

Predicting Soil Behavior Types Along the Danube: An AI-Driven Approach Using CPT Data in the Szigetköz Floodplain Area

Edina Koch^{1#}, Shaymaa Alsamia^{1,2}

¹Széchenyi István University, Department of Structural and Geotechnical Engineering, Egyetem square 1. Győr 9026 Hungary

² University of Kufa, Department of Structures and Water Resources, Faculty of Engineering, Iraq

[#]Corresponding author: koche@sze.hu

ABSTRACT

The Szigetköz (Hungary) is a hotbed of sand boil formation, owing to the combination of a 100-250 m thick gravel layer beneath a relatively thin covering of poor soil with varying thickness. Soil behavior is critical for flood protection in this region. This work proposes a novel way to predict Soil Behaviour Types (SBT) based on detailed CPT data collected from 29 sites in the Szigetköz area using an artificial intelligence (AI) model. The study follows a methodically planned approach that includes data collecting, preprocessing, SBT categorization based on the SBT chart developed by Robertson et al. (1986), and AI model building. The CPT dataset contains critical metrics like cone resistance and friction ratio, which are essential in characterising soil behavior. The AI model, built with powerful machine learning algorithms, is intended to learn complicated associations within data to forecast SBT classifications. Extensive feature selection, hyperparameter tuning, and cross-validation are all necessary steps in model construction to ensure accuracy and generalizability. The results show that the model can accurately forecast SBT classifications for the Szigetköz area, shedding information on the soil's behavior near the Danube River. Spatial distribution visualizations emphasize the region's many SBT categories, giving valuable information for engineering projects, land use planning, and environmental conservation activities. The AI model's interpretability elucidates the major CPT parameters driving SBT forecasts, providing stakeholders with actionable information for decision-making. Furthermore, validation of the model with new, previously unseen CPT data confirms its applicability and robustness in real-world circumstances.

Keywords: Soil Behavior Types; soil classification; cone penetration test; machine learning; soil analysis.

1. Introduction

Investigation for accurate soil behavior is essential for engineering projects where the stability of buildings and infrastructure is crucial, mainly in floodplain areas (Liu et al. 2021). The Szigetköz floodplain and the Danube River offer a typical geological environment with various soil characteristics. Technology development has enabled innovative methods to predict and understand the behavior of soils in this area. The cone penetration test (CPT) is a commonly utilized in-situ technique to assess soil characteristics. This method provides significant benefits compared to conventional approaches for conducting fieldwork site investigations, such as excavation and sampling. It has the advantage of generating a continuous data record that exhibits exceptional repeatability and precision, all at a reasonably affordable cost (Miller et al. 2018). CPT was initially implemented in the Netherlands during the 1930s as a mechanical testing method. Subsequently, in the 1960s, it underwent advancements by integrating electric strain gauge load cells into its design. The contemporary CPTu system comprises a digital cone, sometimes called a piezocone, due to its ability to

measure pore pressures (Grabar et al. 2022, Qiao et al. 2023). This study mainly uses Cone Penetration Test (CPT) data to investigate how artificial intelligence (AI) can be included in geotechnical analysis. We aim to use AI to create predictive models that recognize and classify different types of soil behavior in the Szigetköz floodplain. The proposed AI-driven method, which is Naïve Bayes classifier, promises to improve our understanding of the local soil dynamics and offer insightful information for engineering and construction applications.

1.1. Soil classification

Soil classification is essential in geotechnical engineering, especially when assessing site response to earthquakes (Prasad, 2011). A correct soil classification helps to know the seismic effects on soil behavior and understand the dynamic properties of the soil (Chala and Ray, 2023). According to the pioneering work of Begemann (1965), early research endeavored to forecast the distribution of soil particles utilising CPT readings (Libric et al. 2017). Using artificial intelligence to classify soil types based on CPT data offers significant benefits. AI enables accurate and efficient analysis of

complex, nonlinear data patterns. It combines multiple CPT parameters for a thorough classification. Models of artificial intelligence continuously learn and develop, adapting to new datasets. They excel at processing large volumes of CPT data, making them an effective tool for geotechnical engineering and making informed decisions (Laksa and Liu, 2021, Rauter and Tschuchnigg, 2021, Wu et al. 202). By studying Cone Penetration Test (CPT) data and its correlation with known classifications of soil behavior types, engineers can make well-informed judgments regarding the appropriateness of a site for building, the design of foundations and retaining structures, and the evaluation of geohazards (Dagger et al. 2018).

1.2. Cone Penetration Test (CPT)

The Cone Penetration Test (CPT) is the most widely utilized in-situ soil test worldwide due to its reliable results and minimal site disruption, making it less damaging when compared to boreholes (Berthet, 2019). The Cone Penetration Test (CPT) is a geotechnical investigation method that entails the insertion of a narrow, cone-shaped probe into the ground at a consistent velocity, often facilitated by a hydraulic push system. Throughout this procedure, many sensors situated on the cone consistently gather data, including the cone's resistance to penetration and the pore water pressure (Eslami et al. 2019, Talalay 2023). CPT offers significant contributions to the understanding of geotechnical characteristics of soil, encompassing aspects such as stratigraphy, strength, compressibility, and hydraulic conductivity (Rey and Elbatran, 2020). These data from CPT are essential in several applications, encompassing foundation design, slope stability analysis, liquefaction potential evaluation, and characterization of contaminated sites (Moayed et al. 2020). The study of CPT data can yield a significant interpretation of the characteristics and composition of soil and its geological layering. For instance, a sudden elevation in cone resistance may indicate the existence of compact or hard soil strata, while a decline in cone resistance may imply the presence of soft or loose soils (Bol, 2023).

2. Study Area

The Szigetköz region is located in northwestern Hungary at the confluence of the Danube and Mosoni Danube rivers. Its average width is 7 km, and its length is approximately 50 km. The flooding region is substantially narrower, only reaching more enormous torrents longer than three kilometres Hahn et al. 2011).

The soils found in the Szigetköz region have originated from predominantly alluvial soils. The formations in question can be classified as azonal due to the impeded growth caused by frequent floods, which hindered organic matter accumulation. Sandy or muddy textures predominantly characterize the soil composition. Humic alluvial soil types and their mixtures can be observed in regions characterized by higher elevations. The repeated floods resulted in the deposition of annual muddy layers measuring up to 2 cm. This consistent process ensured a reliable provision of nutrients for the

trees. The topsoil has a mosaic-like pattern in terms of its depth, typically ranging from 50 to 300 cm. Beneath the superficial layer of fine-textured topsoil, a substantial deposit of gravel is present, reaching considerable depths of several hundred meters. This gravel layer may occasionally be intermingled with coarse sand (Somogyi et al. 1999, Guti, 2020).

3. Methodology

Our research presents a novel approach to soil behavior type (SBT) classification using cone penetration test (CPT) data, employing a custom algorithm developed in Python. This algorithm operates on the principle of spatial proximity in a feature space, where each CPT observation is mapped onto a specific class of SBT.

- The process begins by generating huge random points on the chart which was suggested by Robertson et al. (1986). The structure of the chart consists of nine distinct zones, each zone is a specific soil type. Accordingly, nine predefined clusters of labelled datasets are created to be used for the subsequent supervised learning process. These clusters will be the benchmark for the classification algorithm for the soil type.
- Since the clusters are generated in the logarithmic scale proposed by the Robertson chart, the process begins with a set of calibration points that relate logarithmic coordinates to actual physical measurements. To accommodate the non-linear nature of the data, we employ logarithmic transformation for both the cone resistance (q_c) and the friction ratio (R_f), which allows for an interpolation that respects the exponential scaling of the soil properties.
- Data preprocessing is necessary to fill in the missing data due to some technical reasons in certain features like depth or cone penetration resistance. Missing data within a certain feature is straightforward and can be estimated using curve-fitting techniques.
- In the context of machine learning, a system refers to the entire process or model that takes inputs (features) and provides outputs (predictions or classifications). The data itself is not a system; it is a dataset. However, when we use this data within a machine learning framework to predict soil behavior types (SBT), the combination of the data, the algorithm, and the model's predictive capabilities constitute a system.
- To constitute a system, we have created an empty output field called SBT which should consist of nine classes numbered from 1 to 9. Each class represent soil type. In other words, each CPT observation must have a unique indexed class corresponding to a soil behavior.
- A custom algorithm is created that uses the features to classify CPT data observations into predefined SBT clusters. This step involves defining a function that computes the Euclidean distance from the CPT data point in question and all of the points in the SBT cluster. The

classification is occurred by assigning the CPT data observation to the nearest cluster based on the shortest distance.

This custom algorithm automates the classification of CPT data into distinct SBTs, enabling practitioners to quickly and accurately determine the soil profile. Through this automated process, extensive and labor-intensive manual classification is circumvented, streamlining the assessment of subsurface conditions for geotechnical applications.

4. Predefining SBT Clusters

A non-normalized Cone Penetration Test (CPT) Soil Behavior Type (SBT) chart is a graphical representation used in geotechnical engineering to classify soil types based on direct measurements from CPT data without any normalization or correction for overburden pressure, see Fig. 1 and Table 1. It typically plots cone resistance (q_c) against sleeve friction (f_s), with different zones delineated on the chart that corresponds to various soil behaviors. Without normalization, the chart uses raw CPT data, which means that the results can be influenced by site-specific factors like overburden pressure and may not be directly comparable across different sites or depths. However, non-normalized charts can still provide quick and useful insights for on-site assessments and initial soil classification.

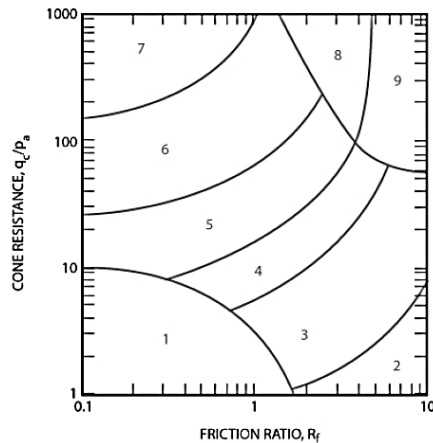


Figure 1. Non-normalized CPT Soil Behavior Type (SBT) chart (Robertson 2015, Collico et al. 2022)

Table 1. Non-normalized CPT Soil Behavior Type (SBT) (Robertson, 2015)

Zone	Soil Behavior Type
1	Sensitive, fine grained
2	Organic soils - clay
3	Clay - silty clay to clay
4	Silt mixtures - clayey silt to silty clay
5	Sand mixtures - silty sand to sandy silt
6	Sands - clean sand to silty sand
7	Gravelly sand to dense sand
8	Very stiff sand to clayey sand*
9	Very stiff fine grained*

Fig. 2 reveals the predefined master clusters that will be the references for the classification of CPT data. The process was described in the methodology section step 1 where we have nine SBT cluster each one has huge number of data points and each SBT cluster correspond to a specific soil type. One thing to know about these clusters is that the centre of the cluster is not important to be filled with random points since this work uses the strategy of the nearest point to the input instead of the nearest centroid.

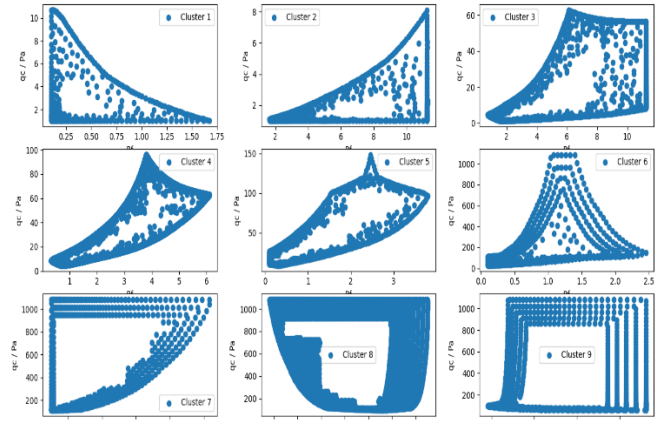


Figure 2. Predefined SBT clusters

5. Classification

Fig. 3 illustrates the soil profile in one of the 29 distinct zones along the Danube river while Fig. 4 represents the distribution of the soil behavior type up to about 26 m depth. Most of the composition in the profile is sand – clean sand to silty sand while the surface layer is sensitive, fine-grained soil with little percentage of sand mixtures, gravelly sand, and very stiff sand to clayey sand.

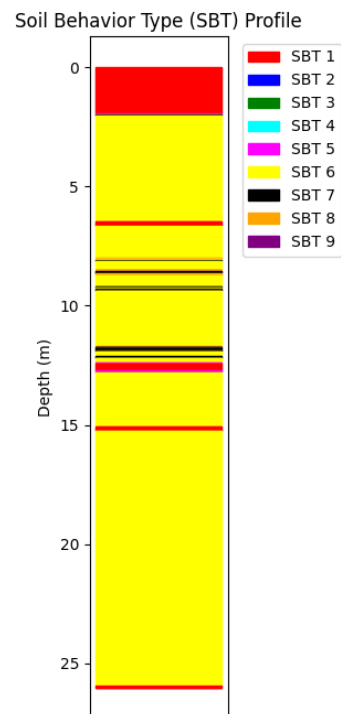


Figure 3. Soil profile in certain locations along the Danube

Also, Fig. 4 reveals that Sensitive, fine-grained and Sands – clean sand to silty sand are distributed along wide ranges of depth. For depths of about 3-14 m a mixture of other soils can be found; sand mixtures, gravelly sand, and very stiff sand to clayey sand.

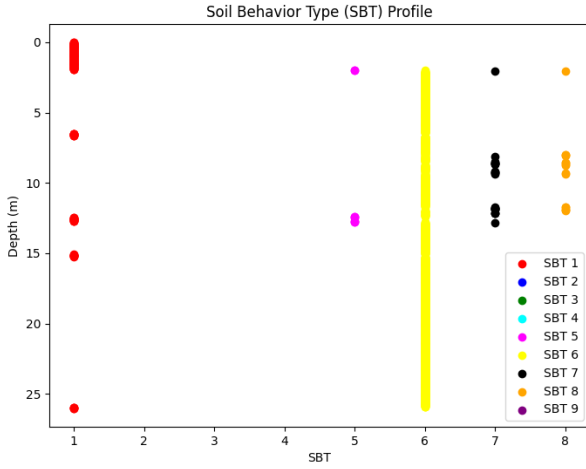


Figure 4. Classified SBT based on CPT data

6. Naive Bayes classification

The Naive Bayes classifier is a probabilistic machine learning model that is used for classification tasks. It is based on Bayes' Theorem, which uses probability to predict the class of an unknown data point. The "naive" assumption in Naive Bayes is that all features are independent of each other, which simplifies the computation, and while this assumption may not hold true in all cases, it often yields surprisingly accurate results. Here's how the Naive Bayes classifier works on the CPT dataset with 5 features and an output of 9 SBT classes.

6.1. Features and Output Description

CPT features (Inputs) are the following:

- depth: Numerical value indicating the depth at which CPT measurements are taken.
- q_c : Cone resistance, reflecting the hardness or density of the soil.
- friction sleeve: The resistance measured against the friction sleeve, related to soil texture.
- pore pressure: The hydrostatic pressure within the soil pores, which can influence soil behavior.
- R_f : Friction ratio, a derived feature from q_c and friction sleeve measurements.

The output (Target) is:

- SBT: Soil Behavior Type, categorized into 9 different classes based on the soil properties inferred from the CPT data.

6.2. Naive Bayes Classifier Mechanics

Regarding the probabilistic model building, for each class of SBT, the classifier will estimate the likelihood of observing each feature value. In the Gaussian Naive Bayes variant, this is done by assuming that the continuous features for each class are distributed according to a Gaussian distribution. It calculates the

mean and variance of each feature for each class. On the other hand, class probability, It computes the prior probability for each SBT class based on the frequency of each class in the training data. This gives the model a baseline to work from before considering the evidence from the features. For a new data point, the likelihood is calculated for each feature within each class. For continuous data like in the CPT dataset, this is typically done by plugging the feature values into the probability density function of the Gaussian (normal) distribution, which has been characterized for each feature in each class. Using Bayes' Theorem, the posterior probability is calculated as Eq. (1):

$$P(\text{Class} | \text{Data}) = \frac{P(\text{Data}|\text{Class}) \times P(\text{Class})}{P(\text{Data})} \quad (1)$$

In a Naive Bayes classifier, this calculation is performed for each class, and the class with the highest posterior probability is typically chosen as the prediction. The following is the twist with Naive Bayes:

- Input Features: Depth, Cone resistance (q_c), Friction sleeve, Pore pressure, Friction ratio (R_f)
- Output Classes: Soil Behavior Types (SBT), which are categorized into 9 different classes based on soil properties inferred from CPT data.

The prior probability of each class (i.e., each SBT), denoted as $P(\text{Class}_k)$, is the initial assumption about the probability of a soil behaviour type before observing any data. It is typically estimated from the frequency of each class in the training dataset (Eq. (2):

$$P(\text{Class}_k) = \frac{\text{Number of instances of Class}_k}{\text{Total number of instances}} \quad (2)$$

The likelihood $P(\text{Data}|\text{Class}_k)$ is the probability of observing the data given a particular class. Under the Naive Bayes assumption, the joint probability of the features given the class is the product of the individual probabilities for each feature, assuming they are conditionally independent (Eq. (3).

$$P(\text{Data} | \text{Class}_k) = P(\text{depth}/\text{Class}_k) \times P(q_c/\text{Class}_k) \times P(\text{friction sleeve}/\text{Class}_k) \times P(\text{pore pressure}/\text{Class}_k) \times P(R_f/\text{Class}_k) \quad (3)$$

For continuous features such as depth, q_c , friction sleeve, pore pressure, and R_f , the probability density is often modeled using a Gaussian distribution as Eq. (4).

$$P(x_i / \text{Class}_k) = \frac{1}{\sqrt{2\pi\sigma_{i,k}^2}} \exp\left(-\frac{(x_i - \mu_{i,k})^2}{2\sigma_{i,k}^2}\right) \quad (4)$$

where x_i is a feature value, $\mu_{i,k}$ and $\sigma_{i,k}$ are the mean and standard deviation of feature i for class k , respectively. Using Bayes' theorem, the posterior probability $P(\text{Class}_k|\text{Data})$ is calculated by updating the prior with the likelihood of the observed data as Eq. (5).

$$P(\text{Class}_k/\text{Data}) = \frac{P(\frac{\text{Data}}{\text{Class}_k}) * P(\text{Class}_k)}{P(\text{Data})} \quad (5)$$

To simplify computation, you can omit the denominator $P(\text{Data})P(\text{Data})$ when comparing which class has the highest posterior probability since it remains constant for all classes as Eq. (6).

$$P(\text{Class}_k|\text{Data}) \propto P(\text{Data}|\text{Class}_k) \times P(\text{Class}_k) \quad (6)$$

Using Eq. (7), the predicted class for a new observation is the class that maximizes the posterior probability.

$$\hat{y} = \arg \max_k P(\text{Class}k | \text{Data}) \quad (7)$$

This representation provides a clear mathematical structure of how the Naïve Bayes classifier is applied to the dataset to predict soil behaviour types. The following procedure can be implemented when classifying a new data point with features depth, qc, friction sleeve, pore water pressure, and R_f . For each SBT class, consider:

- Take the prior probability of the SBT class.
- Multiply it by the likelihood of observing each feature value given that SBT class. For continuous features, you'd plug the feature value into the Gaussian distribution for that feature under the SBT class to get the likelihood.
- Compare the resulting products for each class. The class with the highest product is the one with the highest posterior probability, and hence, is the predicted class for that data point.

The actual probability value of the posterior, if needed, would require calculating or estimating the evidence term and using it to normalize the product of the likelihood and the prior. In summary, here how Naïve Bayes classifier works on the CPT dataset:

- **Feature Independence:** Naïve Bayes simplifies the complexity by assuming that each feature makes an independent and equal contribution to the outcome. For your dataset, this means that each soil measurement is considered independent from the others when calculating the probability of a particular SBT class.
- **Probability Calculation:** For each SBT class, the classifier will calculate the probability that a given data point belongs to that class, based on the features. It does this by looking at the distribution of each feature within each class in the training dataset.
- **Handling Continuous Data:** Since the features in your dataset are continuous, the Gaussian Naïve Bayes variant is used. It assumes that the continuous values associated with each class are distributed according to a Gaussian (normal) distribution.
- **Training the Model:** During training, the algorithm calculates the mean and variance of each feature for each class label. These parameters define the shape of the Gaussian distribution for each class.
- **Class Prediction:** When making predictions, the classifier uses these Gaussian distributions to estimate the probability of the new data point belonging to each SBT class. It does this by plugging the feature values of the new data point into the Gaussian distributions for each class.
- **Probabilistic Output:** The output is a probabilistic statement about the likelihood of the new data point belonging to each of the 9 SBT classes. The classifier picks the class with the highest probability as the prediction.

- **Model Evaluation:** After the model is trained, it's important to evaluate its performance using appropriate metrics. With a multiclass classification problem like this one, you might use a confusion matrix, accuracy, precision, recall, and F1-score.

The process when applied to CPT dataset:

- **Preprocessing:** The dataset is first cleaned and preprocessed. Continuous features may be scaled or normalized if required, although Gaussian Naïve Bayes is less sensitive to this due to its reliance on probabilities rather than distances.
- **Feature Analysis:** For each SBT class, the model computes the summary statistics (mean and variance) for each feature within that class.
- **Probability Estimation:** For a new data point, the model calculates the conditional probability of it belonging to each SBT class based on the observed features.
- **Decision Rule:** The classifier then applies Bayes' Theorem to update the probabilities based on the evidence provided by the new data point's features.
- **SBT Prediction:** The class with the highest posterior probability is chosen as the predicted class for the new data point.

6.3. Predictions

The model was trained on 1304 observations for a CPT dataset. The trained Naïve Bayes classifier has performed excellent accuracy: 0.923 while its confusion matrix was (Eq. (8)).

$$CM = \begin{bmatrix} 29 & 1 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 326 & 6 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 1 \end{bmatrix} \quad (8)$$

The confusion provided matrix is a performance measurement for machine learning classification. It shows the number of correct and incorrect predictions made by the classification model compared to the actual outcomes (ground truth) in the data. Each row of the matrix represents the instances in an actual class while each column represents the instances in a predicted class. Here's how to interpret the confusion matrix:

- The matrix is 5x5, which suggests that we have 5 classes in the dataset.
- The diagonal elements represent the number of points for which the predicted label is equal to the true label.
- The off-diagonal elements are those that were labeled incorrectly by the classifier.
- Rows represent the actual classes, while columns represent the predicted classes.

Class 1: 29 were correctly predicted as class 1 (true positives for class 1), 1 was incorrectly predicted as class 5, and 10 were incorrectly predicted as class 6. Class 5: The model did not make any predictions for class 5, the model did not recognize this class from the features provided. Class 6: 1 was incorrectly predicted as class 1, 1 incorrectly as class 5, 326 correctly as class 6 (true

monitoring”, *Ekologia(Bratislava)/Ecology(Bratislava)*, vol. 18, pp. 59–68, 1999.

Talalay, P. G. “Non-Rotational Drilling and Sampling in Frozen Soils”, in *Geotechnical and Exploration Drilling in the Polar Regions*, Springer, 2023, pp. 105–137. https://doi.org/10.1007/978-3-031-07269-7_4

Wu, S., Zhang, J. M. and Wang, R. “Machine learning method for CPTu based 3D stratification of New Zealand geotechnical database sites”, *Adv. Eng. Informatics*, vol. 50, p. 101397, 2021. <https://doi.org/10.1016/j.aei.2021.101397>