

SOIL Discuss., author comment AC1
<https://doi.org/10.5194/soil-2022-17-AC1>, 2022
© Author(s) 2022. This work is distributed under
the Creative Commons Attribution 4.0 License.

Reply on CC3

Danilo César de Mello et al.

Author comment on "Weathering intensities in tropical soils evaluated by machine learning, clusterization, and geophysical sensors" by Danilo César de Mello et al., SOIL Discuss., <https://doi.org/10.5194/soil-2022-17-AC1>, 2022

Guilherme Oliveira, 31 Aug 2022

This article represents an excellent application of Machine learning techniques in geophysical-soil science data. The manuscript brings an excellent approach to the geophysical relationships with processes that occur in the soil (weathering) and variables of easy acquisition.

A: Thank you for acknowledging our research.

The methodological flowchart is very well defined; However, I had doubts in some specific parts.

A: Thanks for the suggestions.

First of all, the English writing needs to be improved. I suggest submitting to a specific proofreading company.

A: We sent the for a general review of English (American English) to a specialized company, where a geoscience specialist also reviewed the entire manuscript (Proofreading service). A certificate attesting to the new revision of the manuscript was inserted in the "supplementary material" field.

Why did you use the F1-score test instead of the Kappa metric or accuracy metric?

A: The accuracy is an appropriate and good metric to evaluate the model's performance, when the data is balanced. That is not our case. As the geophysical data used, as well as the data from the physicochemical analysis of the soils, are unevenly distributed in terms

of the number of samples, over the different geologies and soil classes, our data were considered unbalanced.

Imbalanced data always becomes one of the classification issues. Imbalanced dataset occurs when one or several classes have less sample (the minority); while another/other classes are the majority, such as agricultural fields. Accuracy is the most intuitive evaluation index to show the performance of a model. However, when the data classes are imbalanced, a supplementary index is required. That index is the F1-Score, which is a harmonic average of Precision and Recall. The F1-Score can accurately evaluate the performance of the model when the data is imbalanced (Lee and Park, 2021; Wardhani et al., 2019). To further support this assertion, we have inserted two citations from recent works in good journals in the text that detail the use of the F-1 Score in terms of accuracy.

Is the data balanced? The accuracy is a good metric when the data is balanced. This is very important because all of these algorithms have a bias with unbalanced data. Moreover, the algorithms used (excepted for RF) required standardized data.

A: No, the geophysical and soil physico-chemical data are unbalanced. We understand the reviewer's point of view and we agree. Due to that, the results and discussion section were based on F-1 Score metrics, instead of accuracy. We only presented the accuracy and kappa for readers to compare. We used the F-1 Score metric due to criticisms and even recommendations from some other reviewers on articles previously submitted for publication in this and other journals. For this same reason, we also used 2 other parameters to complement the F-1 Score (sensitivity and specificity of the models).

Why didn't you employ the covariates in the unsupervised clustering? I dont understand why you compare the number of PCA dimensions with the number of clusters.

A: We did not employ the covariates in the unsupervised clustering because they were used to map the groups of evaluated attributes. If these covariates were used in clustering, they would act as predictor and predicted variables at the same time, generating overestimated performance results that do not match reality.

We did not compare the optimal number of clusters with the cps results. We use the two data together to create the groups using the k-means method. The use of cps instead of pure data is because the cps are not correlated with each other. On the other hand, pure data can present correlation between them.

With proper revisions, I believe it to be an important contribution to the journal.

A: Thank you for all the suggestion. We made a few changes to the text to better clarify future readers and they were certainly great contributions to our work.