

# Predicting the swell potential of expansive clays using machine learning models: a comparative study with experiments

**Uttara D.D Liyanage**, Kumari W.G.P., Jayan S. Vinod

*Faculty of Engineering and Information Sciences, University of Wollongong, NSW 2522, Australia, vinod@uow.edu.au*

Jun Sugawara, Siva Sivakumar

*Department of Transport and Main Roads, Brisbane, Queensland, QLD 4008, Australia.*

Bindumadhava Aery

*Aurecon, Bowen Hills, Queensland, QLD 4006, Australia.*

**ABSTRACT:** Expansive soils can undergo significant volumetric changes in response to variations in moisture content, making it vital to understand this volume change in the field before constructing structures or roads. Laboratory swell tests are time-consuming to conduct to understand the swell characteristics. This study proposes a machine learning approach as an efficient and accurate alternative to traditional empirical methods for predicting swell potential. A dataset of 158 clay soil records was compiled from existing literature, consisting of five soil index properties: liquid limit, plasticity index, clay content, water content, dry density and corresponding swell potential value. Support Vector Regression (SVR), Random Forest (RF), Gradient Boosting Machine (GBM), and Artificial Neural Network (ANN) models were used to build four prediction models. Each model was optimized and tuned with the training dataset and tested separately. Performance of each model was evaluated based on mean squared error (MSE) and the coefficient of determination ( $R^2$ ). From all models, the ANN model performed the best with an MSE of 6.1 and an  $R^2$  of 0.91 for the test dataset. Laboratory results of three independent soil samples were used for further validation, and the ANN model produced the best predictions with an MSE of 11.1 and  $R^2$  of 0.91. The results show the ANN model has the potential to be a reliable method in predicting swell behavior in expansive clays, which can help practicing engineers to readily evaluate the swell potential of widespread natural expansive clays.

**KEYWORDS:** Expansive soil; Clays; Mathematical modelling; Computational geotechnics; Machine learning models

## 1 INTRODUCTION

Expansive soils undergo volume change with moisture content and are primarily found in arid and semi-arid regions where annual evapotranspiration exceeds precipitation (Fityus et al. 2004). Understanding the stability of infrastructure foundations on these soils is crucial due to their volume changes upon wetting and drying. Factors affecting swelling in expansive clay soils include clay mineralogy, plasticity, dry unit weight, initial water content, water content variation, and overburden stress (Çimen et al. 2012, Sridharan and Rao 1973).

Many studies have investigated the influence of these factors on swell potential through laboratory tests. The swell potential can be determined using a one-dimensional swell-consolidation test in an oedometer. These tests are often time-consuming experimental procedures that require expensive equipment and skilled laboratory personnel (Benbouras and Petrisor, 2021). While empirical and semi-empirical correlations have provided reasonable swell potential estimations for decades, they are often limited to a specific soil type. Further, they are unable to make satisfactory predictions that cover a wider range of environmental conditions. (Elbadry, 2017, Erguler and Ulusay, 2003, Taherdangkoo et al. 2023, Vanapalli and Lu, 2012, Yilmaz, 2006). These empirical methods generally depend on using index properties of soil, such as Atterberg limits, clay content, dry density, and moisture content (Sivapullaiah et al. 1996).

Machine Learning (ML) algorithms are becoming more reliable than simple empirical methods due to their ability to identify complex relationships between input parameters and outcomes. They can effectively handle uncertainty in input parameters, making them suitable for geotechnical applications (Chen et al. 2021). Techniques like Artificial Neural Networks (ANN), Support Vector Machine (SVM), and Multivariate Adaptive Regression Splines (MARS) have been used for soil classification and predicting properties such as swell pressure and shear strength (Zhang et al. 2021a, Zhang et al. 2021b).

These models rely on “data-driven decision making,” utilizing training examples with known inputs and outputs (Eyo and Abbey, 2022). Recent studies have developed ML models to predict expansive soil properties using techniques such as ANN, ANFIS and Gene Expression Programming (GEP) to predict swell pressure (Jalal et al. 2021).

Benbouras and Petrisor (2021) explored various methods, including extreme learning machine and deep neural networks, to predict the swell index of Algerian cohesive soils. Chen et al. (2022) used three ML methods: random forests, extreme gradient boosting, and stacked generalization, to predict the uniaxial tensile strength of expansive soil. Teodosio et al. (2023) utilized deep learning to predict the shrink-swell index of Australian cohesive soils. Other applications include pile capacity prediction (Pal and Deswal, 2008), over-consolidation ratio prediction (Samui et al., 2008) and soil moisture estimation (Kashif et al., 2007). More recent studies like Mahdevari et al. (2014) have utilized SVR for predicting tunnel boring machine rates and the bearing capacity of geogrid-reinforced stone columns (Debnath and Dey, 2017).

In this study, a high-quality dataset of 158 soil samples was collected from peer-reviewed literature after preprocessing and removal of missing values and outliers. Soil data included five predictor variables: Liquid Limit (LL), Plasticity Index (PI), Clay Content (CC), Water Content (WC) and Dry Density (DD) and the target variable Swell Potential (SP). Four machine learning models were used, which included two classical models: Support Vector Regression (SVR) and Artificial Neural Networks (ANN), and two ensemble models: Random Forest (RF) and Gradient Boosting Machines (GBM). By testing multiple machine learning algorithms, we aim to identify which model best captures the patterns in swell potential with minimized overfitting, even with a limited dataset. Each model was trained using cross-validation and hyperparameter tuning. Additionally, each model was validated for its performance using independent experimental data obtained from field samples. Model performance was assessed

based on root mean square error (RMSE) and coefficient of determination ( $R^2$ ). The outcome of this study aims to help practicing engineers in making better-informed decisions in civil engineering projects where swell potential in clays is critical for safe and efficient infrastructure management.

## 2 OVERVIEW OF MACHINE LEARNING MODELS

### 2.1 Support Vector Regression (SVR)

The goal of SVR is to find a function that minimizes the error between predicted and actual values. SVR uses the training data to come up with a symmetrical loss function that can ignore high and low outliers within the input data. It then creates a flexible tube around the estimated function, using Vapnik's  $\epsilon$ -insensitive approach, such that any values that are outside the tube are ignored, but those within the tube, either above or below the function, are taken into account (Awad and Khanna, 2015).

For non-linear cases, SVR maps data into a higher-dimensional space using kernel functions that satisfy Mercer's condition, which ensures the kernel function represents an acceptable similarity measure between data points (Clarke et al., 2005). A kernel is typically a nonlinear function that transforms the original input space into a high-dimensional feature space, which will improve linear separability. The optimization problem for non-linear SVR involves minimizing the squared Euclidean norm of the weight vector and a regularization parameter,  $C$ , which controls the trade-off between margin maximization and error minimization (Schölkopf and Smola, 2003).

### 2.2 Artificial Neural Networks (ANN)

ANNs are computational models that are based on how the human brain is built and operates. An ANN is formed by neurons. They are known as the processing units and are linked to one another. These neurons are arranged in layers: an input layer, one or more hidden layers, and an output layer. Every neuron gets inputs ( $x_i$ ), multiplies them by weights that are related to them ( $w_{ji}$ ), adds a bias ( $\theta_j$ ), and then uses an activation function ( $f(\cdot)$ ) to make an output ( $y_j$ ) (Kukreja et al., 2016, Shahin et al., 2001, Zupan, 1994). The process is shown in equation (1) and (2).

$$I_j = \sum w_{ji}x_i + \theta_j \quad (1)$$

$$y_j = f(I_j) \quad (2)$$

In R programming, several parameters can be defined to build an ANN model. The "size" parameter defines the number of neurons in a hidden layer. Higher number of hidden layers will capture complex patterns but can overfit. Using a technique known as L2 regularization, the "decay" parameter will prevent overfitting by gradually pushing the weights to remain small. The maximum number of times the network is permitted to update its weights during training is determined by the "maxit" parameter. A higher "maxit" makes sure the model gets enough training cycles to fully understand the data and perform at its best without stopping the training too soon.

### 2.3 Random Forest (RF)

RF is an ensemble learning method built on the classification and regression trees (CART) algorithm. Unlike Decision trees, RF will build many trees using a random sample from the training data instead of using the same dataset every time. Each tree is built by only considering a random subset of the input variables at each node. This method is used to prevent

overfitting. For regression applications, each tree produces a numerical prediction, and RF will take the average of the predictions to give an output (Breiman, 2001, Segal, 2004).

RF has different hyperparameters that can help optimize the model. The number of predictors chosen at random for every node split is specified by the "mtry" hyperparameter. A smaller "mtry" lowers the chance of overfitting by increasing tree variety. The "ntree" parameter determines the number of trees in the forest; additional trees often boost performance by minimizing variance.

### 2.4 Gradient Boosting Machines (GBM)

GBM is another ensemble method originally used for classification tasks (Friedman, 2001). In order to create a strong model with high accuracy, the basic method is to iteratively combine simple models called "weak learners", which are usually decision trees (Touzani et al., 2018). GBM will check the errors made by the weak learner and build the second model to correct these errors (loss function). This process is repeated while minimizing the loss functions at each step (Natekin and Knoll, 2013).

Hyperparameter "n.trees" controls the number of boosting iterations. Each iteration will add a new tree to the model, and the larger the number of trees, the higher the prediction accuracy. The maximum depth of each tree is determined by "interaction.depth" parameter, which also controls how effectively the model captures interactions between variables; larger values capture complicated connections but raise the possibility of overfitting. Each tree's effect is determined by the "learning rate". Lower values decrease learning and improve generalization, leading to more trees.

## 3 MATERIALS AND METHODS

### 3.1 Dataset

The training dataset comprises 158 records collected from literature (Alazigha et al. 2016; Ashayeri and Yasrebi 2009; Bhuvaneshwari et al. 2010; Çimen et al. 2012; GhavamShirazi and Bilsel 2021; Mishra et al. 2008; Phanikumar and Singla 2016; Puppala et al. 2013; Puppala and Musenda 2000; Rao and Thyagaraj 2003; Sabtan 2005; She et al. 2020; Shelke and Murty 2010; Soltani et al. 2017; Sridharan and Gurtug 2004; Turkoz and Vural 2013; Zamin et al. 2021) with 5 parameters, namely, liquid limit (LL), plasticity index (PI), clay content (CC), water content (WC), and dry density (DD), along with their respective swell potentials (SP). These properties have been extensively investigated by previous researchers and established that the swelling potential can be correlated to these parameters (e.g. Gupta et al. 2008). Moreover, these soil index properties can be easily obtained from conventional, less time-consuming experiments. Considering the importance of data quality, these records were attentively collected from experimental swell test results published in the literature, which followed the standard oedometer swell test procedure (e.g. ASTM D 4546).

Table 1 shows the statistical metrics (mean, standard deviation (S.D.), minimum, and maximum) obtained from the entire dataset of each predictor variable, while Figure 1 shows the histogram frequency distribution. The specified ranges for LL, PI, CC, WC and DD are chosen to focus on a wide range of soil conditions, from low to high plasticity. Although CC, LL and PI are partially correlated, the nonlinear ML models used learn from the combined effect of all five input variables rather than relying on any single parameter. Along with plasticity-related factors (LL, PI, CC), the initial state variables WC and DD are also included so that the predictions can show how moisture condition and density affect swell potential. In the

dataset, swell potential is usually higher when the initial water content is lower. However, there isn't a clear pattern with dry density in the small DD range that was examined. The histograms indicate that the properties are spread across their respective ranges, with WC showing a normal distribution. The data records were split into two groups, where 80% was utilized as training data and 20% for testing.

Table 1. Statistics of all predictors and target variables.

Data type		Mean	S.D.	Min	Max
Predictors	Liquid Limit, %	65.3	16.3	28.0	108
	Plasticity Index, %	40.3	14.8	7.10	71
	Clay Content, %	46.9	15.8	10.0	93
	Water Content, %	19.8	9.0	1.30	40
	Dry Density, kN/m <sup>3</sup>	15.2	1.9	11.5	19.6
Target	Swell Potential, %	11.9	6.9	0.3	34.4

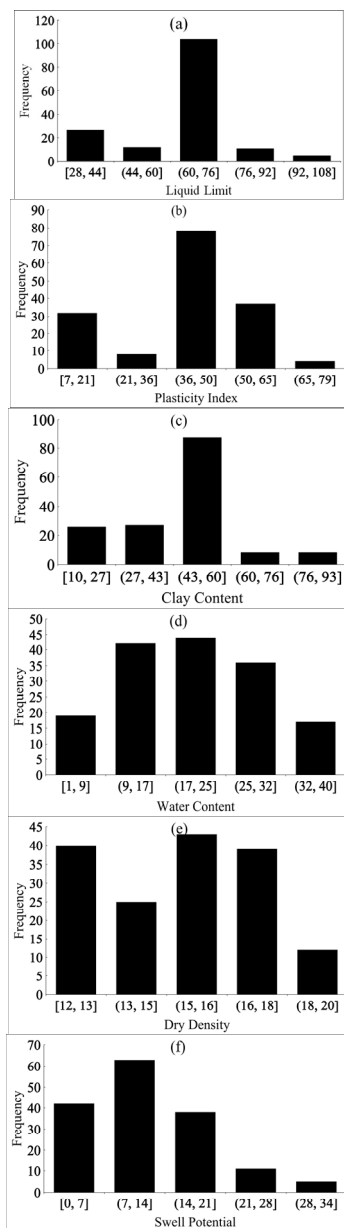


Figure 1. Histograms of frequency distribution for (a) LL (b) PI (c) CC (d) WC (e) DD and (f) SP.

## 3.2 Model Implementation

### 3.2.1 SVR

When implementing a non-linear SVR model, the goal is to find a suitable, optimized value of the objective function. The performance of the SVR model depends upon the selection of a suitable kernel function, the cost parameter “ $C$ ”, and the loss function “ $\epsilon$ ”. Radial kernel was used due to its ability to manage non-linear relationships, unlike other kernels. The optimal values for the hyperparameters  $C$  and  $\epsilon$  were selected using a 10-fold cross-validation test. This test would systematically evaluate different combinations of  $\epsilon$  and  $C$  parameters, by splitting the training dataset into 10 equal parts or folds, iteratively training the model on 9 folds and testing it on the remaining one. On each iteration, the performance metrics are calculated for each fold. The model would utilize these parameters and functions on the training dataset and learn from them to provide accurate predictions.

### 3.2.2 ANN

The ANN model was built using the “nnet” package in R. A standard multilayer neural network was used, where inputs are transmitted forward through a hidden layer to produce a continuous output. The model was tuned by a grid search of various combinations of hidden layer sizes and weight decay values to enhance generalization and prevent overfitting. The optimal ANN configuration was determined based on a repeated cross-validated error analysis. The completed model was assessed using both the training and test datasets.

### 3.2.3 RF

The model was built using the “randomForest” method in the “caret” package in R. A grid search was run to improve performance by varying the “mtry” parameter (the quantity of variables randomly chosen at each split) between 2 and 5. The model was tuned and verified by repeated 10-fold cross-validation with three repeats to ensure optimal hyperparameters were selected, and the final RF model consisted of 500 trees (ntree = 500).

### 3.2.4 GBM

The GBM model was developed using “gbm” component within the “caret” package. A continuous boosting technique was used to construct the model, adding decision trees one after the other to fix errors produced by earlier trees. To maximize performance, the model parameters were adjusted, including the minimum node size, learning rate, tree depth, and number of trees. Ten-fold cross-validation was used for both model training and evaluation to ensure generalization and avoid overfitting.

Section 4 shows the results with performance metrics of each model against the training and testing datasets. The performance of each model was further validated using independent experimental results obtained from laboratory tests done on expansive clay samples.

## 3.3 Experimental program

The expansive soil used in this study was collected from a distressed site in Queensland, Australia. Figure 2 shows the clay samples which were collected from varying depths in the soil layer. Table 2 presents all soil characteristics of the soil samples. For these samples, clay content (CC) was determined using a laser diffraction particle size analyzer on dried soil that had first been dry sieved to a particle size smaller than 1 mm. CC refers to the fraction of particles smaller than 2  $\mu\text{m}$  while LL and PI serve as behavioral indicators that also represent clay mineralogy and activity. Soils with a modest clay content may

have a high PI when the clay minerals are very active (Skempton, 1953). The detailed mineral composition of the clay fraction was not investigated in this study, and this is recognized as a limitation of the current dataset. The ASTM D-4546 procedure was followed to determine the one-dimensional vertical free swell of remolded clay samples. The clay sample was mixed with the respective optimum moisture content and allowed to rest for 24 hours to ensure uniform moisture distribution. Corresponding final swell potential values of each sample are also shown in Table 2.

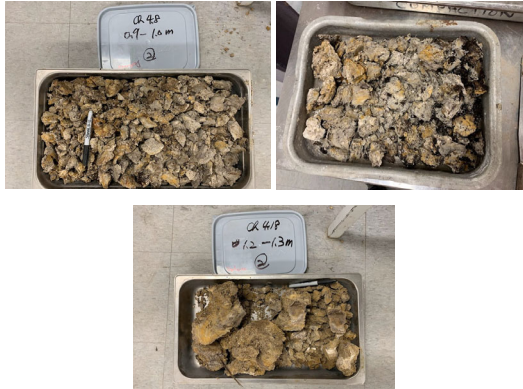


Figure 2. Soil samples obtained from 0.9m, 1.0m and 1.2m depths.

Table 2. Soil index properties of clay samples

Depth (m) / Sample name	0.9 m Sample	1.0 m Sample	1.2 m Sample	Reference
Specific gravity	2.8	2.8	2.8	ASTM D 854-02
Liquid limit (%)	79	81	79	AS 1289.3.9.1-2002
Plasticity index (%)	56	46	49	AS 1289.3.3.1-2009
Maximum dry unit weight (kN/m <sup>3</sup> )	19	19	19	ASTM D698-12
Optimum moisture content (%)	24	24	24	ASTM D698-12
Gravel content (%)	0	0	0	AS 1289.3.6.1-2009
Sand content (%)	6.7	5	2	AS 1289.3.6.1-2009
Clay content (%)	8.0	10.0	6.0	AS 1289.3.6.1-2009
Soil classification	CH	CH	CH	USCS
Swell Potential (%)	5.0	8.0	7.0	ASTM D4546-21

The soil index properties obtained from laboratory tests were used as input parameters for each model, and the predicted swell potential value was compared with the experimentally measured swell potential values. The results of this comparison are shown in the following section.

#### 4 RESULTS AND DISCUSSION

Figure 3 to Figure 6 show the results obtained from each prediction model for the training and testing datasets. All models demonstrated satisfactory performance on training data; however, the results from the testing data showed the ability to generalize varied for each model. The performance with testing data, in terms of R<sup>2</sup> value, varied from the lowest 0.8 (for SVR and RF) to the highest 0.91 (for ANN). ANN's better

performance can be attributed to its ability to model complicated non-linear patterns between input variables and the target variable, with the use of many neurons and activation functions in the hidden layers.

To validate the models further, predictions were compared with swell potential results obtained from laboratory experiments for three clay samples obtained from the field. These soil data were not used in the training or testing of any model, so this acts as an independent dataset. Table 3 shows the actual swell values and the predicted values, while Table 4 shows the average MSE and R<sup>2</sup> values for each model.

Of the four models, ANN showed a significantly better performance with an R<sup>2</sup> of 0.91 and a lower MSE of 11.1. However, evaluating model performance based solely on statistical metrics like R<sup>2</sup> and MSE is not sufficient. It is essential to critically examine the prediction behavior, data ranges, and potential correlations through domain knowledge. For instance, ANN underpredicts Sample 1.0m considerably, suggesting the model did not give clay content proper weight, possibly due to the lack of representative data. The histogram, Figure 1(c), shows that the majority of the CC values are higher than 20%; therefore, the lower values are not represented adequately.

Predicted values of SVR are close together, suggesting a lack of sensitivity to variation in features like PI or clay content. The outputs of RF and GBM are high and clustered, indicating the model overfits to high swell examples in training or is biased by outliers. RF and GBM are models that split data based on decision rules, where only one variable is evaluated at a time.

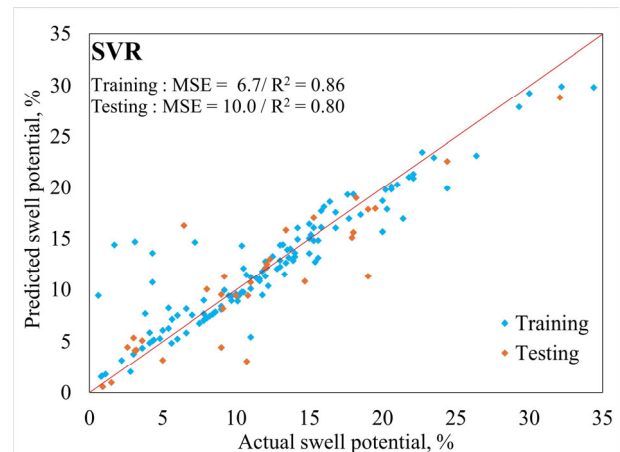


Figure 3. SVR model performance for training and testing datasets.

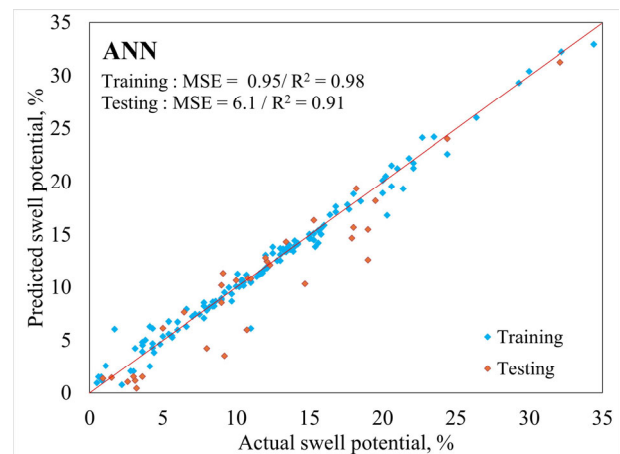


Figure 4. ANN model performance for training and testing datasets.

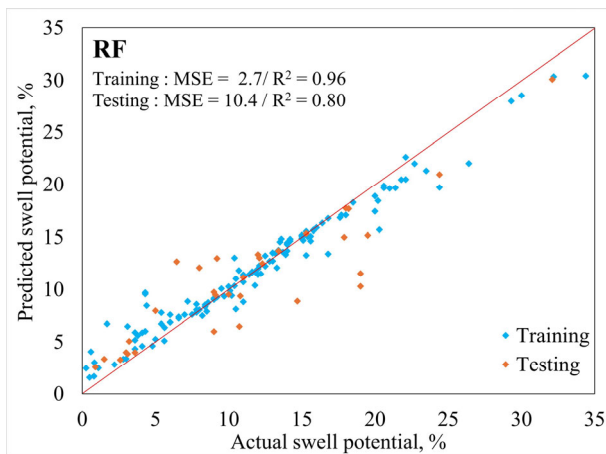


Figure 5. RF model performance for training and testing datasets.

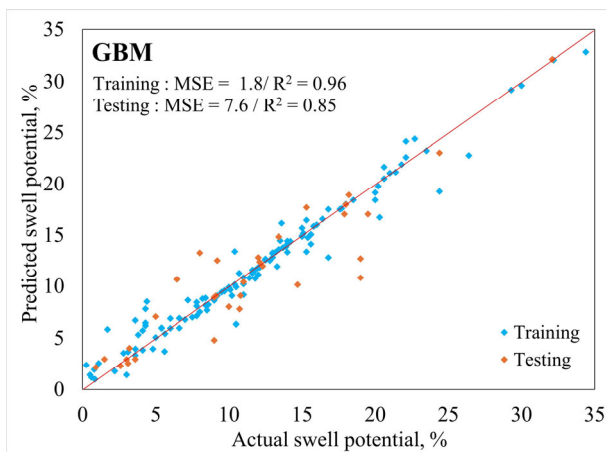


Figure 6. GBM model performance for training and testing datasets.

SVR uses kernels to transform data into a higher space and tries to fit the data with the maximum margin. With these models, there is a chance to miss multiple nonlinear relationships that exist between the predictor variables. Unlike other models, ANN can combine and learn from various variables at once and has the capability of capturing non-linear relationships between variables.

Table 3. Comparison between experimental swell potential values and model predictions for each soil sample

Soil Sample	Experimental SP (%)	SVR	ANN	RF	GBM
0.9 m	5.0	12.6	8.3	17.6	17.3
1.0 m	8.0	11.8	3.3	15.8	13.7
1.2 m	7.0	8.2	6.3	17.6	17.5

Table 4. Performance metrics (RMSE and R<sup>2</sup>) for each machine learning model against experimental values.

ML Model	MSE	R <sup>2</sup>
SVR	24.5	0.13
ANN	11.1	0.91
RF	110.7	0.57
GBM	98	0.52

## 5 CONCLUSIONS

In this study, four machine learning models were developed to predict the swell potential of expansive clays using five soil index properties. Of the four models, ANN performed the best

with both the training and testing datasets showing R<sup>2</sup> values of 0.98 and 0.91, respectively. The models were also validated using an independent dataset obtained from three soil samples. SVR, RF and GBM showed adequate performance for the testing dataset but overpredicted the values for the independent dataset. Mainly because low clay content soils were less in number for training dataset. ANN still outperformed the other models even with this setback. These findings highlight the potential of the ANN model in producing reliable predictions of swell potential. Future work on this topic should include mineralogical data as a predictor variable, as it impacts the swelling behavior of a clay soil and expanding the dataset to have more distribution among the variable values can help improve model performance.

## 6 LIMITATIONS

While the current results are limited by data availability, machine learning models have the potential to outperform traditional methods by learning intricate patterns across diverse conditions. To improve the predictive capabilities, additional relevant features such as mineral composition or soil structure can be considered to capture the full complexity of swell behavior. For practicing engineers, as depicted in the current study, machine learning models like ANN can serve as powerful decision-support tools, especially when integrated with user-friendly front-end interfaces and robust, well-validated models operating in the back end. The effectiveness of such tools depends heavily on the quality of input data; therefore, promoting data sharing and establishing standardized protocols for documentation and collaboration across institutions is essential for developing accurate, generalizable models that can be reliably applied in real-world engineering practice.

## 7 ACKNOWLEDGEMENTS

The authors gratefully acknowledge the support of the UOW Partnership Grant Funding Scheme and Queensland Department of Transport and Main Roads (TMR) for supporting this research project. TMR acknowledge the National Asset Centre of Excellence (NACOE) Project P164, Task 4: Stabilization of Expansive Soil: Assessments and Control Measures. We also highly appreciate the assistance of Mr. Sankalpa Fonseka and the technical staff in the Faculty of Engineering and Information Sciences at the University of Wollongong.

## 8 REFERENCES

- Alazigha, D.P., Indraratna, B., Vinod, J.S. and Ezeajugh, L.E. 2016. The swelling behaviour of lignosulfonate-treated expansive soil. *Proceedings of the institution of civil engineers-ground improvement* 169(3) 182-193.
- Ashayeri, I. and Yasrebi, S. 2009. Free-swell and swelling pressure of unsaturated compacted clays; experiments and neural networks modeling. *Geotechnical and Geological Engineering* 27 137-153.
- Awad, M. and Khanna, R. 2015. Efficient learning machines: theories, concepts, and applications for engineers and system designers.
- Benbouras, M. and Petrisor, A.-I. 2021. Prediction of Swelling Index Using Advanced Machine Learning Techniques for Cohesive Soils. *Applied Sciences* 11(2) 536.
- Bhuvaneshwari, S., Robinson, R., and Gandhi, S. (2010). "Micro-fabric and mineralogical studies on the stabilization of an expansive soil using inorganic additives." *International Journal of Geotechnical Engineering*, 4(3), 395-405.
- Breiman, L. 2001. Random forests. *Machine learning* 45 5-32.
- Chen, J. et al. 2021. Machine learning-based digital integration of geotechnical and ultrahigh-frequency geophysical data for

- offshore site characterizations. *Journal of Geotechnical and Geoenvironmental Engineering* 147(12) 04021160.
- Chen, Y. et al. 2022. Predicting uniaxial tensile strength of expansive soil with ensemble learning methods. *Computers and Geotechnics* 150 104904.
- Çimen, Ö., Keskin, S.N. and Yıldırım, H. 2012. Prediction of Swelling Potential and Pressure in Compacted Clay. *Arabian Journal for Science and Engineering* 37(6) 1535-1546.
- Clarke, S.M., Griebisch, J.H. and Simpson, T.W. 2005. Analysis of support vector regression for approximation of complex engineering analyses. *ASME. J. Mech. Des.*(November 2005) 127(126): 1077–1087.
- Debnath, P. and Dey, A.K. 2017. Bearing capacity of geogrid reinforced sand over encased stone columns in soft clay. *Geotextiles and Geomembranes* 45(6) 653-664.
- Elbadry, H. 2017. Simplified reliable prediction method for determining the volume change of expansive soils based on simply physical tests. *HBRC journal* 13(3) 353-360.
- Erguler, Z.A. and Ulusay, R. 2003. A simple test and predictive models for assessing swell potential of Ankara (Turkey) Clay. *Engineering Geology* 67(3-4) 331-352.
- Eyo, E. and Abbey, S. 2022. Multiclass stand-alone and ensemble machine learning algorithms utilised to classify soils based on their physico-chemical characteristics. *Journal of Rock Mechanics and Geotechnical Engineering* 14(2) 603-615.
- Fityus, S., Smith, D. and Allman, M. 2004. Expansive soil test site near Newcastle. *Journal of Geotechnical and Geoenvironmental Engineering* 130(7) 686-695.
- Friedman, J.H. 2001. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics* 29(5) 1189-1232.
- GhavamShirazi, S. and Bilsel, H. 2021. Characterization of volume change and strength behavior of micro-silica and lime-stabilized Cyprus clay. *Acta Geotechnica* 16(3) 827-840.
- Gupta, R., McCartney, J.S., Nogueira, C.d.L. and Zornberg, J.G. 2008. Moisture migration in geogrid reinforced expansive subgrades. *GEOAMERICAS*. 1-10.
- Jain, A., Fandango, A. and Kapoor, A. 2018. *TensorFlow Machine Learning Projects: Build 13 real-world projects with advanced numerical computations using the Python ecosystem*: Packt Publishing Ltd.
- Jalal, F.E. et al. 2021. Predictive modeling of swell-strength of expansive soils using artificial intelligence approaches: ANN, ANFIS and GEP. *Journal of environmental management* 289 112420.
- Kashif Gill, M., Kemblowski, M.W. and McKee, M. 2007. Soil moisture data assimilation using support vector machines and ensemble Kalman filter 1. *JAWRA Journal of the American Water Resources Association* 43(4) 1004-1015.
- Kukreja, H., Bharath, N., Siddesh, C. and Kuldeep, S. 2016. An introduction to artificial neural network. *Int J Adv Res Innov Ideas Educ* 1(5) 27-30.
- Mahdevari, S., Shahriar, K., Yagiz, S. and Akbarpour Shirazi, M. 2014. A support vector regression model for predicting tunnel boring machine penetration rates. *International Journal of Rock Mechanics and Mining Sciences* 72 214-229.
- Mishra, A.K., Dhawan, S. and Rao, S.M. 2008. Analysis of swelling and shrinkage behavior of compacted clays. *Geotechnical and Geological Engineering* 26 289-298.
- Natekin, A. and Knoll, A. 2013. Gradient boosting machines, a tutorial. *Frontiers in Neuroinformatics* Volume 7 - 2013.
- Pal, M. and Deswal, S. 2008. Modeling pile capacity using support vector machines and generalized regression neural network. *Journal of Geotechnical and Geoenvironmental Engineering* 134(7) 1021-1024.
- Phanikumar, B. and Singla, R. 2016. Swell-consolidation characteristics of fibre-reinforced expansive soils. *Soils and foundations* 56(1) 138-143.
- Puppala, A.J., Manosuthikij, T. and Chittoori, B.C. 2013. Swell and shrinkage characterizations of unsaturated expansive clays from Texas. *Engineering Geology* 164 187-194.
- Puppala, A.J. and Musenda, C. 2000. Effects of fiber reinforcement on strength and volume change in expansive soils. *Transportation research record* 1736(1) 134-140.
- Rao, S.M. and Thyagaraj, T. 2003. Lime slurry stabilisation of an expansive soil. *Proceedings of the Institution of Civil Engineers-Geotechnical Engineering* 156(3) 139-146.
- Sabtan, A.A. 2005. Geotechnical properties of expansive clay shale in Tabuk, Saudi Arabia. *Journal of Asian Earth Sciences* 25(5) 747-757.
- Samui, P., Sitharam, T. and Kurup, P.U. 2008. OCR prediction using support vector machine based on piezocone data. *Journal of Geotechnical and Geoenvironmental Engineering* 134(6) 894-898.
- Schölkopf, B. and Smola, A. 2003. Kernel methods and support vector machines. *Encyclopedia of Biostatistics* 8(2) 5328-5335.
- Segal, M.R. 2004. Machine learning benchmarks and random forest regression.
- Shahin, M.A., Jaksa, M.B. and Maier, H.R. 2001. Artificial neural network applications in geotechnical engineering. *Australian geomechanics* 36(1) 49-62.
- She, J. et al. 2020. Experimental study on the engineering properties of expansive soil treated with Al13. *Scientific Reports* 10(1) 13930.
- Shelke, A. and Murty, D. 2010. Reduction of swelling pressure of expansive soils using EPS geofom. Indian geotechnical conference.
- Sivapullaiah, P., Sridharan, A. and Stalin, V. 1996. Swelling behaviour of soil bentonite mixtures. *Canadian Geotechnical Journal* 33(5) 808-814.
- Soltani, A., Taheri, A., Khatibi, M. and Estabragh, A. 2017. Swelling potential of a stabilized expansive soil: a comparative experimental study. *Geotechnical and Geological Engineering* 35 1717-1744.
- Sridharan, A. and Rao, G.V. 1973. Mechanisms controlling volume change of saturated clays and the role of the effective stress concept. *Géotechnique* 23(3) 359-382.
- Sridharan, A. and Gurtug, Y. 2004. Swelling behaviour of compacted fine-grained soils. *Engineering Geology* 72(1-2) 9-18.
- Taherdangkoo, R. et al. 2023. A Hydro-mechanical Approach to Model Swelling Tests of Clay-Sulfate Rocks. *Rock Mechanics and Rock Engineering* 56(8) 5513-5524.
- Teodosio, B. et al. 2023. Shrink–swell index prediction through deep learning. *Neural Computing and Applications* 35(6) 4569-4586.
- Touzani, S., Granderson, J. and Fernandes, S. 2018. Gradient boosting machine for modeling the energy consumption of commercial buildings. *Energy and Buildings* 158 1533-1543.
- Turkoz, M. and Vural, P. 2013. The effects of cement and natural zeolite additives on problematic clay soils. *Science and Engineering of Composite Materials* 20(4) 395-405.
- Vanapalli, S. and Lu, L. 2012. A state-of-the art review of 1-D heave prediction methods for expansive soils. *International Journal of Geotechnical Engineering* 6(1) 15-41.
- Vieira, S., Pinaya, W. and Mechelli, A. 2017. Using deep learning to investigate the neuroimaging correlates of psychiatric and neurological disorders: Methods and applications. *Neuroscience & Biobehavioral Reviews* 74.
- Yilmaz, I. 2006. Indirect estimation of the swelling percent and a new classification of soils depending on liquid limit and cation exchange capacity. *Engineering Geology* 85(3-4) 295-301.
- Zamin, B. et al. 2021. An experimental study on the geotechnical, mineralogical, and swelling behavior of KPK expansive soils. *Advances in civil engineering* 2021(1) 8493091.
- Zhang, T. et al. 2021a. Improving Convection Trigger Functions in Deep Convective Parameterization Schemes Using Machine Learning. *Journal of Advances in Modeling Earth Systems* 13.
- Zhang, W. et al. 2021b. Prediction of undrained shear strength using extreme gradient boosting and random forest based on Bayesian optimization. *Geoscience Frontiers* 12(1) 469-477.
- Zupan, J. 1994. Introduction to artificial neural network (ANN) methods: what they are and how to use them. *Acta Chimica Slovenica* 41(3) 327.