
Short Note

Machine Learning Application in G.I.S. and Remote Sensing: An Overview

Anjeel Upreti

Capital College and Research Center, Kathmandu, Nepal
practiceriwaz691@gmail.com

Abstract: Machine learning (ML) is a subdivision of artificial intelligence in which the machine learns from machine-readable data and information. It uses data, learns the pattern and predicts the new outcomes. Its popularity is growing because it helps to understand the trend and provides a solution that can be either a model or a product. Applications of ML algorithms have increased drastically in G.I.S. and remote sensing in recent years. It has a broad range of applications, from developing energy-based models to assessing soil liquefaction to creating a relation between air quality and mortality. Here, in this paper, we discuss the most popular supervised ML models (classification and regression) in G.I.S. and remote sensing. The motivation for writing this paper is that ML models produce higher accuracy than traditional parametric classifiers, especially for complex data with many predictor variables. This paper provides a general overview of some popular supervised non-parametric ML models that can be used in most of the G.I.S. and remote sensing-based projects. We discuss classification (Naïve Bayes (NB), Support Vector Machine (SVM), Random Forest (RF), Decision Trees (DT)) and regression models (Random Forest (RF), Support Vector Machine (SVM), Linear and Non-Linear) here. Therefore, the article can be a guide to those interested in using ML models in their G.I.S. and remote sensing-based projects.

Keywords: machine learning; artificial intelligence; pattern; models; classification; regression; GIS; remote sensing

1. Introduction

Machine learning (ML) is a subdivision of artificial intelligence in which the machine learns from machine-readable data and information (Verma & Verma, 2021). It uses data, learns the pattern and predicts the new outcomes (Maxwell, Warner, & Fang, 2018). Its popularity is growing because it helps to understand the trend and provides a solution that can be either a model or a product. There are four types of machine learning approaches: supervised, unsupervised, semi-supervised and reinforcement learning (Sarker, 2021). In supervised learning, the labelled training data is provided; in unsupervised learning, unlabeled training data is provided (Sarker, 2021). The semi-supervised learning approach is a hybrid of both supervised and unsupervised learning where mostly labelled information is provided for the training (Sidey-Gibbons & Sidey-Gibbons, 2019). However, the model is free to figure out the trend in the data on its own. In reinforcement learning, the agent learns from trial and error to make decisions and cope with the interactive environment [4]. A ML project consists of several steps and each step should be planned carefully (Figure 1).

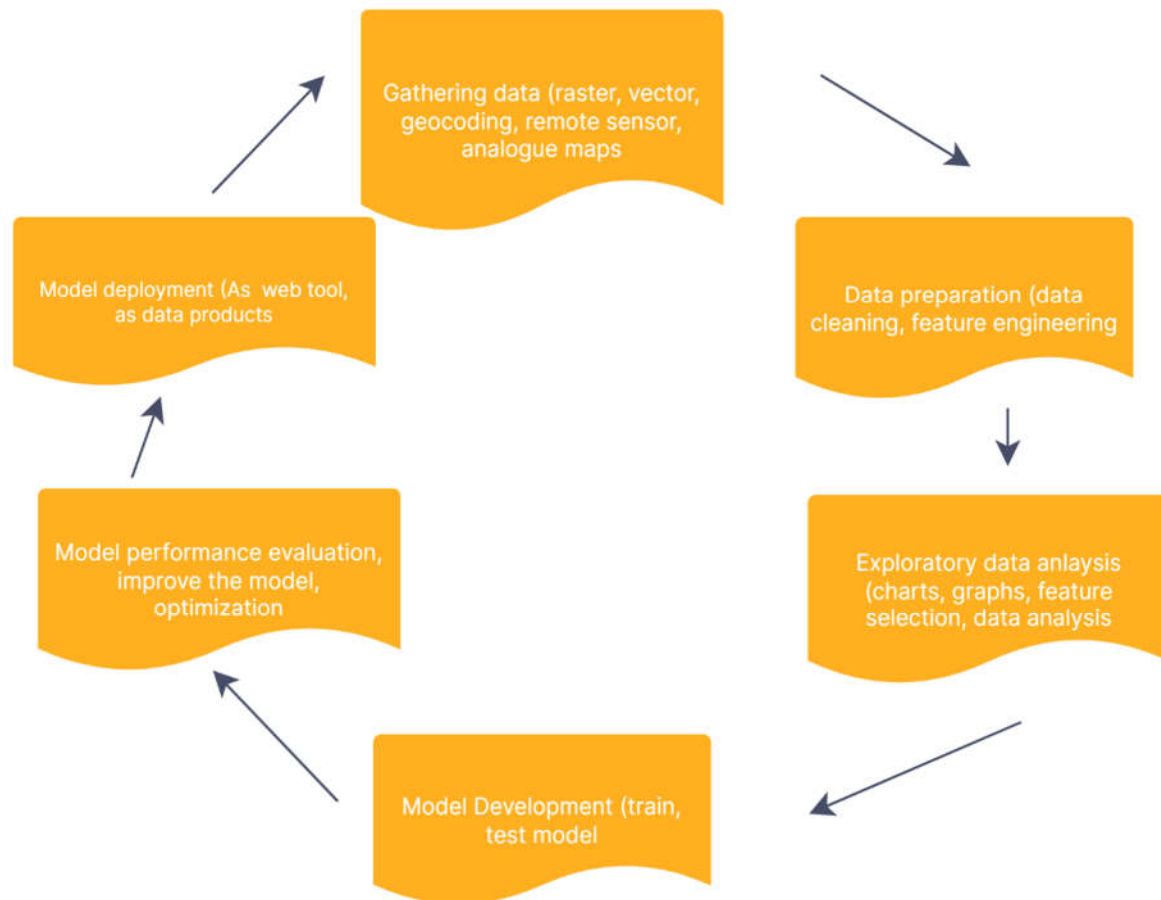


Figure 1. Machine learning workflow.

Applications of machine learning algorithms have increased drastically in G.I.S. and remote sensing in recent years (C. Xu & Jackson, 2019). It has a broad range of applications, from developing energy-based models to assessing soil liquefaction to creating a relation between air quality and mortality (Greener, Kandathil, Moffat, & Jones, 2022). Other examples include qualitative and quantitative evaluation of satellite imagery sensor data for regional and urban scale air quality (Avand & Moradi, 2021), support vector machine approach for longitudinal dispersion coefficients in natural streams (Bahari, Ahmad, & Aboobaidar, 2014), crisis management (Yu, 2017), disaster, linear programming for irrigation scheduling (Sun & Zhu, 2019), global climate change and weather forecast (Ise, Oba, & Al, 2019), the status of land cover classification accuracy assessment (J. Wang, Bretz, Dewan, & Delavar, 2022), air pollutants and sources associated with health effects (Verma & Verma, 2021), settlement detection (Assarkhaniki, Sabri, & Rajabifard, 2021) features such as roads/highways and ditch segments extraction (Avand & Moradi, 2021), identify crops' diseases and their yield estimation, building vegetation indices, natural disaster response, and disease outbreak response (Hossain, Zarin, Sahriar, Haque, & Chemistry of the Earth, 2022). In addition, researchers/users are benefitted from the publicly available remote sensing datasets using which they can develop, test and run their ML models for their research (Das, 2020). Most of the remote sensing datasets are global and unbiased (Palacios Salinas, Baratchi, Rijn, & Vollrath, 2021). This further simplifies the workflow in building accurate ML models in this domain (Odebiri, Odindi, Mutanga, & Geoinformation, 2021). Furthermore, remote sensing-based research is not halted due to natural disasters or unexpected accidents (Das, 2020).

Here, in this paper, we discuss the most popular supervised ML models (classification and regression) in G.I.S. and remote sensing. The motivation for writing this paper is that machine learning models produce higher accuracy than traditional parametric

classifiers, especially for complex data with many predictor variables(Das, Ghosh, Chowdary, Mitra, & Rijal, 2022). Therefore, the article can be a guide to those interested in using ML models in their G.I.S. and remote sensing based projects (Zerrouki, Harrou, Sun, & Hocini, 2019). This paper provides a general overview of 4 supervised non-parametric ML models that can be used in most of the G.I.S. and remote sensing based projects. We discuss classification (Naïve Bayes (NB), Support Vector Machine (SVM), Random Forest (RF), Decision Trees (DT)) and regression models (Random Forest (RF), Support Vector Machine (SVM), Linear, Count and Poisson) here. Binomial and multiclass classification models are more common in G.I.S. and remote sensing-based projects.(Avand & Moradi, 2021).If the classification has two classes, the classifier is known as binomial; if there are more than two classes, the category is multiclass.

											Attributes
...	Long_	Elev	ElevMin	ElevMax	Bname	Shape_Leng	Shape_Area	HydroID	OutletID	geometry	
...	-95.004900	3.857793	1.0	9.0	None	20278.2932	4.692447e+06	300001	100018	POLYGON ((304289.891 3264141.990, 304280.085 3...	
...	-95.028273	3.764508	0.0	13.0	None	28593.5702	1.114920e+07	300002	100001	POLYGON ((303858.438 3263161.415, 303838.826 3...	
...	-95.006282	3.241213	0.0	9.0	None	14571.3466	3.493518e+06	300003	100002	POLYGON ((304515.423 3262916.271, 304515.423 3...	
...	-95.034368	3.574497	0.0	6.0	None	16061.8212	2.994389e+06	300004	100003	POLYGON ((303250.481 3262631.904, 303250.481 3...	

Figure 2. Data engineering in G.I.S. world.

1.1. Supervised machine learning models in G.I.S. and remote sensing

1.1.1. Naïve Bayes Algorithms

These supervised models are the easiest to build, less complex and can be applied to massive datasets.(Liu et al., 2017) It is fast. However, Naïve Bayes classification cannot be used for continuous numerical values(Sitthi, Nagai, Dailey, & Ninsawat, 2016). It ignores noise, hence might lead to inaccurate predictions.(Tien Bui et al., 2018).There are three types of Naïve Bayes: Gaussian, Multinomial, Bernoulli . Gaussian assumes the distribution to be normal.(Chen, Hu, Hua, & Zhao, 2021).Multinomial for discrete counts and Bernoulli for binary outcomes.(Mitchell, 2005) .These classifiers are efficient for multiclass predictions.(El-Magd & Ahmed, 2022). These models can be best utilized in making best management practices models (B.M.P.s), habitat suitability models, weather prediction.

1.1.2. Random Forest Classifier

It's a supervised classification model that can be applied to classification and regression models(Belgiu, Drăguț, & sensing, 2016). It is a collection of decision trees and predicts the results based on the multiple models/sub-models.(Pal, 2005).Therefore it is also known as the ensemble classifier.(Piramanayagam, Schwartzkopf, Koehler, & Saber, 2016).R.F. works on the bagging principle while making models, which means it makes different models based on the subset of training sample data, and the outcome is based on the majority/average of the sub-models.(Berhane et al., 2018). Multiple studies suggest that the number of trees generally does not significantly impact the resulting R.F. classification accuracy, as long as the number is sufficiently large enough(Kulkarni & Lowe, 2016).This is because when the number of trees in the classifier is small, the prediction

accuracy increases as additional trees are added. Still, the accuracy tends to plateau with a large number of trees.(Chan, Huang, Defries, & Sensing, 2001; Chan & Paelinckx, 2008; Pal, 2005; Rodriguez-Galiano et al., 2012). Some common examples of the projects that can be solved using R.F. algorithm include: land use land cover classification(Thapa, Prasai, & Indrawati Municipality, 2022), feature extraction such as ditch segments, roads, settlements, or objects of interest, object detection such as tree species, vehicle, species identification such as tigers, elephants, bird species, insects, habitat classification(R. Prasai, 2021; R. J. N. Prasai, 2021) and modelling related projects such as flood-prone/drought(R. Prasai, 2022a, 2022b), core habitat, classify soil types, diseases, weeds, climate and weather-related model(Dahal & Prasai, 2022; R. Prasai, 2021, 2022a, 2022b, 2022c; R. Prasai et al., 2021; R. J. C. W. Prasai, Energy, & Engineering, 2022; R. J. N. Prasai, 2021; Thapa et al., 2022) and their forecasting.

1.1.3. Support vector machine

It's an ML model that can be applied to classification and regression problems(Mountrakis, Im, Ogole, & Sensing, 2011).It fits the data based on a distinct line known as a hyperplane(Sheykhmousa et al., 2020). As the model is easy to build and robust to outliers, it is widely used in the G.I.S. and remote sensing domains(Cavallaro, Willsch, Willsch, Michielsen, & Riedel, 2020). Building a support vector ML model requires the use to specify the kernel type(Waske, Benediktsson, & Sveinsson, 2009). Some popular kernels in remote sensing are polynomial kernels and the radial basis function (RBF) kernel(C. Huang, Davis, & Townshend, 2002). Classification of satellite based imagery, detection of features like roads, wetlands, grasslands, can be solved using SVM models.

1.1.4. Linear regression

These models are the most popular research models in G.I.S. and remote sensing (Sudalaimuthu & Sudalayandi, 2019).Linear regression helps to identify and evaluate the relationship between two or more factors/covariates when we leverage the power of space in our analysis using the distance features, for example, the influence of distance to water in habitat selection.(Mansouri, Feizi, Jafari Rad, & Arian, 2018). This ML model helps to address the questions like:

- Is there a linear relationship between diameter at breast height and crown diameter of trees?
- What demographic factors contribute to the use of high rates of public transport?
- What factors contribute to the high spread of COVID in geographical regions?
- What is the relation between environmental factors and the cyanobacteria population?
- What variables affect gender-specific leadership?
- What is the relation between climate change and migration?

There are 3 types of linear regression commonly used in GIS and remote sensing based projects. They are Continuous (Gaussian), Logistic and Poisson distribution. The distribution should be normal(Susiluoto, Spantini, Haario, Härkönen, & Marzouk, 2020) for the Gaussian distribution linear regression. It is also called continuous because the dependent variable can take a wide range of values such as temperature, rainfall, and tree diameter(Avand & Moradi, 2021).If the dependent variable is not normally distributed, we can change it to binary values using reclassify function(Shi, Li, & Zhao, 2020). Binary is also known as logistic regression models, which builds models with only two outputs:-pass/fail, presence/absence.(Ghosh et al., 2022). We use count/Poisson regression models if the dependent variables are the counts/number of occurrences of an event.(Graff et al., 2020).

The dependent variable cannot be negative or decimal values(Graff et al., 2020). These models are generally used for species distribution models and understanding event patterns.

1.1.5. Non-linear regression

The regression in which the predictor and response variable has a non-linear relationship is known as non-linear regression(Liang et al., 2022). Since most relationships in G.I.S. and remote sensing are non-linear, it is widely used in this sector(Adsuara et al., 2019). Due to its flexibility, a wide variety of models can be built using these models(Hsieh, 2020).For example, study the crops and soil processes, study the real estate price and immigration relation, study the relation between diameter and canopy cover.

1.2. Methods to improve the accuracy of the ML models

1.2.1. Feature engineering

Feature engineering is most prevalent in predictive models(Paulson et al., 2022). It is the process of filtering the most logical and influential variables/covariates in the models from the less important/influential variables, in ML terms, it is known as feature reduction(Sarith Divakar, Sudheep Elayidom, & Rajesh, 2022).It requires domain knowledge and understanding of the requirements of the projects(Paulson et al., 2022).Researchers run exploratory data analyses to observe the relationship between different variables/covariates and extract only the best variables to make an ML model(Song, Yang, Dai, Yuan, & Engineering, 2020).

1.2.2. Boosting

Boosting is a method used in machine learning to reduce errors in predictive data analysis(Schapire & classification, 2003).Data scientists train machine learning software, called machine learning models, on labelled data to make guesses about unlabeled data(Mayr, Binder, Gefeller, & Schmid, 2014).A single machine learning model might make prediction errors depending on the accuracy of the training dataset(Jafarzadeh, Mahdianpari, Gill, Mohammadimanesh, & Homayouni, 2021).For example, if a cat-identifying model has been trained only on images of white cats, it may occasionally misidentify a black cat. Boosting tries to overcome this issue by training multiple models sequentially to improve the accuracy of the overall system. Boosting improves machine models' predictive accuracy and performance by converting multiple weak learners into a single robust learning model. Machine learning models can be vulnerable learners or strong learners:

1.2.3. Hyperparameter optimization

Hyperparameter tuning depends on several factors: sample size, classifier/regression models used,and model type.(Audebert, Le Saux, Lefèvre, & magazine, 2019; S. Xu, Zhao, Wang, & Shi, 2022; Yang & Shami, 2020) It's an additional step to improve the accuracy and performance of the model(Pannakkong, Thiwa-Anont, Singthong, Parthanadee, & Buddhakulsomsiri, 2022). For example, selection of the best polynomial features in linear regression models, number of trees in a random forest, number of layers and neurons in a neural network, maximum depth in decision trees, and learning rate for gradient descent(Pannakkong et al., 2022). Some common hyper parameter tuning techniques are grid search, randomized search, Bayesian optimization, sequential model-based optimization, and genetic algorithms (Yang & Shami, 2020).

1.3. Overfitting and Underfitting in ML models

Overfitting occurs when the model learns the noise and unwanted details in the learning data, which negatively impacts predicting the new data(Gu et al., 2016).Underfitting refers to a model that neither models the training data nor generalizes to new data(Bashir, Montañez, Sehra, Segura, & Lauw, 2020; Guyon & Yao, 1999; Jabbar, Khan, & Devices, 2015; Van der Aalst et al., 2010).In comparing classifiers, we emphasize that more than just overall accuracy should be considered; the user's and producer's accuracies for individual classes should also be considered(Alnaim, Sun, & Tong, 2022). This

is particularly true if the mapping focuses on rare classes (i.e. classes of limited extent in the image data). Rare classes tend to have little effect on the overall accuracy but may nevertheless be vital in determining the usefulness of the classification (Bogner, Seo, Rohner, & Reineking, 2018). However, if it is not feasible to test a variety of classifiers, SVM and R.F. generally appear to be reliable classification methods (Bogner et al., 2018). Some common approaches to reduce overfitting and underfitting of ML models are to use cross entropy, cross validation, early stopping and regularization approach.

1.3.1. Cross-Entropy and Cross-Validation:

Entropy is a measurable physical quality most usually linked with disorder, unpredictability, or uncertainty (E.-W. Huang et al., 2022). The smallest average encoding size per transmission with which a source can efficiently convey a message to a destination without losing any data is defined as entropy (Janik, 2019). The difference between two probability distributions for a given random variable or set of occurrences is measured by cross-entropy. (Brochet, Lapuyade-Lahorgue, Bougleux, Salaün, & Ruan, 2021; Gordon-Rodriguez, Loaiza-Ganem, Pleiss, & Cunningham, 2020; Pacheco, Ali, & Trappenberg, 2019; Ruby & Yendapalli, 2020; Z. Z. Wang & Goh, 2022). As a loss function, cross-entropy is extensively employed in ML (Juszczuk et al., 2021). Each example has a known class label with a probability of 1.0, whereas all other labels have a probability of 0.0 in classification (Ho & Wookey, 2019). In this case, the model determines the probability that a given example corresponds to each class label (Singh, 2013). Cross-entropy can then be used to calculate the difference between two probability distributions. (Gordon-Rodriguez et al., 2020)

1.3.2. Cross-Validation

Cross-validation is a technique in which we train our model using the subset of the dataset and then evaluate it using the complementary subset of the dataset (Tougui, Jilbab, & El Mhamdi, 2021). It is useful when there is a limited amount of data available (Battula & Technology, 2021). An example of cross-validation is K-fold cross-validation (Learn, 2022). The data is divided into K parts, where 1 part is used as a validation dataset and the other remaining as a training dataset (Phinzi, Abriha, & Szabó, 2021). And this process is repeated K times to reduce the biases and produce an effective model (A. Ramezan, A. Warner, & E. Maxwell, 2019).

1.3.3. Early Stopping and Regularization

Early stopping and regularization are other techniques used to reduce the overfitting of the data (J. J. U. h. m. c. e.-s.-t.-a.-n.-n.-m. Brownlee, 2018). The early stopping technique stops the training on ML models once the ML model's performance starts dropping and then increasing (Behnke & Guo, 2021). The regularization technique can be applied in multiple ways. Their examples are L1, L2, and Dropout regularization (Alem & Kumar, 2022)

1.4. Model performance calculation

1.4.1. Confusion matrix (Popular for classification models)

A confusion matrix, also known as an error matrix, is used for classification models (Piscini, Carboni, Del Frate, & Grainger, 2014). These matrices help in evaluating, monitoring and managing models (Jeong, Ko, Shin, & Yeom, 2022). From these matrices, we can develop metrics like accuracy, precision, recall, specificity, and F1 score (Jeong et al., 2022). When we create a confusion matrix, positive observation is known as Positives (P) (Weaver, Moore, Reith, McKee, & Lunga, 2018), negative observation is known as Negative (N) ((Weaver et al., 2018)), an outcome where the model correctly predicts the positive class is called True Positives (T.P.). In this outcome, the model correctly predicts the negative classes are, known as True Negatives (T.N.). The model incorrectly predicts the positive class when negative, also called a type 1 error are False Positive (F.P.) (Bosman, Liotta, Iacca, & Wörtche, 2013) An outcome where the model incorrectly predicts the negative class

when it is positive also called a type 2 error, is known as a False Negative (F.N.) (Bosman et al., 2013). We should learn about the accuracy, precision, recall, specificity, F1 score to read and interpret the output of the confusion matrix.

- Accuracy

Accuracy can be calculated by using the following formula

$$\text{Accuracy} = \frac{\# \text{ of correct predictions}}{\text{total \# of predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$$

- Precision

Precision can be calculated by using the following formula

$$\text{Precision} = \frac{TP}{TP + FP}$$

- Recall

Recall, also known as the sensitivity, hit rate, or the true positive rate (T.P.R.) (Alakus, Turkoglu, & Fractals, 2020) answers the question, "What proportion of actual positives were identified correctly?"

$$\text{Recall} = \frac{TP}{TP + FN}$$

- Specificity

Specificity, also known as the true negative rate (TNR), measures the proportion of actual negatives that are correctly identified as such. (Erickson & Kitamura, 2021). It is the opposite of recall.

$$\text{Specificity} = \frac{TN}{TN + FP}$$

- F1 score

The F1 score measures a test's accuracy — it is the harmonic mean of precision and recall. (Chicco & Jurman, 2020)

It can have a maximum score of 1 (perfect precision and recall) and a minimum of 0. Overall, it measures the preciseness and robustness of your model. (Goutte & Gaussier, 2005)

$$\text{F1 score} = \frac{2 * (\text{precision} * \text{recall})}{\text{precision} + \text{recall}} = \frac{2TP}{2TP + FP + FN}$$

- Receiver operator characteristic curve

Receiver operator characteristic (R.O.C.) analysis is a quantitative method for determining a binary classification based on a threshold (cut-off) value usually calculated from continuous data. (Carter, Pan, Rai, & Galandiuk, 2016; Mbizvo & Lerner, 2021; Søreide, 2009). Plotting the true positive rate (T.P.R.) against the false positive rate (F.P.R.) at various threshold levels yields the R.O.C. curve (J. J. M. l. m. Brownlee, 2018). Sensitivity, recall, and the chance of detection are all terms used to describe the true-positive rate. (LeDell, Petersen, & van der Laan, 2015; Statnikov, Aliferis, Hardin, & Guyon, 2013). The likelihood of a false alarm is also known as the false-positive rate, and it can be computed as (1-specificity). It's also known as a plot of the power as a function of the decision rule's Type I Error (when the performance is calculated from just a sample of the population, it can be

thought of as estimators of these quantities). As a result, the R.O.C. curve represents sensitivity or recall as a function of fall-out.

1.4.1. For regression models:

Classifying the accuracy of the regression models is a little different than the classification models because, in regression models, we are not only concerned with the model predicting right or wrong but also how accurately the models have predicted the actual value (Goldstein, 2005). For example, when we use regression models to forecast the temperature, if the model gives the value as 43 C and the actual value is 43.5 C, the model is better and vice versa. (Haack & Rafter, 2010). We measure the accuracy of the regression models using explained variance and mean squared error

- Explained variance

Explained variance is the amount of variation in the original dataset that our model can explain (Estrella, Gilerson, Foster, & Groetsch, 2021; Goldstein, 2005).

- Mean squared error

It is the average of the squared differences between the predicted and actual output. R^2 coefficient represents the proportion of variance in the outcome that our model can predict based on its features (Mittlböck & Schemper, 1996).

1.5. Factors to consider while selecting the ML models in GIS and remote sensing based projects

- No rule of thumb
- Experiment with multiple classifiers
- Hyper tuning parameters for the accuracy
- Use random forest classifiers for the weak datasets and Decision trees when simple and fast models are needed
- The default value for the number of trees in R.F. can be 500; for kernel size in SVM, it can be polynomial kernels and radial basis kernels
- Visualize the relationships between the input and predictors to evaluate their relationship and find if there is any band that can help in predicting things better
- Normalize the rare classes/imbalanced datasets
- Computation time also depends on user-defined parameters, classifier chosen, sample size
- If parameters cannot be tuned, R.F. should be used, setting the number of trees to 500 to provide
- Balance the datasets/data normalization. The classes with few samples/rare classes can be affected
- Computational complexities of different ML models which is the amount of resources to run a ML model.

N =number of training examples, m =number of features, n' =number of support vectors, k =number of neighbors, k' = number of trees (Majeed, 2019)

Table 1. Computational complexity of discussed ML models.

S.N.	Model	Train time complexity	Test time complexity	Space complexity
1	Linear regression	$O(n*m^2 + m^3)$	$O(m)$	$O(m)$
2	Logistic regression	$O(n*m)$	$O(m)$	$O(m)$
3.	Support Vector Machine	$O(n^2)$	$O(n'*m)$	$O(n*m)$
4.	Decision tree	$O(n*\log(n)*m)$	$O(m)$	$O(\text{depth of tree})$
5.	Random forest	$O(k' * n * \log(n) * m)$	$O(m*k')$	$O(k' * \text{depth of tree})$
6.	Naïve Bayes	$O(n*m)$	$O(m)$	$O(c*m)$

2. Conclusion

In recent years, ML models are increasingly being used in GIS and remote sensing based projects. ML models helps in solving GIS and remote sensing problems by identifying the underlying patterns, for example satellite based image classification, detection of features likes roads, wetlands, grasslands, image segmentation. We discuss few popular ML models and methods of their application in GIS and remote sensing based projects here. Researchers can use this paper as a reference while starting a ML based project. There are other ML models which can be learnt easily after learning above discussed models.

Conflicts of Interest: The authors declare no conflict of interest.

References

- A. Ramezan, C., A. Warner, T., & E. Maxwell, A. J. R. S. (2019). Evaluation of sampling and cross-validation tuning strategies for regional-scale machine learning classification. 11(2), 185.
- Adsuara, J. E., Pérez-Suay, A., Muñoz-Marí, J., Mateo-Sanchis, A., Piles, M., Camps-Valls, G. J. I. T. o. G., & Sensing, R. (2019). Nonlinear distribution regression for remote sensing applications. 57(12), 10025-10035.
- Alakus, T. B., Turkoglu, I. J. C., Solitons, & Fractals. (2020). Comparison of deep learning approaches to predict COVID-19 infection. 140, 110120.
- Alem, A., & Kumar, S. J. A. A. I. (2022). Transfer Learning Models for Land Cover and Land Use Classification in Remote Sensing Image. 36(1), 2014192.
- Alnaim, A., Sun, Z., & Tong, D. J. R. S. (2022). Evaluating Machine Learning and Remote Sensing in Monitoring NO2 Emission of Power Plants. 14(3), 729.
- Assarkhaniki, Z., Sabri, S., & Rajabifard, A. J. B. E. D. (2021). Using open data to detect the structure and pattern of informal settlements: an outset to support inclusive SDGs' achievement. 5(4), 497-526.
- Audebert, N., Le Saux, B., Lefèvre, S. J. I. g., & magazine, r. s. (2019). Deep learning for classification of hyperspectral data: A comparative review. 7(2), 159-173.
- Avand, M., & Moradi, H. J. J. o. H. (2021). Using machine learning models, remote sensing, and GIS to investigate the effects of changing climates and land uses on flood probability. 595, 125663.
- Bahari, N. I. S., Ahmad, A., & Aboobaider, B. M. (2014). Application of support vector machine for classification of multispectral data. Paper presented at the IOP Conference Series: Earth and Environmental Science.
- Bashir, D., Montañez, G. D., Sehra, S., Segura, P. S., & Lauw, J. (2020). An information-theoretic perspective on overfitting and underfitting. Paper presented at the Australasian Joint Conference on Artificial Intelligence.
- Battula, K. J. I. J. o. E., & Technology, A. (2021). Research OF machine learning algorithms using K-fold cross validation. 8(6S), 215-218.
- Behnke, M., & Guo, S. J. P. M. (2021). Comparison of early stopping neural network and random forest for in-situ quality prediction in laser based additive manufacturing. 53, 656-663.
- Belgiu, M., Drăguț, L. J. I. j. o. p., & sensing, r. (2016). Random forest in remote sensing: A review of applications and future directions. 114, 24-31.
- Berhane, T. M., Lane, C. R., Wu, Q., Autrey, B. C., Anenkhonov, O. A., Chepinoga, V. V., & Liu, H. J. R. s. (2018). Decision-tree, rule-based, and random forest classification of high-resolution multispectral imagery for wetland mapping and inventory. 10(4), 580.
- Bogner, C., Seo, B., Rohner, D., & Reineking, B. J. P. o. (2018). Classification of rare land cover types: Distinguishing annual and perennial crops in an agricultural catchment in South Korea. 13(1), e0190476.
- Bosman, H. H., Liotta, A., Iacca, G., & Wörtche, H. J. (2013). Anomaly detection in sensor systems using lightweight machine learning. Paper presented at the 2013 IEEE International Conference on Systems, Man, and Cybernetics.

- Brochet, T., Lapuyade-Lahorgue, J., Bougleux, S., Salaün, M., & Ruan, S. J. I. (2021). Deep Learning Using Havrda-Charvat Entropy for Classification of Pulmonary Optical Endomicroscopy. 42(6), 400-406.
- Brownlee, J. J. M. I. m. (2018). How to use ROC curves and precision-recall curves for classification in Python. 30.
- Brownlee, J. J. U. h. m. c. e.-s.-t.-a.-n.-n.-m. (2018). A gentle introduction to early stopping to avoid overtraining neural networks.
- Carter, J. V., Pan, J., Rai, S. N., & Galandiuk, S. J. S. (2016). ROC-ing along: Evaluation and interpretation of receiver operating characteristic curves. 159(6), 1638-1645.
- Cavallaro, G., Willsch, D., Willsch, M., Michielsen, K., & Riedel, M. (2020). Approaching remote sensing image classification with ensembles of support vector machines on the d-wave quantum annealer. Paper presented at the IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium.
- Chan, J. C.-W., Huang, C., Defries, R. J. I. T. o. G., & Sensing, R. (2001). Enhanced algorithm performance for land cover classification from remotely sensed data using bagging and boosting. 39(3), 693-695.
- Chan, J. C.-W., & Paelinckx, D. J. R. S. o. E. (2008). Evaluation of Random Forest and Adaboost tree-based ensemble classification and spectral band selection for ecotope mapping using airborne hyperspectral imagery. 112(6), 2999-3011.
- Chen, H., Hu, S., Hua, R., & Zhao, X. J. E. J. o. A. i. S. P. (2021). Improved naive Bayes classification algorithm for traffic risk management. 2021(1), 1-12.
- Chicco, D., & Jurman, G. J. B. g. (2020). The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. 21(1), 1-13.
- Dahal, P., & Prasai, R. (2022). Representative Service Providers and Their Selection in Cloud Computing Domain; A Comprehensive Overview. doi: 10.20944/preprints202207.0190.v1
- Das, M. (2020). Online prediction of derived remote sensing image time series: An autonomous machine learning approach. Paper presented at the IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium.
- Das, M., Ghosh, S. K., Chowdary, V. M., Mitra, P., & Rijal, S. J. R. S. (2022). Statistical and Machine Learning Models for Remote Sensing Data Mining—Recent Advancements. In (Vol. 14, pp. 1906): MDPI.
- El-Magd, A., & Ahmed, S. J. A. J. o. G. (2022). Random forest and naïve Bayes approaches as tools for flash flood hazard susceptibility prediction, South Ras El-Zait, Gulf of Suez Coast, Egypt. 15(3), 1-12.
- Erickson, B. J., & Kitamura, F. J. R. A. I. (2021). Magician's corner: 9. Performance metrics for machine learning models. 3(3).
- Estrella, E. H., Gilerson, A., Foster, R., & Groetsch, P. J. J. o. A. R. S. (2021). Spectral decomposition of remote sensing reflectance variance due to the spatial variability from ocean color and high-resolution satellite sensors. 15(2), 024522.
- Ghosh, S. S., Dey, S., Bhogapurapu, N., Homayouni, S., Bhattacharya, A., & McNairn, H. J. R. S. (2022). Gaussian Process Regression Model for Crop Biophysical Parameter Retrieval from Multi-Polarized C-Band SAR Data. 14(4), 934.
- Goldstein, H. J. E. o. s. i. b. s. (2005). Heteroscedasticity and complex variation. 2, 790-795.
- Gordon-Rodriguez, E., Loaiza-Ganem, G., Pleiss, G., & Cunningham, J. P. (2020). Uses and abuses of the cross-entropy loss: Case studies in modern deep learning.
- Goutte, C., & Gaussier, E. (2005). A probabilistic interpretation of precision, recall and F-score, with implication for evaluation. Paper presented at the European conference on information retrieval.
- Graff, C. A., Coffield, S. R., Chen, Y., Foufoula-Georgiou, E., Randerson, J. T., Smyth, P. J. I. T. o. G., & Sensing, R. (2020). Forecasting daily wildfire activity using poisson regression. 58(7), 4837-4851.
- Greener, J. G., Kandathil, S. M., Moffat, L., & Jones, D. T. J. N. R. M. C. B. (2022). A guide to machine learning for biologists. 23(1), 40-55.
- Gu, Y., Wylie, B. K., Boyte, S. P., Picotte, J., Howard, D. M., Smith, K., & Nelson, K. J. J. R. s. (2016). An optimal sample data usage strategy to minimize overfitting and underfitting effects in regression tree models based on remotely-sensed data. 8(11), 943.
- Guyon, X., & Yao, J.-f. J. J. o. M. A. (1999). On the underfitting and overfitting sets of models chosen by order selection criteria. 70(2), 221-249.
- Haack, B., & Rafter, A. J. G. I. (2010). Regression estimation techniques with remote sensing: a review and case study. 25(1), 71-82.

- Ho, Y., & Wookey, S. J. I. A. (2019). The real-world-weight cross-entropy loss function: Modeling the costs of mislabeling. 8, 4806-4813.
- Hossain, M. T., Zarin, T., Sahriar, M. R., Haque, M. N. J. P., & Chemistry of the Earth, P. A. B. C. (2022). Machine learning based modeling for future prospects of land use land cover change in Gopalganj District, Bangladesh. 126, 103022.
- Hsieh, W. W. J. a. p. a. (2020). Improving predictions by nonlinear regression models from outlying input data.
- Huang, C., Davis, L., & Townshend, J. J. I. J. o. r. s. (2002). An assessment of support vector machines for land cover classification. 23(4), 725-749.
- Huang, E.-W., Lee, W.-J., Singh, S. S., Kumar, P., Lee, C.-Y., Lam, T.-N., . . . Reports, E. R. (2022). Machine-learning and high-throughput studies for high-entropy materials. 147, 100645.
- Ise, T., Oba, Y. J. F. i. R., & AI. (2019). Forecasting climatic trends using neural networks: an experimental study using global historical data. 32.
- Jabbar, H., Khan, R. Z. J. C. S., Communication, & Devices, I. (2015). Methods to avoid over-fitting and under-fitting in supervised machine learning (comparative study). 70.
- Jafarzadeh, H., Mahdianpari, M., Gill, E., Mohammadimanesh, F., & Homayouni, S. J. R. S. (2021). Bagging and boosting ensemble classifiers for classification of multispectral, hyperspectral and PolSAR data: a comparative evaluation. 13(21), 4405.
- Janik, R. A. J. a. p. a. (2019). Entropy from Machine Learning.
- Jeong, S., Ko, J., Shin, T., & Yeom, J.-m. J. S. R. (2022). Incorporation of machine learning and deep neural network approaches into a remote sensing-integrated crop model for the simulation of rice growth. 12(1), 1-10.
- Juszczuk, P., Kozak, J., Dzikowski, G., Głowania, S., Jach, T., & Probiez, B. J. E. (2021). Real-World Data Difficulty Estimation with the Use of Entropy. 23(12), 1621.
- Kulkarni, A. D., & Lowe, B. (2016). Random forest algorithm for land cover classification.
- Learn, S. J. I. A. h. s.-l. o. s. m. c. v. h. c. (2022). Cross-validation: evaluating estimator performance.
- LeDell, E., Petersen, M., & van der Laan, M. J. E. j. o. s. (2015). Computationally efficient confidence intervals for cross-validated area under the ROC curve estimates. 9(1), 1583.
- Liang, D., Frederick, D. A., Lledo, E. E., Rosenfield, N., Berardi, V., Linstead, E., & Maoz, U. J. B. I. (2022). Examining the utility of nonlinear machine learning approaches versus linear regression for predicting body image outcomes: The US Body Project I. 41, 32-45.
- Liu, R., Chen, Y., Wu, J., Gao, L., Barrett, D., Xu, T., . . . Yu, J. J. R. a. (2017). Integrating entropy-based naïve Bayes and GIS for spatial evaluation of flood hazard. 37(4), 756-773.
- Majeed, A. J. A. o. D. S. (2019). Improving time complexity and accuracy of the machine learning algorithms through selection of highly weighted top k features from complex datasets. 6(4), 599-621.
- Mansouri, E., Feizi, F., Jafari Rad, A., & Arian, M. J. S. E. (2018). Remote-sensing data processing with the multivariate regression analysis method for iron mineral resource potential mapping: a case study in the Sarvian area, central Iran. 9(2), 373-384.
- Maxwell, A. E., Warner, T. A., & Fang, F. J. I. J. o. R. S. (2018). Implementation of machine-learning classification in remote sensing: An applied review. 39(9), 2784-2817.
- Mayr, A., Binder, H., Gefeller, O., & Schmid, M. J. M. o. i. i. m. (2014). The evolution of boosting algorithms. 53(06), 419-427.
- Mbizvo, G. K., & Lamer, A. J. J. N. D. M. (2021). Receiver operating characteristic plot and area under the curve with binary classifiers: pragmatic analysis of cognitive screening instruments. 11(05), 353-360.
- Mitchell, T. M. J. M. L. (2005). Logistic regression. 10, 701.
- Mittlböck, M., & Schemper, M. J. S. i. m. (1996). Explained variation for logistic regression. 15(19), 1987-1997.
- Mountrakis, G., Im, J., Ogole, C. J. I. J. o. P., & Sensing, R. (2011). Support vector machines in remote sensing: A review. 66(3), 247-259.
- Odebiri, O., Odindi, J., Mutanga, O. J. I. J. o. A. E. O., & Geoinformation. (2021). Basic and deep learning models in remote sensing of soil organic carbon estimation: A brief review. 102, 102389.

- Pacheco, A. G., Ali, A.-R., & Trappenberg, T. J. a. p. a. (2019). Skin cancer detection based on deep learning and entropy to detect outlier samples.
- Pal, M. J. I. j. o. r. s. (2005). Random forest classifier for remote sensing classification. 26(1), 217-222.
- Palacios Salinas, N. R., Baratchi, M., Rijn, J. N. v., & Vollrath, A. (2021). Automated machine learning for satellite data: integrating remote sensing pre-trained models into AutoML systems. Paper presented at the Joint European Conference on Machine Learning and Knowledge Discovery in Databases.
- Pannakkong, W., Thiwa-Anont, K., Singthong, K., Parthanadee, P., & Buddhakulsomsiri, J. J. M. P. i. E. (2022). Hyperparameter Tuning of Machine Learning Algorithms Using Response Surface Methodology: A Case Study of ANN, SVM, and DBN. 2022.
- Paulson, N. H., Kubal, J., Ward, L., Saxena, S., Lu, W., & Babinec, S. J. J. J. o. P. S. (2022). Feature engineering for machine learning enabled early prediction of battery lifetime. 527, 231127.
- Phinzi, K., Abriha, D., & Szabó, S. J. R. S. (2021). Classification efficacy using k-fold cross-validation and bootstrapping resampling techniques on the example of mapping complex gully systems. 13(15), 2980.
- Piramanayagam, S., Schwartzkopf, W., Koehler, F., & Saber, E. (2016). Classification of remote sensed images using random forests and deep learning framework. Paper presented at the Image and signal processing for remote sensing XXII.
- Piscini, A., Carboni, E., Del Frate, F., & Grainger, R. G. J. A. i. R. S. (2014). A Neural Network Algorithm to Detect Sulphur Dioxide Using IASI Measurements.
- Prasai, R. (2021). Distribution of Bengal Tiger (*Panthera tigris*) and Their Main Prey Species in Chitwan National Park, Nepal. Tarleton State University,
- Prasai, R. (2022a). Earth engine application to retrieve long-term terrestrial and aquatic time series of satellite reflectance data. doi: 10.54660/anfo.2022.3.3.11
- Prasai, R. (2022b). An open-source web-based tool to perform spatial multicriteria analysis. doi: 10.54660/anfo.2022.3.3.19
- PRASAI, R. (2022c). Using Google Earth Engine for the complete pipeline of temporal analysis of NDVI in Chitwan National Park of Nepal. doi: 10.21203/rs.3.rs-1633994/v3
- Prasai, R., Schwertner, T. W., Mainali, K., Mathewson, H., Kafley, H., Thapa, S., . . . Drake, J. J. E. I. (2021). Application of Google earth engine python API and NAIP imagery for land use and land cover classification: A case study in Florida, USA. 66, 101474. doi: 10.1016/j.ecoinf.2021.101474
- Prasai, R. J. C. W., Energy,, & Engineering, E. (2022). Pre and Post Effects Assessment of Marine Ranch Construction in Chlorophyll-a Concentration Using MODIS Data and a Web-Based Tool. A Case Study in Zhelin Bay, China. 11(3), 85-92. doi: 10.4236/cweee.2022.113005
- Prasai, R. J. N. (2021). Distribution of Bengal Tiger (*Panthera tigris*) and Their Main Prey Species in Chitwan National Park.
- Rodriguez-Galiano, V. F., Ghimire, B., Rogan, J., Chica-Olmo, M., Rigol-Sanchez, J. P. J. I. j. o. p., & sensing, r. (2012). An assessment of the effectiveness of a random forest classifier for land-cover classification. 67, 93-104.
- Ruby, U., & Yendapalli, V. J. I. J. A. T. C. S. E. (2020). Binary cross entropy with deep learning technique for image classification. 9(10).
- Sarith Divakar, M., Sudheep Elayidom, M., & Rajesh, R. (2022). Feature Engineering of Remote Sensing Satellite Imagery Using Principal Component Analysis for Efficient Crop Yield Prediction. In *Evolutionary Computing and Mobile Sustainable Networks* (pp. 189-199): Springer.
- Sarker, I. H. J. S. C. S. (2021). Machine learning: Algorithms, real-world applications and research directions. 2(3), 1-21.
- Schapire, R. E. J. N. e., & classification. (2003). The boosting approach to machine learning: An overview. 149-171.
- Sheykhmousa, M., Mahdianpari, M., Ghanbari, H., Mohammadimanesh, F., Ghamisi, P., Homayouni, S. J. I. J. o. S. T. i. A. E. O., & Sensing, R. (2020). Support vector machine versus random forest for remote sensing image classification: A meta-analysis and systematic review. 13, 6308-6325.
- Shi, X., Li, Y., & Zhao, Q. J. R. S. (2020). Flexible hierarchical Gaussian mixture model for high-resolution remote sensing image segmentation. 12(7), 1219.

- Sidey-Gibbons, J. A., & Sidey-Gibbons, C. J. J. B. m. r. m. (2019). Machine learning in medicine: a practical introduction. 19(1), 1-18.
- Singh, V. P. (2013). Entropy theory and its application in environmental and water engineering: John Wiley & Sons.
- Sitthi, A., Nagai, M., Dailey, M., & Ninsawat, S. J. S. (2016). Exploring land use and land cover of geotagged social-sensing images using naive bayes classifier. 8(9), 921.
- Song, H., Yang, W., Dai, S., Yuan, H. J. M. B., & Engineering. (2020). Multi-source remote sensing image classification based on two-channel densely connected convolutional networks. 17(6), 7353-7378.
- Søreide, K. J. J. o. c. p. (2009). Receiver-operating characteristic curve analysis in diagnostic, prognostic and predictive biomarker research. 62(1), 1-5.
- Statnikov, A., Aliferis, C. F., Hardin, D. P., & Guyon, I. (2013). Gentle Introduction To Support Vector Machines In Biomedicine, A-Volume 2: Case Studies And Benchmarks: World Scientific Publishing Company.
- Sudalaimuthu, K., & Sudalayandi, K. (2019). Development of linear regression model to predict ground elevation from satellite elevation–statistical approach. Paper presented at the AIP Conference Proceedings.
- Sun, L., & Zhu, Z. (2019). Linear programming Monte Carlo method based on remote sensing for ecological restoration of degraded ecosystem. Paper presented at the IOP Conference Series: Earth and Environmental Science.
- Susiluoto, J., Spantini, A., Haario, H., Härkönen, T., & Marzouk, Y. J. G. M. D. (2020). Efficient multi-scale Gaussian process regression for massive remote sensing data with satGP v0. 1.2. 13(7), 3439-3463.
- Thapa, N., Prasai, R. J. M., & Indrawati Municipality, N. (2022). Impacts of Floods in Land Use Land Cover Change:-a Case Study of Indrawati and Melamchi River, Melamchi, and Indrawati Municipality, Nepal. doi: 10.2139/ssrn.4104357
- Tien Bui, D., Shahabi, H., Shirzadi, A., Chapi, K., Hoang, N.-D., Pham, B. T., . . . Bin Ahmad, B. J. R. S. (2018). A novel integrated approach of relevance vector machine optimized by imperialist competitive algorithm for spatial modeling of shallow landslides. 10(10), 1538.
- Tougui, I., Jilbab, A., & El Mhamdi, J. J. H. i. r. (2021). Impact of the choice of cross-validation techniques on the results of machine learning-based diagnostic applications. 27(3), 189-199.
- Van der Aalst, W. M., Rubin, V., Verbeek, H., van Dongen, B. F., Kindler, E., Günther, C. W. J. S., & Modeling, S. (2010). Process mining: a two-step approach to balance between underfitting and overfitting. 9(1), 87-111.
- Verma, V. K., & Verma, S. J. M. T. P. (2021). Machine learning applications in healthcare sector: An overview.
- Wang, J., Bretz, M., Dewan, M. A. A., & Delavar, M. A. J. S. o. T. T. E. (2022). Machine learning in modelling land-use and land cover-change (LULCC): Current status, challenges and prospects. 153559.
- Wang, Z. Z., & Goh, S. H. J. A. G. (2022). A maximum entropy method using fractional moments and deep learning for geotechnical reliability analysis. 17(4), 1147-1166.
- Waske, B., Benediktsson, J. A., & Sveinsson, J. R. (2009). Classifying remote sensing data with support vector machines and imbalanced training data. Paper presented at the International Workshop on Multiple Classifier Systems.
- Weaver, J., Moore, B., Reith, A., McKee, J., & Lunga, D. (2018). A comparison of machine learning techniques to extract human settlements from high resolution imagery. Paper presented at the IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium.
- Xu, C., & Jackson, S. A. J. G. b. (2019). Machine learning and complex biological data. In (Vol. 20, pp. 1-4): Springer.
- Xu, S., Zhao, Y., Wang, M., & Shi, X. J. E. J. o. S. S. (2022). A comparison of machine learning algorithms for mapping soil iron parameters indicative of pedogenic processes by hyperspectral imaging of intact soil profiles. 73(1), e13204.
- Yang, L., & Shami, A. J. N. (2020). On hyperparameter optimization of machine learning algorithms: Theory and practice. 415, 295-316.
- Yu, X. J. N. H. (2017). Disaster prediction model based on support vector machine for regression and improved differential evolution. 85(2), 959-976.

Zerrouki, N., Harrou, F., Sun, Y., & Hocini, L. J. I. S. J. (2019). A machine learning-based approach for land cover change detection using remote sensing and radiometric measurements. 19(14), 5843-5850.