemble Learning Approach. *IEEE Access*, Vol. 12, 71535–71554. <https://ieeexplore.ieee.org/document/10534069>
- [38] Wen, H. T., Wu, H. Y. and Liao, K. C., (2022). Using XGBoost Regression to Analyze the Importance of Input Features Applied to an Artificial Intelligence Model for the Biomass Gasification System. *Inventions*, Vol. 7(4). <https://doi.org/10.3390/inventions7040126>.