# Machine learning applied to the prediction of rockfall slope probability

Aprendizado de máquina aplicado à predição da probabilidade de quedas de blocos em taludes

Aprendizaje de máquinas aplicado a la predicción de la probabilidad de caída de rocas

**Larissa Regina Costa Silveira**
ORCID: https://orcid.org/0000-0001-7403-4463
Federal University of Ouro Preto, Brazil
E-mail: larissaeng07@gmail.com
**Milene Sabino Lana**
ORCID: https://orcid.org/0000-0001-9077-279X
Federal University of Ouro Preto, Brazil
E-mail: milene@ufop.edu.br
**Tatiana Barreto dos Santos**
ORCID: https://orcid.org/0000-0001-5484-6675
Federal University of Ouro Preto, Brazil
E-mail: tatiana.santos@ufop.edu.br

**Abstract**
The objective of this work is to propose a predictive model of rockfall slope probability in rock slopes using the K-Nearest Neighbors (KNN) method. A dataset composed by 220 rock slopes was used, whose variables are related to the presence of water, characteristics of the rock mass, degree of overhang, among others. For each slope of the dataset, rockfall probability (high, medium, or low) is known and determined by cluster analysis. The number of the nearest neighbors (k) ranged from 1 to 20. The obtained average accuracy of the tested predictive models was equal to 78.4%. The models produced satisfactory results in the prediction of the rockfall probability, since the area under the ROC curve was equal to 0.80. The best model was selected based on the k value with the highest accuracy and the highest area under the ROC curve. The selected model had a k value equal to 7.
**Keywords:** Rockfall; Machine learning; K-Nearest neighbors; Rock slope stability.

**Resumo**
O objetivo desse trabalho é propor um modelo de predição da probabilidade de queda de blocos em taludes rochosos utilizando o método K-Nearest Neighbors (KNN). Foi utilizado um banco de dados composto por 220 taludes rochosos, cujas variáveis estão relacionadas à presença de água, características do maciço rochoso, descalçamento de blocos, entre outras. Para cada talude do banco de dados, a probabilidade de queda de blocos (alta, média ou baixa) é conhecida e foi determinada através de análise de agrupamento. O número de vizinhos mais próximos (k) variou entre 1 e 20. A acurácia média obtida dos modelos de predição testados foi igual a 78,4%. Os modelos produziram resultados satisfatórios na previsão da probabilidade de queda de blocos, uma vez que a área sob a curva ROC foi igual a 0,80. O melhor modelo foi selecionado com base no valor de k com maior acurácia e maior área sob a curva ROC. O modelo selecionado teve um valor de k igual a 7.
**Palavras-chave:** Queda de blocos; Aprendizado de máquina; *K-Nearest Neighbors*; Estabilidade de taludes em rocha.

**Resumen**
El objetivo de este trabajo es proponer un modelo de predicción de la probabilidad de caída de bloques en taludes rocosos utilizando el método K-Nearest Neighbors (KNN). Se utilizó una base de datos con 220 taludes rocosos, cuyas variables están relacionadas con la presencia de agua, características del macizo rocoso, la presencia de bloques sueltos en los taludes, entre otras. Para cada talud del conjunto de datos, se conoce la probabilidad de caída de rocas (alta, media o baja) y se determinó a través del análisis de conglomerados. El número de vecinos más cercanos (k) se varió entre 1 y 20. La precisión promedio obtenida de los modelos de predicción probados fue igual a 78,4%. Los modelos arrojaron resultados satisfactorios en la predicción de la probabilidad de caída de rocas, ya que el área bajo la curva ROC fue igual a 0,80. El mejor modelo se seleccionó en función del valor k con mayor precisión y el área más alta bajo la curva ROC. El modelo seleccionado fue el que tenía un valor de k igual a 7.
**Palabras clave:** Caída de rocas; Aprendizaje automático; *K-Nearest Neighbors;* Estabilidad de taludes de roca.

## 1. Introduction

Machine learning and multivariate statistics are widely used by researches of several areas of knowledge, especially

those with the aim of finding dataset patterns. The identification of these patterns provides the possibility of predicting the behavior of new individuals in the model.

Mascarenhas et al. (2020) applied machine learning algorithms to propose an automatic classification system of Specialized Knowledge of Physics Teachers based on a pre-classified database. Ossani et al. (2020) carried out unsupervised classification learning techniques to find clustering patterns of specialty coffees and compared the obtained clusters with the original ones. Subsequently, Ossani et al. (2021) used supervised machine learning techniques to classify specialty coffees and they compared the performance of each used technique. Silva et al. (2021) used artificial neural networks and linear regression to build a tool for predicting the spatio-temporal distribution of viruses transmitted by Aedes aegypti. Fernandes et al. (2021) compared different artificial neural networks architectures to evaluate their behavior into predicting charges in an electrical system. Pessoa et al. (2021) used artificial neural networks to predict the load capacity of foundation.

In geotechnical engineering, there are worldwide methodologies for rock mass classification and excavation stability analysis, such as Rock Mass Rating (RMR) (Bieniawski, 1989), Q-System (Barton et al., 1974), Slope Mass Rating (SMR) (Romana, 1985), Q-Slope (Bar & Barton, 2017). However, assessment models and systems of rock mass classifications have often a high degree of uncertainty and subjectivity, since they are based only on the field survey experience and general empirical rules.

Taking into account the successful application of machine learning and multivariate statistical techniques to create prediction models, Santos et al. (2021) applied these techniques to predict the class of a rock mass according to a modified Rock Mass Rating (RMR). In the model, only the relevant variables were considered, since they were determined through multivariate factor analysis, which reduces the subjectivity inherent to rock mass classification problems.

Subsequently, Santos et al. (2022) compared machine learning techniques to make predictions of classes in rock mass using the same dataset of Santos et al. (2021). Regarding the use of multivariate statistical techniques and machine learning for rock slope stability analysis, Santos et al. (2019) and Naghadehi et al. (2013) proposed models to predict the stability condition classification of rock mine slopes.

Rockfall is a complex slope mass movement, difficult to predict. A trigger is not always necessary for a rockfall movement, differently from soil failures, which have in the precipitation an example of a common trigger. In this context, monitoring of geological risk areas is common in periods of high precipitation. It does not always occur in relation to rockfalls, which can result in catastrophic events in urban areas, highways and mining.

Some methodologies were developed in order to assess rockfall hazard, such as Rockfall Hazard Rating System (RHRS) (Pierson & Van Vickle, 1993) and Colorado Rockfall Hazard Rating System (CRHRS) (Santi et al., 2009). These methodologies were proposed for highway slopes; but they are not able to predict rockfall probability. These methods rank the evaluated slopes in more hazardous and less hazardous, according to the sum of the scores attributed to the variables related to rock mass and traffic conditions. They have the goal of determining the slopes where the intervention is more urgent.

Therefore, methodologies capable of predicting rockfall probability are needed to solve the aforementioned problems. However, because of the uncertainties inherent in field surveys and rockfall movements, these methodologies must be optimized, as accurate as possible, and must be able to quantify the errors of prediction. The objective of this research is to propose a classification model of rockfall probability through K-Nearest Neighbors (KNN). The number of nearest neighbors (k) was varied and the best classification model to predict the class of any rock slope was proposed. A dataset with 220 slopes was used, whose rockfall probability classification (high, medium and low) was determined through cluster analysis, which is an unsupervised multivariate statistical technique.

**1.1 Cluster Analysis**

Cluster Analysis is an unsupervised multivariate statistical technique used to group individuals in homogeneous clusters without any prior labeling of individuals. Among the various clustering techniques presented in the literature, an example is the non-hierarchical method Kmedoids (Kaufman & Rousseeuw, 1990). Partitioning Around Medoids algorithm (PAM) can be used to perform cluster analysis through kmedoids method.

According to Kassambara (2017), the steps of the PAM algorithm are:

1st - the algorithm randomly selects k individuals to become the medoids. The medoid is the representative individual of each group, so the number of groups is equal to k;

2nd - the dissimilarity matrix is calculated and every individual in the dataset is assigned to a cluster, according to the distance between it and the medoid;

3rd - if any individual in any cluster is able to reduce the dissimilarity coefficient, this individual becomes the new medoid of this cluster and the algorithm repeats the steps mentioned above. If not, the algorithm ends.

Dissimilarity matrix can be computed using any statistical distance, like Euclidean distance and the Manhattan distance. Manhattan distance may be applied when the database contains outliers.

**1.2 K-Nearest Neighbors**

K-Nearest Neighbors (KNN) is a supervised machine learning technique used to predict the class of an individual according to the similarities between this individual and the individuals pre-classified in specific classes. A way to evaluate the similarity between individuals is through the statistical distance measures. An individual will be classified in a class where the distances between it and a k number of the nearest neighbors (labeled individuals) are the smallest. This distance measure can be the Euclidean distance, Minkowski distance or the Mahalanobis distance (Kubat, 2017).

The steps of the algorithm to establish the class of a new individual are:

1st - the distances between the individuals are calculated;

2nd - the labeled individuals closer to the new individual are found;

3rd - the new individual is classified according to the class of the most k nearest neighbors (k labeled individuals with the shortest distances).

The number of the near neighbors whose distances will be evaluated, the number k, must be provided. If k is equal to 1, the individual will automatically be classified in the class where its nearest neighbor is allocated. If k is equal to 3, the distance between the new individual and its three closest neighbors will be evaluated and the new individual is classified according to the class of the k nearest neighbors. In Figure 1, when k is equal to 1, the new individual is classified in Class B; when k is equal to 3, the individual is classified in Class A.

**Figure 1 -** Classification of an individual through KNN.



Source: Authors.

To obtain the best value for k, tests of models varying the value of k must be done. The optimal k is related to the model with the best validation metrics, such as the apparent error (Equation 1) and the accuracy of the model (Equation 2). Another important and widely used metric to evaluate the performance of the model is the area under the ROC curve (AUC). An AUC equal to 1 represents a perfect model, without errors. Therefore, the closer the AUC is to 1, the better the model.

$$Aparent\ error = \frac{samples\ incorrectly\ classified}{total\ number\ of\ samples} \tag{1}$$

$$Accuracy = \frac{samples\ correctly\ classified}{total\ number\ of\ samples} \tag{2}$$

## 2. Methodology

### 2.1 Softwares and repositories

Both KNN and cluster analysis were performed using freeware R version 4.0.2 (R Core Team, 2020). Methodology scripts are available on the GitHub platform:

https://github.com/larissarcs/Prediction-of-rockfall-probability-using-KNN/blob/main/KNN

https://github.com/larissarcs/PAM_cluster_Likelihood/blob/main/PAM_cluster_Likelihood.

### 2.2 Dataset

The dataset used in this research is part of the dataset used by Santi et al. (2009) to generate CRHRS. It is composed of 220 rock slopes and the variables were surveyed in highway slopes of Colorado (USA). Although the database refers to highway slopes in the state of Colorado, all the variables and characteristics considered in this study could be easily surveyed on rock slopes located in any place of the world. These parameters are traditionally used in rock mass classifications and slope stability analysis, and some of the variables included in this research are evaluated in a similar way in classification methodologies, already established in geotechnical engineering practice, such as the RMR (Bieniawski, 1989).

All variables related to the rock mass in CRHRS method were considered in this study; except the number of discontinuity sets and the weathering degree of the intact rock. The number of sets does not vary in the dataset used to propose the model. Weathering degree of the intact rock was not considered, as rockfalls occur even if the intact rock is fresh. Weathering degree and infilling of the joints were considered. Table 1 shows the eight independent variables (P1 to P8) used in the proposed model, where each variable received scores ranging from 1 to 4, according to its characteristics (1 represents a safe condition and 4 a critical condition).

**Table 1 -** Independent variables.

| VARIABLES | Characteristics/scores | | | |
|---|---|---|---|---|
| | 1 | 2 | 3 | 4 |
| Water condition (P1) | Dry | Dump/wet | Dripping | Running water |
| Rock character (P2) | Homogeneous/ massive | Small faults/strong veins | Schist/ shear zones < 15 cm | Weak pegmatite/micas/ shear zones >15 cm |
| Degree of overhang (P3) | 0 to 0.3 m | 0.3 to 0.6 m | 0.6 to 1.2 m | > 1.2 m |
| Block size/ volume (P4) | <0.3m/ <0.75m³ | 0.3 to 0.6m/ 0.75 to 2.3m³ | 0.6 to 1.5m/ 2.3 to 7.6m³ | > 1.5m/ > 7.6m³ |
| Discontinuities: Persistence and orientation (P5) | < 3m and dips into slope | > 3m dips to slope | < 3m and daylights out of the slope | > 3m and daylights out of the slope |
| Discontinuities: Aperture (P6) | 0 | 0.1 to 1 mm | 1 to 5 mm | > 5mm |
| Discontinuities: Weathering condition (P7) | Fresh | Surface staining | Granular infilling | Clay infilling |
| Discontinuities: Friction (P8) | Rough | Undulating | Planar | Slickenside |

Source: Adapted from Santi *et al*. (2009).

As KNN is a supervised machine learning technique, the status or dependent variable must be known. The dependent variable was determined by cluster analysis. The obtained clusters were labeled and each group was classified as high, medium and low rockfall probability. Non-hierarchical method *kmedoid* was used. Cluster analysis was carried out through PAM algorithm, using *factoextra* package (Kassambara & Mundt, 2020) from R software (R Core Team, 2020).

Cluster analysis grouped the dataset into three clusters based on the sum of scores assigned to variables. As there are eight variables and the score ranges from 1 to 4, the minimum possible final score is 8, representing a safe slope, with low rockfall probability. The maximum possible score is 32, representing an unsafe slope, with high rockfall probability. Table 2 presents the range of the scores for each class, used to classify the 220 slopes of the dataset. Table 3 presents a part of the dataset with scores ranging from 1 to 4, assigned to the independent variables and the dependent variable, obtained through *kmedoid*.

**Table 2 -** Classification criteria from cluster analysis.

| Probability class | Sum of the scores |
|---|---|
| Low | 8 - 16 |
| Medium/Transition Zone | 17 - 20 |
| High | 21 - 32 |

Source: Authors.

**Table 3 -** First six slopes of the dataset.

| P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | Status |
|----|----|----|----|----|----|----|----|--------|
| 2  | 3  | 3  | 3  | 2  | 4  | 3  | 1  | high   |
| 3  | 4  | 3  | 3  | 4  | 4  | 3  | 1  | high   |
| 3  | 3  | 2  | 2  | 4  | 3  | 3  | 2  | high   |
| 2  | 3  | 4  | 3  | 2  | 4  | 3  | 1  | high   |
| 1  | 3  | 3  | 2  | 2  | 4  | 3  | 1  | medium |
| 1  | 3  | 3  | 2  | 2  | 4  | 3  | 1  | medium |

Source: Authors.

## 2.3 Applied methodology

The developed methodology is summarized in the flowchart shown in Figure 2. The explanation of each step summarized in the flowchart appears next.

**Figure 2 -** Methodology Flowchart for rockfall probability model.



Source: Authors.

Before applying machine learning techniques, standardization of variables must be applied to solve scale problems, especially when the variables have different measurement units or large differences in their magnitude (Kubat, 2017). As the dataset used in this work is ordinal and all variables can only receive integer values between 1 and 4, this problem does not occur. However, according to Laurence (1992), for the application of artificial neural networks (ANN) and other machine learning techniques in ordinal data, it is convenient to convert the data to a percentile to keep the value below 1. Therefore, the dataset was normalized on a scale between 0 and 1. Table 4 shows part of the dataset already pre-processed, ready for application of the technique.

**Table 4 -** Part of the dataset already prepared for the application of KNN (P1 to P8 represent the independent variables and Status represents the dependent variable or class).

| P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | Status |
|------|------|------|------|------|------|------|------|--------|
| 0.33 | 0.67 | 0.67 | 0.67 | 0.33 | 1.00 | 0.67 | 0.00 | high |
| 0.67 | 1.00 | 0.67 | 0.67 | 1.00 | 1.00 | 0.67 | 0.00 | high |
| 0.67 | 0.67 | 0.33 | 0.33 | 1.00 | 0.67 | 0.67 | 0.33 | high |
| 0.33 | 0.67 | 1.00 | 0.67 | 0.33 | 1.00 | 0.67 | 0.00 | high |
| 0.00 | 0.67 | 0.67 | 0.33 | 0.33 | 1.00 | 0.67 | 0.00 | medium |
| 0.00 | 0.67 | 0.67 | 0.33 | 0.33 | 1.00 | 0.67 | 0.00 | medium |

Source: Authors.

In order to validate the proposed rockfall probability model, a randomly subsampling of the dataset was carried out. The 220 samples were randomly divided into 70% for training and 30% for test to validate the model.

To apply KNN algorithm, *class* package from R software was used (Venables & Ripley, 2002; R Core Team, 2020). To use KNN algorithm, the k number of neighbors must be predetermined. K was varied between 1 and 20, and the model with best metrics (apparent error, accuracy and AUC) was chosen. The performance was evaluated in training and test samples, in order to find models with overfitting and choose the most suitable model to predict rockfall probability.

After the choice of the most suitable rockfall probability model, it was applied to two new slopes in order to determine their rockfall probability. These two slopes are located in a quartzite mine in São Thomé das Letras city, Minas Gerais State (Brazil).

## 3. Results and Discussion

### 3.1. Determination of the best (suitable) model

Figure 3 presents a graph with the error rate of the test sample, considering the variation of k between 1 and 20.

**Figure 3 -** Apparent error for each value of k, considering the test sample.



Source: Authors.

Each point in Figure 3 represents the value of the obtained apparent error; upper and lower dashed horizontal lines represent, respectively, the largest and the smallest error; vertical lines indicate the k value referring to the smallest errors.

The lowest error rate (16.7%) refers to k values equal to 1, 3 and 7, see Figure 3. 11 slopes of the total of 66 (test sample is 30% of the dataset) were incorrectly classified. The highest error rate was equal to 30.3%, when the k value was equal to 19. It is important to emphasize that a value of k equal to 1 is not statistically significant, because the algorithm will force the classification of the individual in the class where its nearest neighbor is allocated. The algorithm uses only one sample to classify the individual. According to Kubat (2017), KNN classifier's performance should improve for k > 1, because the effect of the noisy nearest neighbors may be eliminated.

Depending on the number of neighbors selected, overfitting can occur. Overfitting occurs when the error rate of the test sample is high and the error rate of the training sample is low. Thus, the error rates of the test and training sample were evaluated in order to evaluate if there is overfitting. The error values obtained through KNN application in the test and training sample are presented in Figure 4.

**Figure 4 -** Apparent error for each value of k in the test and training samples.



Source: Authors.

Overfitting was not observed for k between 2 and 15. For k values between 16 and 20, the error rate in the test sample increases (Figure 4). In case of k equal to 3 and 7, the error rates of the test and training samples were small and close. Training samples presented an error rate equal to 11% for k equal to 3 and 12.3% for k equal to 7. Thus, considering error rate evaluation, both the model with k equal to 3 and the model with k equal to 7 are suitable, as they present acceptable and close error rates in the training and test samples, with an accuracy of 83.3 % in the test sample. In the model for k equal to 7, the error rates in the training and test samples are closer than in the model for k equal to 3.

The area under the ROC curve (AUC) was also analyzed in order to verify the most suitable/best model. Figure 5 shows the results for the AUC when KNN was applied to the test and training samples.

**Figure 5 -** Area under the ROC curve (AUC).



Source: Authors.

Observing the Figure 5, the model with the highest AUC, considering the test sample, was the model with k equal to 7 (0.858). Therefore, considering the error rates, and the AUC values, the suitable/best model is the one whose number of neighbors is equal to 7.

In general, the models obtained through KNN method to determining rockfall probability were satisfactory. Considering the test sample, the average error of the models was equal to 21.6% (78.4% of accuracy). For the training sample, the average error of the models was equal to 15% (85% of accuracy). The average AUC values for the test and training samples were 0.80 and 0.86, respectively. Considering the uncertainties arising from geotechnical field surveys and the difficulty of predicting rockfall probability, the obtained prediction models for rockfall probability classes are quite satisfactory.

**3.2 Validation of rockfall probability classes**

After obtaining the best model (k equal to 7), the behavior of this model regarding its errors was verified, since an error of 16.7% is acceptable, but not negligible. Therefore, it is necessary to know the type of error of the model. Thus, slopes incorrectly classified by KNN model in the test samples were analyzed. Table 5 shows these slopes and their classifications.

**Table 5.** Analysis of the test sample errors.

| Sample | Sum of scores | Class predefined by cluster analysis | Class defined by KNN | Probability of sample belonging to class defined by KNN |
|---|---|---|---|---|
| 28 | 21 | High | Medium | 0.63 |
| 40 | 21 | High | Medium | 0.50 |
| 76 | 22 | High | Medium | 0.57 |
| 78 | 21 | High | Medium | 0.71 |
| 82 | 20 | Medium | High | 0.67 |
| 91 | 17 | Medium | Low | 0.50 |
| 166 | 21 | High | Medium | 0.71 |
| 180 | 17 | Medium | Low | 0.67 |
| 181 | 17 | Medium | Low | 0.67 |
| 182 | 17 | Medium | Low | 0.67 |
| 211 | 17 | Medium | Low | 1.00 |

Source: Authors.

All incorrect classifications obtained through KNN model are related to the transition zone (medium probability). In addition, the probability of the sample belonging to the class determined by KNN model is smaller than 70% in 8 samples. Therefore, the KNN confirms that there are uncertainties regarding classification of slopes whose sum of scores is ranging from 17 to 21 and its borderline zone.

Among the eleven slopes in Table 5, the classification provided by KNN model was less conservative than the classification proposed in Table 2 in ten of them, i.e., medium probability rockfall was classified by KNN as low probability rockfall; high probability rockfall was classified as medium probability. This type of error means underestimating the rockfall probability, so that slopes with need of intervention or monitoring could not receive the proper treatment.

Despite the errors shown in Table 5, the KNN model can be considered adequate, as an error of 16.7% calls the attention for the uncertainties of the variables, which were scored according to a situation observed at the field or a range of values. Their values are associated with a description of a situation that better represents the rock mass behavior, in the point of view of the geologist or the engineer.

## 3.3 Determination of the rockfall probability for 2 new slopes

After validating the optimal model with k equal to 7 and understanding the type of classification errors of KNN in rockfall probability, the model was used to classify two new slopes whose probability classes were unknown. These slopes are located in a quartzite mine in São Thomé das Letras city, Minas Gerais State (Brazil).

The slope 1 (Figure 6a) is composed of a homogeneous fresh quartzite. There is one set of discontinuity (the quartzite foliation), whose persistence is higher than 3m and the spacing varies between 3cm and 20cm, being the smallest spacing predominant. The aperture is in the 0.1mm to 1mm range, with granular infilling. The foliation is practically perpendicular to the slope face, with an average orientation of 11/240. The slope face has an orientation equal to 86/206 (dip/dip direction) and the surface is regular, without overhangs; no evidence of rockfall or sliding was observed. Water dripping in the slope face and in the discontinuities was observed.

The slope 2 (Figure 6b) is also composed of a homogeneous quartzite. There are three sets of discontinuities; one of them is the foliation (set 1). The foliation is practically perpendicular to the slope face, with an average orientation of 08/265. The set 2 is the more critical set, because it daylights out of the slope and can cause rock sliding; the average orientation is

64/120. The slope face has an orientation equal to 75/138 (dip/dip direction) and the surface is irregular, with a degree of overhang ranging from 0.6 to 1.2 m; evidences of rock sliding were observed (Figure 6c). The persistence of the discontinuities is higher than 3m. The spacing of the critical set varies between 20cm and 80cm and this set is planar. Discontinuities with aperture higher than 1cm, with granular infilling were observed. The slope was dry during the field surveys and operations were paralyzed on this front, due to evidence of rockfall hazard. Table 6 presents the scores of the variables P1 to P8 for each slope, according to the described characteristics.

**Figure 6 -** Slopes 1 and 2 and rock sliding evidences in slope 2.



Source: Authors.

**Table 6 -** Mine slopes scores.

| Slope | Scores assigned to independent variables | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 |
| 1 | 3 | 1 | 1 | 1 | 2 | 2 | 3 | 3 |
| 2 | 1 | 1 | 3 | 3 | 4 | 4 | 3 | 3 |

Source: Authors.

The sum of the scores of the Slope 1 is 16, thus according to Table 2, the expected rockfall probability is low. The sum of the scores of the Slope 2 is 22, so the expected rockfall probability is high. The predicted rockfall probability by KNN for Slope 1 is low and the probability of this slope belonging to the low class is 87.50 %. The predicted rockfall probability by KNN for Slope 2 is high and the probability of this slope belonging to the high class is 77.80%. Thus, for these two slopes, the algorithm was able to make correct predictions, consistent with the observations in the field.

## 4. Conclusion

This article presented a complete assessment of the performance of the KNN for predicting the rockfall probability in rock slopes, through the analysis of error, accuracy and AUC for different values of k. The choice of the optimal model considered the error rates, overfitting trends and AUC; the suitable/best model is the one that presented the best metrics.

The suitable/best model is the one whose number of neighbors is equal to 7. This model presented an apparent error of 16.7%, accuracy of 83.3% and AUC of 0.858; the highest AUC among all the models tested. The average error rate considering all the tested models was 21.6% and the average AUC was 0.80, which shows that in general, KNN, a simple machine learning technique, presents good results in predicting rockfall probability in rock slopes.

Analyzing the 11 slopes of the test sample incorrectly classified using the best KNN model; it was observed that all errors involved the medium class of rockfall probability. It was also possible to verify that in most of these cases, the KNN achieved a probability of less than 70% that these slopes belonged to the predicted class, proving that there is an uncertainty or transition zone in this type of analysis. As the classification errors were concentrated in the borderline zone of the classes, it can be considered that the trained KNN model is suitable to predict the rockfall probability.

In view of the results presented and the efficiency of KNN to predict the rockfall probability, it is suggested that the research continues through more robust machine learning techniques, such as Artificial Neural Networks, Decision Trees and Random Forest, in order to compare the results with KNN and understand which variables have the greatest impact on the results and which have little impact.

## Acknowledgments

## References

Bar, N., & Barton, N. (2017). The Q-Slope Method for Rock Slope Engineering. *Rock Mechanics and Rock Engineering* 50, 3307–3322 (2017). https://doi.org/10.1007/s00603-017-1305-0.

Barton, N., Lien, R.., & Lunde, J. (1974). Engineering classification of rock masses for the design of rock support. *Rock Mechanics and Rock Engineering*. 6:189–236

Bieniawski, Z. (1989). *Engineering rock mass classifications: a complete manual for engineers and geologists in mining, civil, and petroleum engineering.* 1st ed. John Wiley & Sons.

Fernandes, A. T. dos R., Fonseca, J. L. T., Silva, I. L. P. da, Agamez Arias, P. D. M., Ramos, R. A., & Oliveira, W. D. de. (2021). Avaliação da influência das tensões de barra na previsão de cargas via redes neurais. *Research, Society and Development*, 10(12), e600101220917. https://doi.org/10.33448/rsd-v10i12.20917

Kassambara, A., & Mundt, F. (2020). *Factoextra: Extract and Visualize the Results of Multivariate Data Analyses*. R package version 1.0.7. https://CRAN.R-project.org/package=factoextra

Kassambara, A. (2017). *Practical guide to cluster analysis in R unsupervised machine learning*. STHDA.

Kaufman, L., & Rousseeuw, P. J. (1990). *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley Series in Probability and Statistics. ISBN: 0-47 1-73578-7.

Kubat, M. (2017). *An Introduction to Machine Learning*. Springer Cham. ISBN 978-3-319-63913-0 (eBook). https://doi.org/10.1007/978-3-319-63913-0

Lawrence, L. (1992). *Data Preparation for a Neural Network*. Neural Network. Special Report: A Miller Freeman Publication.

Mascarenhas, T. A. T., Moriel Junior, J. G., Gomes, R. de S. R., & Mello, G. J. (2020). Aplicação de algoritmos de aprendizado de máquina na classificação de Conhecimentos Especializados de Professores de Física. *Research, Society and Development*, 9(11), e86191110584. https://doi.org/10.33448/rsd-v9i11.10584

Naghadehi, M. Z., Jimenez, R., KhaloKakaie, R., & Jalali, S. M. E. (2013). A new open-pit mine slope instability index defined using the improved rock engineering systems approach. *International Journal of Rock Mechanics and Mining Sciences*. https://doi.org/10.1016/j.ijrmms.2013.01.012

Ossani, P. C., Rossoni, D. F., Cirillo, M. Ângelo, & Borém, F. M. (2021). Classificação de cafés especiais usando técnicas de aprendizado de máquina. *Research, Society and Development*, 10(5), e13110514732. https://doi.org/10.33448/rsd-v10i5.14732

Ossani, P. C., Rossoni, D. F., Cirillo, M. Â., & Borém, F. M. (2020). Unsupervised classification of specialty coffees in Homogeneous sensory attributes through machine learning. *Coffee Science*, 15, e151780. 10.25186/cs.v15i1.1780

Pessoa, A. D., Sousa, G. C. L. Araujo, R. da C. de, & Anjos, G. J. M. dos. (2021). Modelo de rede neural artificial para previsão da capacidade de carga de estacas cravadas. *Research, Society and Development*, 10(1), e12210111526. https://doi.org/10.33448/rsd-v10i1.11526

Pierson, L. A., & Van Vickle, R, (1993). *Rockfall Hazard Rating System*. Transportation Research Record N° 1343, National Research Board, Washington, D.C., pp 6-19.

R Core Team (2020). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

Romana, M. (1995). *The geomechanical classification SMR for slope correction*. In: Proceedings of the 8th ISRM congress on rock mechanics, vol 3. 8 p.

Santi, P. M., Russel, C. P., Higgins, J. D., & Spriet, J. I. (2009). Modification and statistical analysis of the Colorado Rockfall Hazard Rating System. *Engineering Geology* 104: 55–65. 10.1016/j.enggeo.2008.08.009

Santos, A. E. M., Lana, M. S. & Pereira, T. M. (2022). Evaluation of machine learning methods for rock mass classification. *Neural Computing and Applications* **34,** 4633–4642 (2022). https://doi.org/10.1007/s00521-021-06618-y

Santos, A. E. M., Lana, M. S., & Pereira, T. M. (2021). Rock mass classification by multivariate statistical techniques and artificial intelligence. *Geotechinical and Geological Engineering* 39:2409–2430. https://doi.org/10.1007/s10706-020-01635-5

Santos, T. B., Lana, M. S., Pereira, T. M., & Canibulat, I. (2019). Quantitative hazard assessment system (Has-q) for open pit mine slopes. *International Journal of Mining Science and Technology*. Volume 29, Issue 3, Pages 419-427, ISSN 2095-2686, https://doi.org/10.1016/j.ijmst.2018.11.005

Silva, C. C. da, Lima, C. L. de, Silva, A. C. G. da, Moreno, G. M. M., Musah, A., Aldosery, A., Dutra, L., Ambrizzi, T., Borges, I. V. G., Tunali, M., Basibuyuk, S., Yenigün, O., Jones, K., Campos, L., Massoni, T. L., Silva Filho, A. G. da, Kostkova, P., & Santos, W. P. dos. (2021). Predição de casos de Dengue, Chikungunya e Zika em Recife, Brasil: uma abordagem espaço-temporal com base em condições climáticas, notificações de saúde e aprendizado de máquina. *Research, Society and Development*, *10*(12), e452101220804. https://doi.org/10.33448/rsd-v10i12.20804

Venables, W. N. & Ripley, B. D. (2002) Modern Applied Statistics with S. Fourth Edition. Springer.