

The Use of Classification Machine learning for Cone Penetration Test Interpretation in a Singapore Context

Tan, E, Loo, A
Mott Macdonald, Singapore

ABSTRACT: The idea and usage of machine learning has become increasingly popular over the last decade. Leveraging on this movement, we aim to apply these machine learning techniques to geotechnics workflow to increase work efficiency and accuracy. In this paper, we will be evaluating the use of multiclass random forest (RF) and support-vector machine (SVM) algorithms in predicting geological classification from cone penetrometer test (CPT) outputs. This study will predominantly be focused on the following geological units sampled in the east of Singapore - Fill, Made Ground, Marine Clay (M), Estuarine Materials (E), Fluvial Sand (F1) and Fluvial Clay (F2) and the Old Alluvium (OA). At this experimental and exploration stage, this machine learning tool is capable of making contextualized predictions, which when coupled with the geologist's technical experience, can aid in geological interpretation with time savings.

1 INTRODUCTION

1.1 Introduction

The use of machine learning algorithms for data analysis has surged exponentially over the recent years. Coupled with the increasing availability and digitisation of site investigation data, machine learning techniques presents us with the opportunity to predict geological information, verify incoming data and optimize geotechnical design processes.

Current day practices of Cone Penetrometer Test (CPT) interpretation rely on both raw outputs (Begemann, 1965), like tip resistance (Q_c), skin friction (F_s) and porewater pressure (U_2), and empirical correlations (Robertson, 2009; Scheider, 2008) to characterise the material. This, combined with a geologist's knowledge of local ground conditions and adjacent borehole information, help to integrate CPT data into the local ground models. While software are available to deduce materials based on their extracted parameters, these are often general classifications (e.g. sand, silt, or clay). Work is still required after this stage to interpret these general classifications to adhere to that of the local geology.

Hence, this paper aims to introduce machine learning as a tool for speeding up CPT interpretation processes, whereby users can pass CPT data through the model to get specific localized soil units unique to their project location. For example, instead of classifying materials as 'SAND' which is general, the model would identify the material as 'FILL' or 'F1', based on contextualized geology in Singapore. It is intended as an additional tool for geologists to refer to when interpreting CPT, which saves time.

Our study builds on previous studies which apply machine learning to site investigation. Machine learning and other soft computing techniques have been used to derive geotechnical parameters (Sushama & Bindhu, 2016; Ardakani & Kohestani, 2015; Javadi, et al., 2012), interpret well-logs data (Timur, et al., 2018) and classify CPT data (Carvalho & Ribeiro, 2019; Ghaderi, et al., 2018; Cho, et al., 2019). In particular, Ghaderi et. al., Cho et al. and Carvalho and Ribeiro have applied artificial neural networks (ANN), decision tree classifier, and distance-weighted nearest neighbour (DNN) algorithms to classify CPT data. Unlike their studies which focus on the broader soil type or index classifications, we have

developed a locally contextualized model and will be evaluating its capability to classify specifically Singapore geology. The two supervised classification machine learning techniques that we will be using are Random Forest Classifier (RF) and Support Vector Machine (SVM). These algorithms were chosen due to their transparency, robustness, simplicity (Üstüner, et al., 2016; Javadi, et al., 2012). Our approach to obtain contextualized geological predictions was conducted due to the limited availability of literature in this scope.

1.2 Site Geology

Our study utilizes CPT data taken from a project based in the eastern region of Singapore. The geology identified includes Anthropogenic Fill (FILL), Made Ground (MG), Kallang Formation/Kallang Group (consisting of Marine Clay (M), Estuarine Materials (E), Fluvial Sand (F1) and Fluvial Clay (F2)) and the Old Alluvium/Bedok Formation (consisting of O(A)/BD(A), O(B)/BD(B), O(C)/BD(C), O(D)/BD(D) and O(E)/BD(E)).

Based on BS 5930:2015, Fill is often referred to for materials that are artificially deposited for engineering purposes. In Singapore, Fill is often used to classify sand deposits from land reclamation projects. While this material is usually found at shallow depth, it has been spotted at depths up to 20m.

Made Ground (MG), unlike Fill, is placed without engineering control or purpose (BSI, 2015). These materials consist of bricks, debris, sand-clay mixture, random clay pockets and any other materials that are not often associated with Fill or Kallang Formation/Group. This material tends to overlie Kallang Formation/Kallang Group and has been found at depths of up to 30m.

The Kallang Formation (currently known as the Kallang Group based on the 2020 classifications (Chua, et al., 2020)) consists of estuarine (E), alluvial sand (F1), alluvial clay (F2), marine clay (M) and beach sand (B). All units except beach sands (B) have been observed in our study area. The Kallang Formation is usually extensive in the east and have been recorded at up to nearly 50m depth. This material was speculated to have been deposited from the late Pleistocene until present day (DSTA, 2009).

Finally, underlying all the other material is the Bedok Formation (Chua, et al., 2020). The formation was previously known as the Old Alluvium (DSTA, 2009). This layer is composed of 5 different weathering grades. Table 1 below summarizes the different weathering grades. The usual site practice in Singapore classifies Old Alluvium into its weathering grade based on its SPT N values.

Table 1. Bedok formation weathering groups (DSTA, 2009)

Bedok Formation (2020 Classification)	Old Alluvium (2009 Classification)	SPT N
BD(E)	0-10	
BD(D)	O(D)	10-30
BD(C)	O(C)	30-50
BD(B)	O(B)	50-100
BD(A)	O(A)	100

Note that for this paper, we will be using the older convention to address these geological units.

2 METHODOLOGY

2.1 Data inputs

For our study, a total of 54 CPTs were used to train and test the model. The data was obtained from a project based in the east of Singapore. The CPT has a sampling interval of 0.1m, which amounts to a total of 5343 samples in our dataset. Each CPT reading has been independently interpreted by a geologist with the help of the CPeT-IT processing software, while referencing neighbouring boreholes.

As per standard practice when training supervised-learning algorithms, 30% of the total data was isolated

as the holdout dataset, defined as the final dataset used for evaluating model performance. Of the remaining 70% of the data, 70% was used to train the algorithm and 30% was set aside as the testing set

Table 2. Diagram of Training set vs Testing set vs Holdout set

Training + Testing set (70%)		Holdout set (30%)
Training set (70%)	Testing set (30%)	

The trained model would predict the output class of the testing set, based on the testing set's input features. In the case of our study, the output class would be geology, while input features are CPT test result readings (refer to Table 3). The predictive outputs of the trained model were then evaluated against interpreted geologies of the testing set. The weighted average F1-score, precision and recall were used as a measure of the model's performance.

Table 3. Sample of input and output data

Input (features)				Output (class)
Z (m)	Qc (MPa)	Fs (MPa)	U2(MPa)	Geology
-6.5	13.904	0.0385	0.084	FILL
-6.6	6.4129	0.0336	0.085	FILL
-6.7	2.3222	0.0697	0.0901	FILL
-6.8	0.7809	0.0276	0.1196	M
-6.9	0.7809	0.0146	0.1356	M

Many iterations of the model were produced in attempt to increase the performance of the predictive model. These include varying the algorithms used, experimenting with different combinations of input features, feature engineering, and parameter tuning. We present our findings and models with highest performance in Section 3.

2.2 Machine Learning Algorithms and Data Split

Two algorithms were employed in this study – random forest classifier (RF) and Support Vector Machine (SVM).

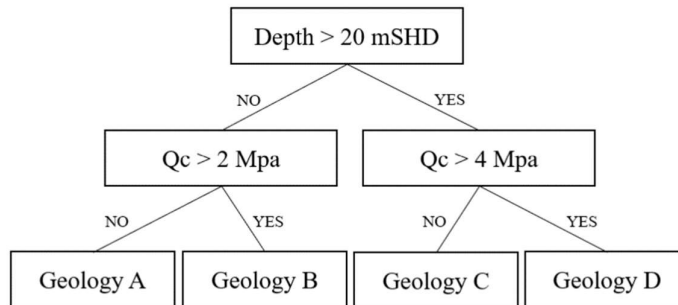


Figure 1. Sample of a simplified decision tree (adapted from: Carnegie Mellon University, n.d.)

Random forest (RF) is a supervised classification algorithm that utilises bagging (Breiman, 2001). This classifier makes predictions by constructing a “forest” comprising of many decision trees (Adam, et al., 2014). A decision tree as seen Figure 1 breaks the sample population into groups as it moves down the tree. At the final branch, it classifies the sample into its outputs class, in this case geology. There are several parameters that controls RF – the number of trees, sampling method, depth of the tree, minimum samples in the terminal node (branch), minimum sample for branching to occur. With the exception of sampling method, these were parameters are analysed and chosen manually, to prevent overfitting. All RF models, in this study, utilizes bootstrapping.

Support Vector Machine (SVM) is a classification model which establishes a hyperplane to divide non-linearly separable data based on maximum margin principles (determining the maximum distance between support vectors) (Petropoulos, et al., 2012). SVM are developed to search for a hyperplane in a multidimensional space using a kernel function where samples in low-dimension space can be reordered in high-dimension space (Karatzoglou, et al., 2006), as illustrated in Figure 2.

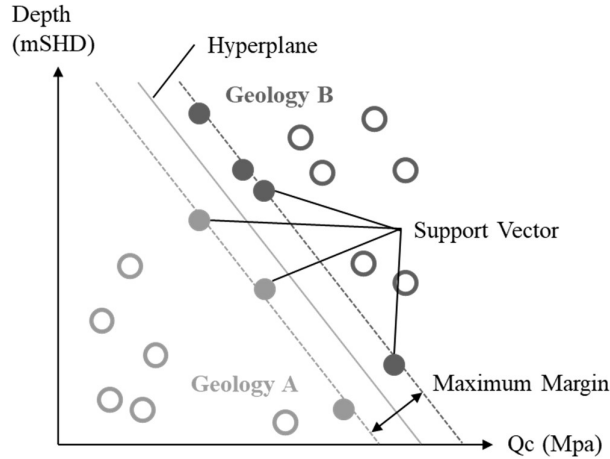


Figure 2. Sample of a 2-dimensional SVM classifier (adapted from: Java T Point, 2021)

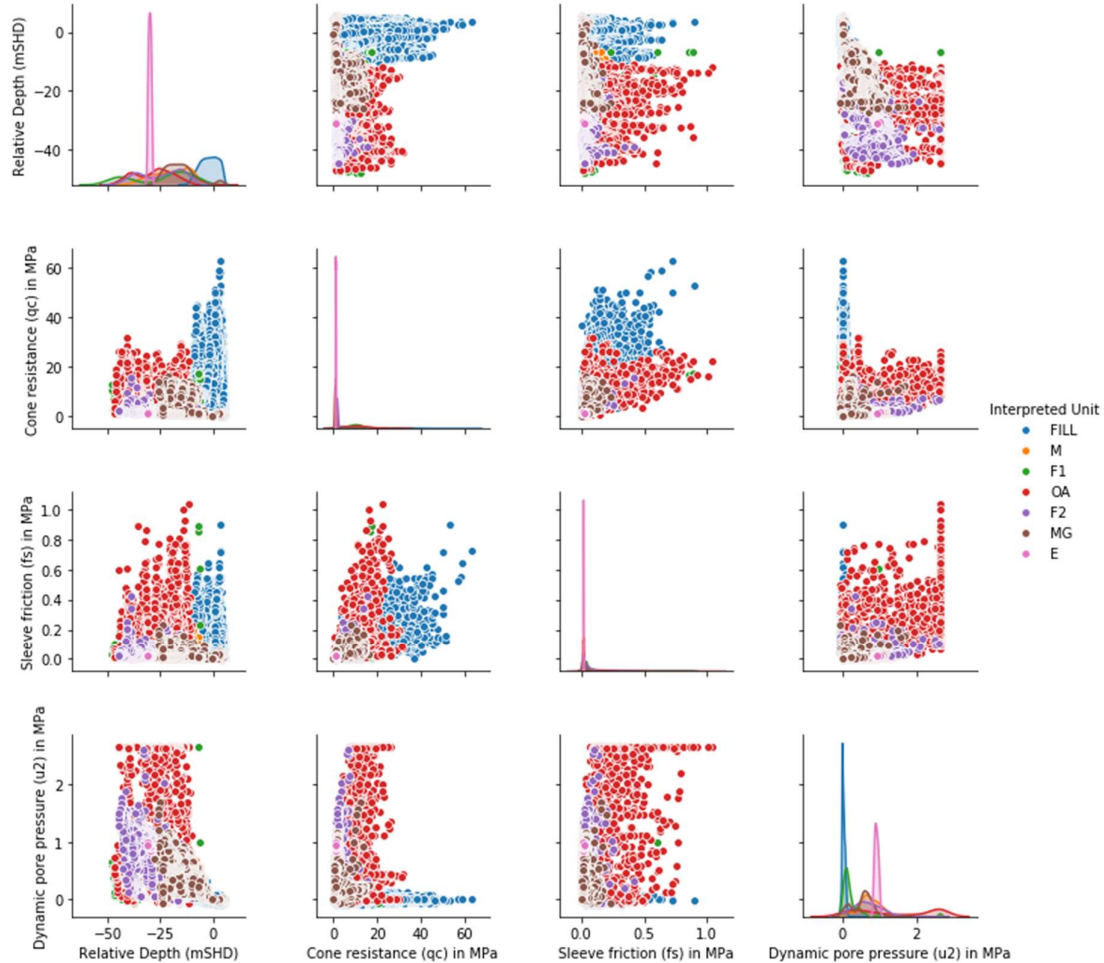


Figure 3. Cross-plots for Qc, Fs, U2 and Z

With reference to cross-plots of our datasets in Figure 3, our preliminary observation indicates the absence of a clear-cut hyperplane to simply classify our data. Hence, non-linear SVM offers a solution that is worth applying to our study.

2.3 Success Metrics

To evaluate our models, we utilised precision, recall and F1-score as our success metrics (Bonnin, 2017). These success metrics are reviewed at every iteration as we continually refined our models during parameter tuning or feature engineering. To understand these metrics better, we first need to understand the confusion matrix.

Table 4. Confusion Matrix

		Predicted	
		True	False
Actual	True	True Positive	False Negative
	False	False Positive	True Negative

The confusion matrix in Table 4 depicts the four possible outcomes of each prediction. In the event where the predicted outcome corresponds to the actual value, we term this as true positive, and vice-versa in the case of true negative. When both true and predicted values do not match, depending on results (refer to Table 4) it can be either false positive or negative.

The success metrics leverages on this confusion matrix for calculation. Precision, for a start, allows us to evaluate the performance of our model is when the prediction results returns positive (Bonnin, 2017). The formula of precision is given as:

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} \quad (1)$$

Recall, on the other hand, evaluates the model in terms of its ability to predict a positive value correctly (Bonnin, 2017). The formula of the recall is given as:

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative} \quad (2)$$

Finally, we have F1-score. F1-score is basically the weighted average of precision and recall (Bonnin, 2017). This is the metrics that will be referred to for most part of our paper.

$$F1\ Score = 2 \frac{Precision * Recall}{Precision + Recall} \quad (3)$$

3 RESULT AND DISCUSSION

3.1 Imbalanced datasets contribute to skewed results

The confusion matrix for the RF and SVM are shown in Figure 4 and Figure 5 respectively.

Both models utilised tip resistance (Q_c), skin friction (F_s), pore water pressure (U_2) and depth in mSHD (Z) as input features. From the confusion matrices, the models are consistent in correctly predicting FILL, M, and Old Alluvium material, which also holds the highest population percentages in the dataset. These classes with high populations in the dataset are usually referred to as the majority class. Conversely, among minority classes which are less observed in the dataset such as E and F2, misclassification commonly occurs. Having varying class sizes in our population poses a major issue for classification algorithms, as the smaller pool is often neglected (Awad & Khanna, 2015).

In our holdout dataset, we have 2032 occurrences of FILL, but only 6 occurrences of E. For future work, more considerations have to be put into creating a holistic population. Proposed improvements will be further explained in Section 4.1.

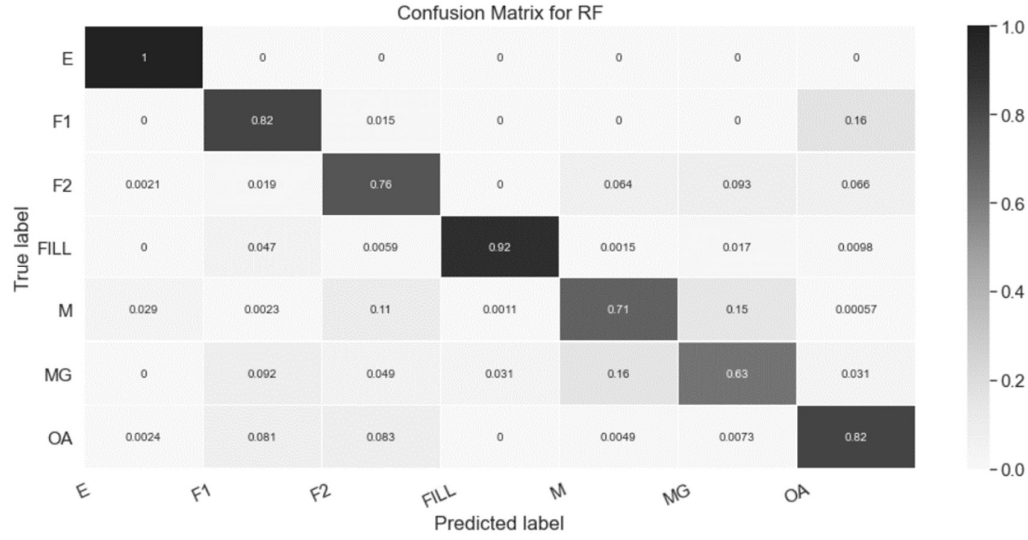


Figure 4. Confusion Matrix for Random Forest Model

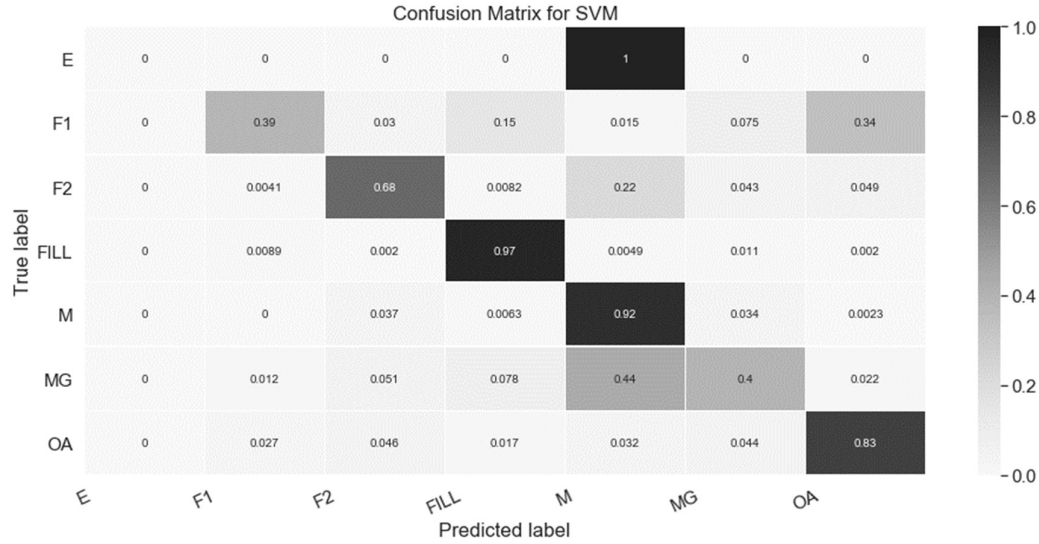


Figure 5. Confusion Matrix for SVM

The second limitation lies in the similarity of the various geological units. M, E and F2 have similar clayey properties, while F1, FILL and MG in the area are predominantly sand. These geological units are often distinguished in the field by its colour, and lab tests results such as Atterberg limits and moisture content. These characteristics, however, were not well captured by CPT data. Hence, feature engineering may be required to better constrain these properties.

3.2 Incorporating Relative depth (Z) returns stronger predictions

We also observed that incorporating relative depth (Z) as a feature results in considerably stronger predictions (0.05 to 0.09 increase in weighted average F1-score for Case 1 vs 2 of Table 5 and Table 6). A similar observation was also made in Carvalho & Ribeiro (2019) study.

Including depth as an additional feature helps the model distinguish seemingly similar units. Taking Fill, F1 and Old Alluvium units as an example, these three materials are often described as predominantly sand. Thus, their CPT readings tend to show similar Q_c and F_s values which makes it difficult for the model to differentiate between units. To improve data separation, we leverage on their known depositional sequences. Out of the three units, the oldest and deepest layer - Old Alluvium was deposited before the late Pleistocene epoch (DSTA, 2009). Overlying this layer is the F1 alluvial sand within the Kallang Formation, which was laid down between the late Pleistocene and Holocene (Chua, et al., 2020). Subsequently, Fill was anthropogenically deposited on natural ground over the last century, making it the youngest and shallowest unit (DSTA, 2009). Hence, by inputting depth as a proxy for geological ages, the models are able to achieve better data separation and produce better predictions. This is evidenced by confusion matrix in the earlier section where Fill, F1 and Old Alluvium attained relatively higher prediction scores. Despite its better performance, it is also imperative to consider the wider applications and scalability of the model. Given the inherent spatial variability of geology from site to site, less emphasis should be placed on relative depth as an input feature, especially if the model is intended for use across multiple sites with drastically different ground conditions. To counter this limitation of our model, we discuss the possibility of incorporating other machine learning techniques as an improvement in Section 4.

3.3 Merging the weathering units for Old Alluvium

Another option explored was merging the weathering units for Old Alluvium. The results are presented in Table 7 and Table 8 for RF and SVM respectively.

Referring to the RF classifier, there was a slight increase in weighted-average precision scores and a decrease in the weighted-average recall values. The F1-score values also showed a marginal decrease by 0.01. This marginal change can be deemed to be inconclusive as they could easily be overshadowed by model variation, where each run yields slightly varying results, usually differing by a value of 0.01 to 0.02. In contrast, the SVM model saw a 0.04 increase in weighted-average F1-score. As this increase is greater than a natural model variation of 0.2, we will deem this result to be conclusive.

When we merge individual Old Alluvium weathering units, the goal was to collectively increase its population class. However, merging the Old Alluvium units gives rise to greater variability of input features within the class. For example, the Q_c value of Old Alluvium for one sample may be 20 MPa as it was based on unweathered O(B), while another sample may be 1 MPa as it was based on extremely weathered O(E). This variability appears to be better accommodated in SVM, as evidenced by the higher F1-score in comparison to results recorded in RF.

3.4 Comparison of RF vs SVM algorithms

The weighted-average and macro-average F1 scores of Case 2 have been summarized in Table 9 and Table 10 respectively. SVM classifier (0.84) recorded a slightly higher weighted-average F1-score in comparison to RF (0.81). In terms of its macro-average F1 score, SVM (0.61) fell short of RF (0.62) by 0.01. These subtle differences are evident in the confusion matrix shown in Figure 4 and 5. Looking at the matrix, RF showed promising predictions across all 7 output classes while SVM tends to perform better for the majority classes. SVM, however, fell short in prediction of minority classes.

From our observations, the RF classifier is more robust in accommodating minority factions and outliers. Similar studies on RF have also highlighted the robustness and consistency of these ensemble models (Houunkpatin, et al., 2018). On the other hand, SVM classifiers are more capable in classifying for the majority classes, often at the expense of the minority classes. A study conducted by Adam et al. (2014) have also suggested that SVM often reduce classification errors without accounting for their distribution. This was evident in our case.

On the surface, while it seems both models are relatively comparable in terms of F1-scores, there are underlying differences when we evaluate performance in predicting majority vs minority classes.

Table 5. Comparison between Random Forest Classifier Models with data inputs Qc, Fs, U vs. Qc, Fs, U, Z

Case	Bedok Formation (Old Alluvium) Inputs	Data Input	Number of Trees	Tree Depth	Minimum number of sam- ples to split	Minimum number per leaf	Macro Aver- age Preci- sion	Weighted Average Precision	Macro Aver- age Recall	Weighted Average Recall	Macro Aver- age F1	Weighted Average F1
1	Grouped Bedok Formation	Qc, Fs, U	200	15	100	5	0.54	0.8	0.71	0.74	0.56	0.76
2	Grouped Bedok Formation	Qc, Fs, U, Z	200	15	100	5	0.59	0.85	0.81	0.8	0.62	0.81

Table 6. Comparison between SVM Models with data inputs Qc, Fs, U vs. Qc, Fs, U, Z

Case	Bedok For- mation (Old Allu- vium) Inputs	Data Input	C	gamma	Macro Average Precision	Weighted Av- erage Preci- sion	Macro Average Recall	Weighted Av- erage Recall	Macro Average F1	Weighted Av- erage F1
1	Grouped Bedok Formation	Qc, Fs, U	1000	0.1	0.58	0.77	0.48	0.79	0.48	0.75
2	Grouped Bedok Formation	Qc, Fs, U, Z	1000	1	0.63	0.84	0.6	0.85	0.61	0.84

Table 7. Comparison between Random Forest Classifier Models with grouped and ungrouped Bedok Formation weathering units

Case	Bedok Formation (Old Alluvium) Inputs	Data Input	Num- ber of Trees	Tree Depth	Mini- mum number of sam- ples to split	Mini- mum number per leaf	Macro Average Precision	Weighted Average Precision	Macro Average Recall	Weighted Average Recall	Macro Average F1	Weighted Average F1
2	Grouped Bedok Formation	Qc, Fs, U, Z	200	15	100	5	0.59	0.85	0.81	0.79	0.62	0.81
3	Bedok Formation in its Weather Grades	Qc, Fs, U, Z	200	15	50	1	0.55	0.84	0.72	0.81	0.59	0.82

Table 8. Comparison between SVM Models with grouped and ungrouped Bedok Formation weathering units

Case	Bedok Formation (Old Alluvium) Inputs	Data Input	C	gamma	Macro Av- erage Pre- cision	Weighted Aver- age Precision	Macro Av- erage Re- call	Weighted Aver- age Recall	Macro Av- erage F1	Weighted Aver- age F1
2	Grouped Bedok Formation	Qc, Fs, U, Z	1000	1	0.63	0.84	0.6	0.85	0.61	0.84
3	Bedok Formation in its Weather Grades	Qc, Fs, U, Z	1000	0.1	0.65	0.82	0.55	0.83	0.57	0.80

Table 9. Comparison between Random Forest Classifier and SVM Models

Case	Algorithm	Bedok Formation (Old Alluvium) Inputs	Data Input	Weighted Average Precision	Weighted Average Recall	Weighted Average F1
2	Random Forest Classifier	Grouped Bedok Formation	Qc, Fs, U, Z	0.85	0.79	0.81
2	SVM	Grouped Bedok Formation	Qc, Fs, U, Z	0.84	0.85	0.84

Table 10. Comparison between Random Forest Classifier and SVM Models

Case	Algorithm	Bedok Formation (Old Alluvium) Inputs	Data Input	Macro Aver- age Preci- sion	Macro Average Recall	Macro Average F1
2	Random Forest Classifier	Grouped Bedok Formation	Qc, Fs, U, Z	0.59	0.81	0.62
2	SVM	Grouped Bedok Formation	Qc, Fs, U, Z	0.63	0.6	0.61

4 FUTURE WORKS

4.1 Sampling and Pre-processing in response to Class Imbalance

To enhance the robustness of the models, we should ideally be exposing our models to a larger pool of evenly distributed data classes. One suggestion would be to incorporate data collected from other project sites into the model building process. This serves to increase the data population size, and increase the variability of the data samples to provide a better and wider representation of ground conditions. However, it is worthy to note that expanding the dataset size does not necessarily translate to having a well distributed sample. In Singapore, Kallang units like estuarine layers tend to be rare relative to marine clay, fill or the old alluvium layers. Hence, distribution gaps will always persist.

A common strategy which Kubat & Matwin recommends would be balancing a dataset, this includes decreasing the size of the majority class through random undersampling. However, this results in information loss for the majority class, and should be mobilized at discretion (i.e. when cost of information loss for minority class outweighs that of majority class).

4.2 Elevating and Expanding on Model Inputs

To increase the performance of the model, we could either improve on the current features or increase the number features used. A possible improvement that can be made to our current model would be the replacement of raw data with corrected data. Tip resistance (Qc) could be replaced with correct tip resistance values (Qt).

Feature engineering can also be employed to generate new and insightful features. Possible features include the use of SBTn values and distance-based variables. Carvalho & Ribeiro applied distance-based machine learning techniques to soil classification systems (Carvalho & Ribeiro, 2019). The study utilises K-nearest neighbour and distance-weighted nearest neighbour techniques, along with Qc, Fs, U2 and Z, to characterise the data.

4.3 Exploring other Possible Machine Learning Algorithms

Apart from RF and SVM, other algorithms can be explored. These includes artificial neural networks (ANN), logistic regression, and other sampling methods (like boosting) for ensemble classifiers.

5 CONCLUSIONS

In this study, we analysed and evaluated the performance of two models, RF and SVM, with varying input features. Using F1-scores as a success metric, we observed that both RF and SVM model performed better when depth was included as a feature. We also analysed the models with merged Old Alluvium units and saw an increase in model performance for SVM model. Similar improvements were not reflected in the RF model. Lastly, we compared the two models and determined that the RF models are more robust and consistent than the SVM models. The choice of which algorithms to use should be governed by the end goal. RF should be used if the focus is to capture the minority class (like E, F2 and F1) for the project. In the event where the majority class (like M, Old Alluvium and Fill) is the focus, SVM might prove to be a more effective option. Despite its limitations, these models can and should still be referred to as an alternative tool to aid and speed up geotechnical interpretation especially for CPT where neighbouring boreholes are not available.

6 REFERENCES

- Adam, E. M., Onisimo, M., John, O. & Elfatih, A.-R., 2014. Land-use/cover classification in a heterogeneous coastal landscape using RapidEye imagery: evaluating the performance of random forest and support vector machines classifiers. *International Journal of Remote Sensing*, 25(10).
- Ardakani, A. & Kohestani, V., 2015. Evaluation of liquefaction potential based on CPT results using C4.5. *Journal of AI and Data Mining*, 3(1), pp. 85-92.
- Awad, M. & Khanna, R., 2015. Support Vector Machines for Classification. In: *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers*. s.l.:Apress, pp. 39-66.
- Begemann, H., 1965. *The friction jacket cone as an aid in determining the soil profile*. Montréal, 6th International Conference on Soil Mechanics and Foundation Engineering .
- Bonnin, R., 2017. Model Implementation and Results Interpretation. In: *Machine Learning for Developers*. s.l.:Packt Publishing Ltd, pp. 54-56.
- Breiman, L., 2001. Random Forests. *Machine Learning*, Volume 45, pp. 5-32.
- BSI, 2015. *Code of practice for ground investigations BS 5930:2015+A1:2020*, s.l.: BSI.
- Carnegie Mellon University, n.d. *Decision Trees*. [Online] Available at: <https://www.cs.cmu.edu/~bhiksha/courses/10-601/decisiontrees/> [Accessed 24 August 2021].
- Carvalho, L. & Ribeiro, D., 2019. Soil Classification System from Cone Penetration Test Data Applying Distance-Based Machine Learning Algorithms. *Soils and Rocks*, 42(2), pp. 167-178.
- Cho, S., Cho, B., SeungMin, K. & Kim, H., 2019. Development of locally specified soil stratification method with CPT data based on machine learning techniques. *Geotechnics for sustainable infrastructure development*, pp. 1287-1296.
- Chua, S. et al., 2020. A new Quaternary stratigraphy of the Kallang River Basin, Singapore: Implications for urban development and geotechnical engineering in Singapore. *Journal of Asian Earth Sciences*, Volume 200.
- DSTA, 2009. *Geology of Singapore*. 2nd ed. Singapore: DSTA.
- Ghaderi, A., Shahri, A. A. & Larsson, S., 2018. An artificial neural network based model to predict spatial soil type distribution using piezocone penetration test data (CPTu). *Bulletin of Engineering Geology and the Environment*, Volume 78, pp. 4579-4588.
- Gupta, N., 2013. Artificial Neural Network. *Network and Complex Systems*, 3(1).
- Houkpatin, K. O. L. et al., 2018. Predicting reference soil groups using legacy data: A data pruning and Random Forest approach for tropical environment (Dano catchment, Burkina Faso). *Scientific Reports*, 8(9959).
- Java T Point, 2021. *Support Vector Machine Algorithm*. [Online] Available at: <https://www.javatpoint.com/machine-learning-support-vector-machine-algorithm> [Accessed 24 August 2021].
- Javadi, A., Ahangar-Asr, A., Johari, A. & Faramarzi, A. T. D., 2012. Modelling stress-strain and volume change behaviour of unsaturated soils using an evolutionary based data mining technique, an incremental approach. *Engineering Applications of Artificial Intelligence*, Volume 25, pp. 926-933.

- Karatzoglou, A., Meyer, D. & Hornik, K., 2006. Support Vector Machines in R. *Journal of Statistical Software*, 15(9).
- Kubat, M. & Matwin, S., 1997. *Addressing the Curse of Imbalanced Training sets: One-sided Selection*. Ottawa, Department of Computer Science, University of Ottawa.
- Nichols, G., 2009. *Sedimentology and Stratigraphy*. 2nd ed. s.l.:John Wiley & Sons Ltd.
- Petropoulos, G., Kalaitzidis, C. & Krishna, P. V., 2012. Support vector machines and object-based classification for obtaining land-use/cover cartography from Hyperion hyperspectral imagery. *Computers and Geosciences*, Volume 41, pp. 99-107.
- Robertson, P., 2010. *Soil Behaviour Type from the CPT: an update*, Signal Hill, California: Gregg Drilling & Testing Inc..
- Robertson, P. & Cabal, K., 2014. *Guide to Cone Penetration Testing for Geotechnical Engineering*. Signal Hill, California: Gregg Drilling & Testing, Inc..
- Robertson, R., 2009. Interpretation of cone penetration test - a unified approach. *Canadian Geotechnical Journal*, pp. 1337-1355.
- Scheider, e. a., 2008. Analysis of Factors Influencing Soil Classification Using Normalized Piezocone Tip Resistance and Pore Pressure Parameters.. *Journal of Geotechnical and Geoenvironmental Engineering*, pp. 1569-1586.
- Sushama, K. & Bindhu, L., 2016. Modelling of soil shear strength using neural network approach. *International Journal of Earth Sciences and Engineering*, 8(05).
- Timur, M., Rassul, Y. & Amirgaliyev, Y., 2018. *Machine Learning Algorithms for Classification Geology Data from Well Logging*. s.l., Conference: 2018 14th International Conference on Electronics Computer and Computation (ICECCO).
- Üstüner, M., Sanli, F. B. & Abdikan, S., 2016. *BALANCED VS IMBALANCED TRAINING DATA: CLASSIFYING RAPIDEYE DATA WITH SUPPORT VECTOR MACHINES*. Prague, Conference: XXIII ISPRS Congress.