



UNIVERSITY OF LEEDS

This is a repository copy of *Geotechnical correlation field-informed and data-driven prediction of spatially varying geotechnical properties*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/218604/>

Version: Accepted Version

Article:

Chen, W., Ding, J., Shi, C. et al. (2 more authors) (2024) Geotechnical correlation field-informed and data-driven prediction of spatially varying geotechnical properties. *Computers and Geotechnics*, 171. 106407. ISSN 0266-352X

<https://doi.org/10.1016/j.compgeo.2024.106407>

© 2024, Elsevier. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>. This is an author produced version of an article published in *Computers and Geotechnics*. Uploaded in accordance with the publisher's self-archiving policy.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Geotechnical Correlation Field-informed and Data-Driven Prediction of Spatially Varying Geotechnical Properties

Weihang Chen ^a, Jianwen Ding ^{a,*}, Chao Shi ^b, Tengfei Wang ^{c,d}, David P. Connolly ^e

^a School of Transportation, Southeast University, Nanjing 210096, China

^b School of Civil and Environmental Engineering, Nanyang Technological University, Singapore
639798, Singapore

^c School of Civil Engineering, Southwest Jiaotong University, Chengdu 610031, China

^d MOE Key Laboratory of High-Speed Railway Engineering, Southwest Jiaotong University,
Chengdu 610031, China

^e School of Civil Engineering, University of Leeds, Leeds LS2 9JT, UK

* Corresponding author

Email: jwding@seu.edu.cn

1 **Abstract**

2 Geotechnical measurements are often limited, leading to the use of interpolation techniques for
3 interpreting spatial variations in geotechnical properties from sparse geo-data. Traditional
4 geostatistical methods suffer from significant computational complexity. On the other hand, data-
5 driven approaches often lack integration with geotechnical domain knowledge, potentially
6 oversimplifying or complicating predictions related to the spatial variability of geotechnical properties.
7 This study introduces a novel framework that combines geotechnical knowledge with data-driven
8 methods to model inherent soil spatial variability incorporating Geotechnical Correlation Field (GCF)
9 that reflects domain knowledge. The GCF, influenced by Autocorrelation Function (ACF) types and
10 Scale of Fluctuation (SOF), provides a flexible basis for accurately representing spatially varying
11 geotechnical properties. Using a large synthetic database comprising known ACF types and SOFs, we
12 constructed a series of specialized neural networks. These networks identify random field parameters
13 at different sites based on sparse data, and the estimated parameters can be directly used to calculate
14 GCFs for a given site. The performance of the proposed method is validated using a set of synthetic
15 data and a real case history in New Zealand. The results demonstrate the model can accurately predict
16 random field parameters for irregularly spaced geo-data, even with limited information. Significantly,
17 the GCFs offer improved physical interpretations and enhance the performance of subsurface modeling.
18 The computational complexity of this method is independent of the number of soil cells, making it
19 highly efficient and scalable.

20 **Keywords:** Spatial variability; Data-driven Method; Random field theory; Site investigation; Neural
21 Network

22 **1 Introduction**

23 The spatial variability of soil properties is a significant source of uncertainty in geotechnical
24 engineering (Phoon et al., 2022; Shi and Wang, 2023; Uzielli et al., 2005). This challenge is prevalent
25 across various fields, such as accurately assessing the extent and concentration of contaminated sites
26 in environmental engineering or predicting mineral reserves and distribution densities in mining
27 engineering (Gu et al., 2023; Wang and Shi, 2023; Zhang et al., 2020; Zhao et al., 2018). Precise and
28 high-resolution geographic information can assist engineers in analysis and design optimization (Chen
29 et al., 2023; Wang et al., 2020; Zhao et al., 2020). However, obtaining high-density in-situ
30 measurements using expensive and time-consuming testing equipment is often impractical.
31 Engineering projects typically retrieve sparse measurements from limited locations. For instance, it is
32 common to use Cone Penetration Tests (CPT) at intervals of 25-100 m along the ground surface (Guan
33 et al., 2020) to characterize subsurface ground conditions (Collico et al., 2024). Although CPT tests
34 provide continuous measurements with depth, the interpretation of spatial distribution of soil properties
35 is challenging, particularly in the horizontal direction (Xie et al., 2022b, 2022a; W. Zhang et al., 2022).

36 The use of interpolation or data-driven methods to predict soil properties at unsampled locations
37 is an active area of study (Shi and Wang, 2021a; Wang et al., 2018; Xie et al., 2024). In this context,
38 several techniques, including Geostatistical methods, Distance-based methods, Bayesian analysis, and
39 data-driven methods, have found wide application (Hu and Wang, 2024; Wang and Chen, 2023; Zou
40 et al., 2017). For instance, Kriging provides direct estimates and uncertainty assessments of soil
41 properties at unsampled locations, but when geotechnical measurements are sparse and exhibit
42 significant variations, Kriging predictions often tend to capture the average trend and disregard the

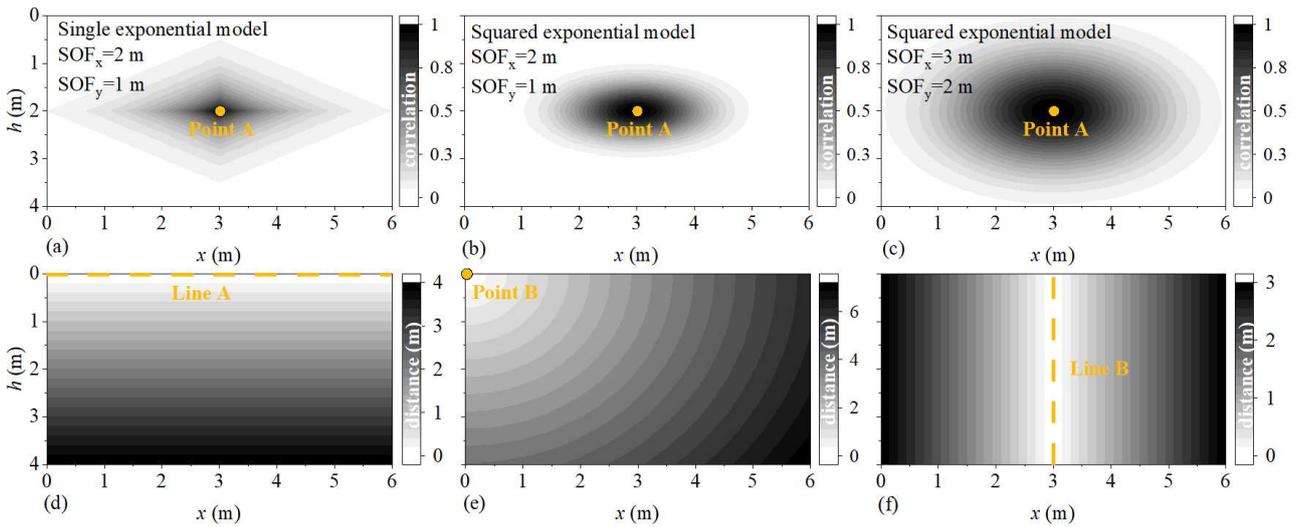
43 spatial variability of soil properties (Nag et al., 2023). Notably, when simulating large-scale or fine-
44 resolution random fields, traditional conditional random field methods may suffer from low
45 computation due to excessively large correlation matrices (Yang et al., 2021). Recently, Yang and
46 Ching (2021) proposed an efficient method for simulating conditional random fields (CRFs) by
47 utilizing the Kronecker-product to decompose the large correlation matrix, and further extended this
48 approach to simulate multivariate cross-correlated CRFs (Z. Yang et al., 2022). Bayesian Compressive
49 Sensing (BCS) is a non-parametric and data-driven method that can be directly applied to non-
50 stationary random fields (Wang and Zhao, 2017). However, BCS does not incorporate specific basis
51 function forms relevant to geotechnical modeling (Cami et al., 2020; Phoon and Wang, 2019). Inverse
52 Distance Weighting (IDW) and the Geotechnical Distance Field (GDF) method both attempt to recover
53 the soil properties at unsampled points based on the "distance" between unsampled points and sampled
54 points (Xie et al., 2022b). It's important to note that soil properties exhibit location-specific
55 dependencies, implying that soil properties within a certain scale of fluctuation (SOF, δ) in a particular
56 subsurface stratigraphy are correlated (Phoon et al., 2003). Therefore, relying solely on relative
57 distances to infer soil properties at unsampled points may overlook this fact.

58 Although data-driven methods can be straightforward to use, they tend to overlook the
59 geotechnical expertise, such as random field theory. This may lead to overly simplified or complex
60 subsurface modeling results. These 'black-box' models may impede effective collaboration between
61 engineers and the models, compromising the prediction performance. The incorporation of
62 geotechnical knowledge can steer the prediction towards correct solutions. Therefore, a random field-
63 informed and data-driven model is introduced to predict spatially varying soil property from sparse

64 site-specific measurements. In this approach, random field theory is embedded in the data-driven
65 model through Geotechnical Correlation Field (GCF). Essentially, GCF is derived from the
66 decomposition of the correlation matrix (C) used in random field theory to describe the correlation ρ
67 between soil cells. For instance, Fig. 1 (a)-(c) illustrate three GCFs for point A, where darker colors
68 correspond to a higher correlation between the respective soil cell and point A. The advantages of
69 employing GCFs are threefold: (1) Soil properties at point A are only correlated with those within a
70 specific range, strictly adhering to the fact that soil properties exhibit location-specific dependencies;
71 (2) The correlation between point A and other points is affected by the type of autocorrelation function
72 (ACF)(Ching et al., 2019), providing a flexible tool for accurately describing soil properties; (3) The
73 impact of point A on other points is constrained by the Scale of Fluctuation (SOF), facilitating the
74 geotechnical engineer to enhance further collaboration with the model by controlling the SOF. In
75 contrast, Fig. 1(d)-(f) represent three Geotechnical Distance Fields (GDFs) corresponding to the
76 distances from the sampling point to the ground surface (Line A), the corner point (Point B), and the
77 exploration location (Line B). Compared to methods using two-dimensional coordinates as input
78 features, GDFs significantly enhance feature dimensions and modeling accuracy. However, GDFs
79 overlook the fact that soil properties exhibit location-specific dependencies. Furthermore, in
80 comparison to BCS and GDFs, GCFs can comprehensively consider multiple ACF types, providing
81 flexible input features, as opposed to relying on a single base function or fixed feature type.

82 It is noteworthy that the generation of GCF depends on the specific site's ACF types and
83 corresponding SOFs. Estimating random field parameters based on sparsely distributed measurement
84 data is challenging (Dasaka and Zhang, 2012; Qi and Liu, 2019; Yan et al., 2023). Traditional methods

85 include the method of moments (Lloret-Cabot et al., 2014; Onyejekwe et al., 2016), maximum
 86 likelihood estimation (Xiao et al., 2018) and Bayesian analysis(Cami et al., 2020; Ching et al., 2018).
 87 However, these methods are influenced by human experience and involve certain assumptions about
 88 describing SOFs, mainly used for estimating vertical SOFs. Recently, some machine learning-based
 89 methods for estimating horizontal and vertical SOFs have been proposed, with prediction models based
 90 on Convolutional Neural Networks (CNN) demonstrating good performance and efficiency (Zhang et
 91 al., 2021; 2022). Nevertheless, current CNN methods cannot assess the optimal ACF type based on
 92 measurement data.



93
 94 Fig. 1 Illustration of GCFs and GDFs: GCFs (a-c); GDFs (d-f)

95 To address these challenges, this study draws inspiration from the generation process of traditional
 96 conditional random field to construct a framework for the data-driven model. This framework provides
 97 support for the proposed subsurface modeling method, by allowing the developed data-driven model
 98 to integrate geotechnical knowledge. The framework employs random field samples generated from
 99 specified GCFs serve as training data to train a neural network that is used for forward prediction of
 100 spatially varying soil properties and estimation of random field parameters. The performance of the

101 proposed method is illustrated using a set of synthetic data and a real case study in New Zealand. The
102 remainder of this study is organized as follows: In section two, the proposed subsurface modeling
103 framework is introduced. The construction and validation of the random field parameter estimation
104 model is in section three. Subsequently, the performance of the proposed GCF-based subsurface
105 modeling approach is validated by a set of synthetic data. In section five, further validation is
106 conducted using a real case study in New Zealand, followed by conclusions.

107 **2 Proposed Methods for Predicting Spatially Varying Geotechnical Properties**

108 The essence of data-driven subsurface modeling lies in utilizing machine learning (ML)
109 techniques to learn the correlation between sampling positions (coordinates) and the corresponding
110 soil properties, and then using the trained ML model to infer soil properties at unsampled positions.
111 Notably, this study effectively integrates random field theory into the process of subsurface modeling
112 using GCF, aiming to enhance modeling reliability and reduce modeling uncertainty.

113 In GCF, instead of using 2D or 3D coordinates to represent the positions of soil cells, it employs
114 the correlation between each soil cell and every sampled soil cell. Therefore, the framework proposed
115 in this study first estimates the random field parameters of the site (e.g., types of ACFs, horizontal SOF,
116 and vertical SOF). Subsequently, based on random field theory and predicted parameters,
117 corresponding GCFs are computed. Finally, a machine learning (ML) model is employed to learn the
118 complex relationship between GCFs and observed soil properties, enabling predictions of soil
119 properties at unsampled locations. As illustrated in [Fig. 2](#), the geotechnical subsurface modeling
120 process in this study involves six key steps, using the case of three CPTs as an illustrative example:

121 (1) Collect CPT data (e.g., cone tip resistance q_c) along with corresponding horizontal coordinates

122 (L).

123 (2) Utilize the ACFs type prediction model (Model #1) to assess the probability P of measured
124 CPT data belonging to each ACF type. This study considers seven common ACF types, ensuring that
125 $\sum P_i=1, i=1-7$.

126 (3) Utilizing the horizontal and vertical SOFs prediction models (Model #2 and Model #3) to
127 accurately estimate the measured CPT data's horizontal and vertical SOFs. Notably, SOFs predicted
128 based on different ACF types can vary. Therefore, both Model #2 and Model #3 incorporate seven sub-
129 models (corresponding to the seven ACF types), resulting in seven sets of predicted horizontal and
130 vertical SOFs corresponding to different ACF types.

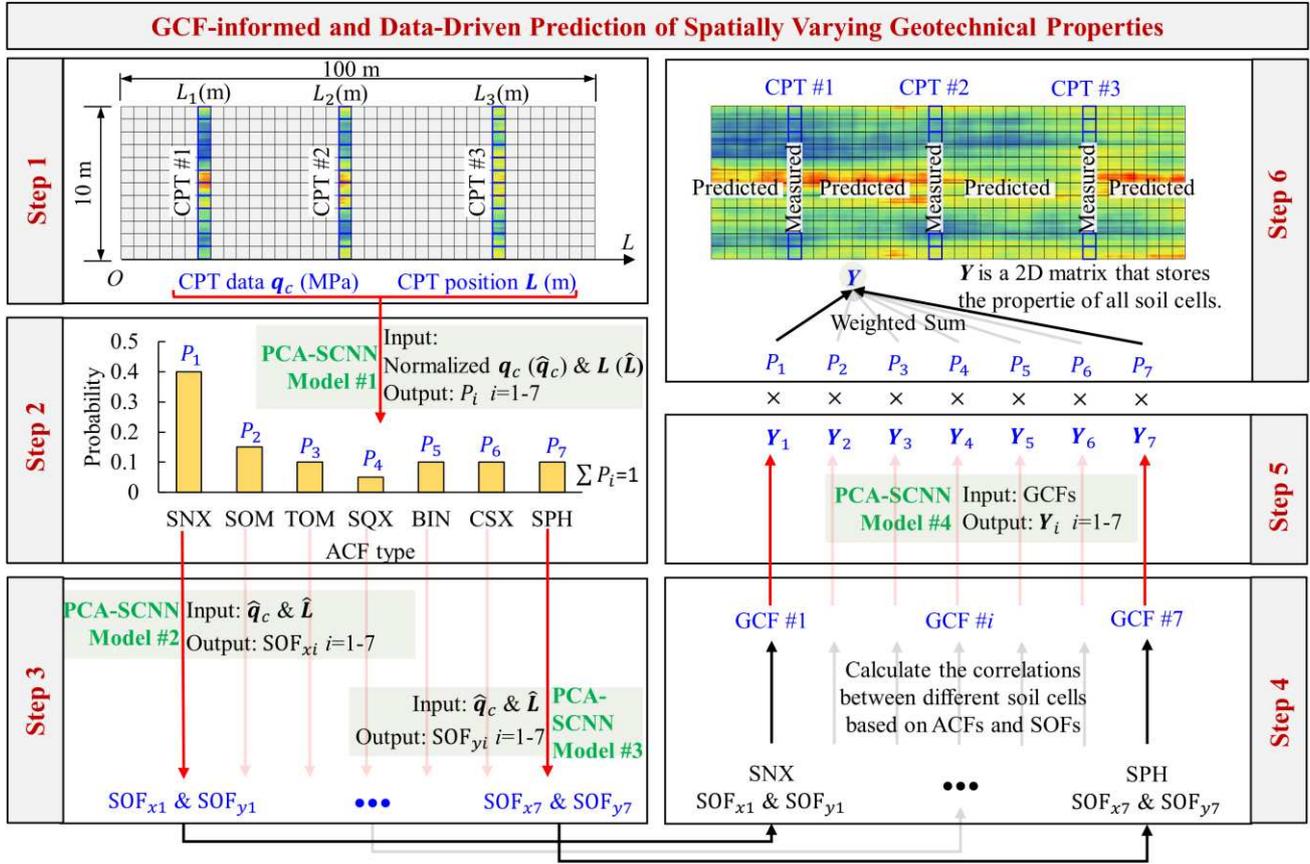
131 (4) Employing seven sets of random field parameters (ACFs and corresponding horizontal and
132 vertical SOFs) to calculate the seven sets of GCFs, as outlined in Section 2.3.

133 (5) Constructing the subsurface prediction model (Model #4) involves using GCFs as input and
134 the geotechnical properties of soil cells as output. Once Model #4 is well-trained, it can be utilized to
135 predict soil properties at unexplored locations. Notably, since Step (4) generates seven sets of GCFs,
136 it allows for the creation of seven sub-models within Model #4.

137 (6) Weighted summation of the predicted results (Y_i) from the seven sub-Model #4, based on the
138 probabilities (P_i), yields the subsurface modeling outcome ($Y = \sum P_i Y_i$) closely correlated with
139 observed data.

140 Notably, in Steps (2), (3), and (5), independent data-driven models (PCA-SC neural networks) are
141 employed for predicting ACFs types, estimating horizontal and vertical SOFs, and subsurface
142 modeling. Specifically, the ACFs types and SOFs prediction models are trained using a large synthetic

143 q_c database containing known ACFs types and SOFs. Detailed descriptions of Models #1-4 are
 144 provided in the subsequent sections.



145
 146 Note: \hat{q}_c and \hat{L} represent the standardized cone tip resistance q_c and horizontal coordinates L . Standardization of
 147 data is necessary for neural network models to expedite the training process and enhance model performance. In Step
 148 (2), the seven types of ACFs can be referenced from Table 1. Since the values of GCF range between 0 and 1, larger
 149 GCF values indicate stronger correlation between two soil cells.

150 Fig. 2. Framework for Proposed Soil Property Recovery for 3 CPTs

151 **2.1 PCA–Shortcut Connection Neural Network**

152 It is noteworthy that a universal PCA-SCNN structure, as depicted in Fig. 3, is utilized for ACFs
 153 classification (Model #1), horizontal SOF estimation (Model #2), vertical SOF estimation (Model #3),
 154 and subsurface modeling (Model #4). While the training datasets fed into the PCA-SCNN differ across
 155 tasks, resulting in different inputs and outputs for each model, the remaining structure of the models is
 156 consistent.

157 In Fig. 2, Models #1-3 take standardized q_c and L as input, while Model #4 utilizes GCFs. It's
158 worth noting that CPT tests provide almost continuous soil information in the vertical direction,
159 leading to higher-dimensional q_c and GCFs. However, handling high-dimensional input features
160 increases model complexity, extends training time, and requires substantial memory. This study
161 proposes a solution: the PCA-NN structure. PCA, an unsupervised dimensionality reduction method,
162 transforms multiple potentially correlated features into a smaller set of linearly uncorrelated principal
163 components (PCs). This process eliminates redundancy in input data and reduces the number of
164 neurons in the neural network's input layer. By simplifying the data-driven model structure, this
165 enhances model convergence speed (Bai et al., 2023; He et al., 2016). The dimension of PCs is set to
166 100 in this study.

167 As the depth of a neural network increases, its non-linear representation strengthens (Hong et al.,
168 2024; Zhang et al., 2023). But deeper networks introduce challenges like gradient vanishing and
169 exploding. To tackle these issues, shortcut connections are incorporated into the model. Research by
170 Li et al. (2018) indicates that these connections promote loss surface minimization and prevent chaotic
171 behavior. Therefore, a shortcut connection is established between the input layer and the last hidden
172 layer, concatenating input features with those of the final hidden layer to improve training efficiency
173 and overall performance.

174 Fig. 3 displays the model structure. Notably, Monte-Carlo dropout (MC dropout) is applied after
175 the hidden layers. Dropout randomly deactivates a portion of neuron connections during training,
176 preventing the model from becoming too dependent upon specific neurons and thus reducing
177 overfitting (P. Zhang et al., 2022). MC dropout extends this concept by randomly deactivating neuron

178 connections not only during training but also during prediction. While the network structure and weight
179 parameters are fixed after training, MC dropout introduces randomness to the model structure. Through
180 multiple predictions, an output distribution is obtained, assisting in the assessment of prediction
181 uncertainty. The model formulation is as follows:

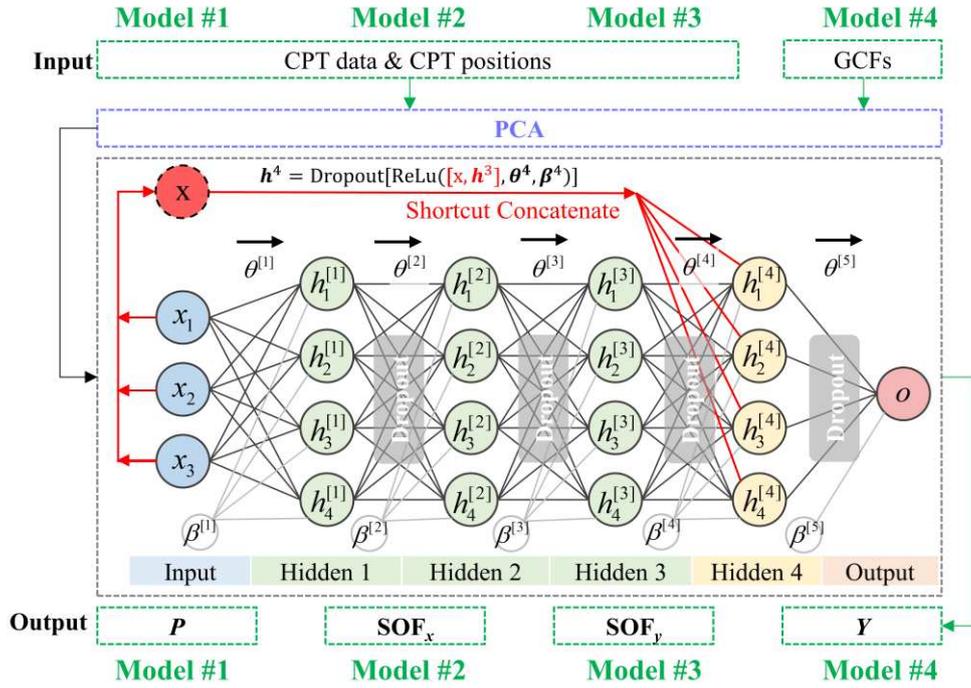
$$182 \quad \mathbf{x} = \text{PCA}(\text{Input}) \quad (1)$$

$$183 \quad \mathbf{h}^i = \text{Dropout}[\text{ReLU}(\mathbf{x}, \boldsymbol{\theta}^i, \boldsymbol{\beta}^i)] \quad i = 1, 2, 3 \quad (2)$$

$$184 \quad \mathbf{h}^4 = \text{Dropout}[\text{ReLU}(\text{concatenate}(\mathbf{x}, \mathbf{h}^3), \boldsymbol{\theta}^4, \boldsymbol{\beta}^4)] \quad (3)$$

$$185 \quad \mathbf{o} = \text{Dropout}[\text{ReLU}(\mathbf{h}^4, \boldsymbol{\theta}^4, \boldsymbol{\beta}^4)] \quad (4)$$

186 where \mathbf{x} denotes the input feature vector following PCA preprocessing; \mathbf{h} is the hidden layer feature
187 vector; $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ are the weight and bias vectors for each layer, respectively. concatenate refers to the
188 direct merging of two vectors. For activation functions, ReLU is frequently employed due to its rapid
189 convergence, computational simplicity, and absence of gradient vanishing, defined as $\text{ReLU}(v) =$
190 $\max(0, v)$. It is essential to note that, the model is used for both regression and classification, and in
191 classification models, the activation function in Eq. (4) should be changed as Softmax to ensure that
192 the sum of probabilities across all classifications equals 1.



193

194

Fig. 3 Architecture of the PCA–Shortcut Connection Neural Network (PCA-SCNN)

195 2.2 Classification of ACF Types and Estimation of SOFs

196

In this study, the training samples for Models #1-3 consist of synthetic zero-mean stationary

197

normal random fields q_c (Zhang et al., 2021). The CPT data depth is 10 m with a resolution of 0.05 m

198

thus ensuring sufficient information in the synthetic sample path. The CPT ranges from 0 to 100 m

199

horizontally with a resolution of 0.25 m. Given that the number of CPTs within a site is variable, four

200

scenarios of CPT quantities are considered (2, 3, 4, and 5). The analytical modeling process for

201

different CPT quantities follows a similar approach. To avoid redundancy, an example with 3 CPTs is

202

illustrated in Fig. 4. Apart from the CPT quantity, this study considers three other variables:

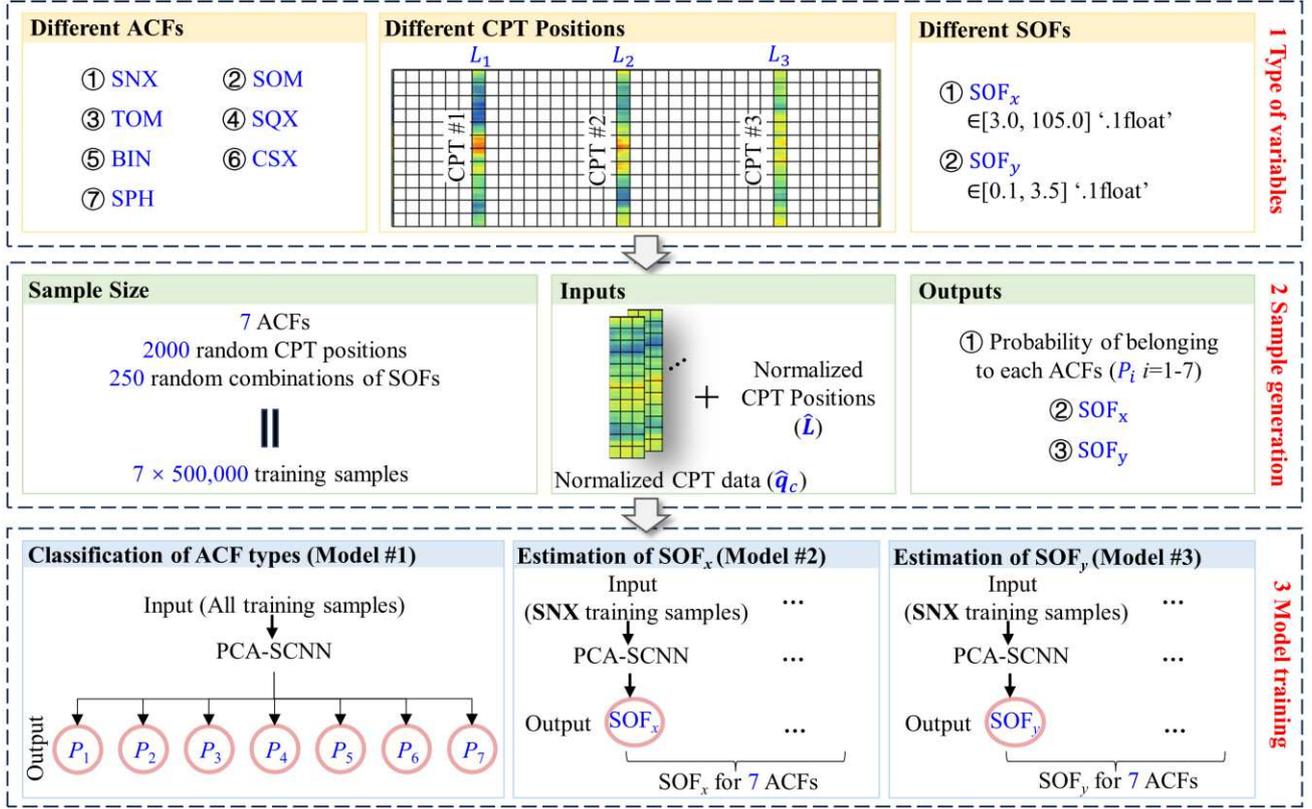


Fig. 4 Flow of ACF Type Classification and SOFs Estimation for 3 CPTs

(1) **Different ACF types.** Table 1 presents expressions for common ACFs (Cami et al., 2020),

where τ_h and τ_v represent the horizontal and vertical distances of two soil cells, δ_h and δ_v are the horizontal and vertical SOFs. This study opts for seven commonly used one-parameter ACFs prevalent in geotechnical engineering. While some two-parameter ACFs that allow defining smoothness have been developed (Ching et al., 2018), their application in practical cases is relatively limited.

Table 1. Frequently Used Autocorrelation Functions (ACFs) - Adapted from Cami et al. (2020)

Autocorrelation model	Autocorrelation function $\rho(\tau_h, \tau_v)$	Frequency of usage
Single exponential (SNX)	$\exp\left[-2\left(\frac{ \tau_h }{\delta_h} + \frac{ \tau_v }{\delta_v}\right)\right]$	47 %
Second-order Markov (SOM)	$\left(1 + 4\frac{ \tau_h }{\delta_h}\right)\left(1 + 4\frac{ \tau_v }{\delta_v}\right)\exp\left[-4\left(\frac{ \tau_h }{\delta_h} + \frac{ \tau_v }{\delta_v}\right)\right]$	4 %
Third-order Markov (TOM)	$\left(1 + \frac{16 \tau_h }{3\delta_h} + \frac{256}{27}\left(\frac{ \tau_h }{\delta_h}\right)^2\right)\left(1 + \frac{16 \tau_v }{3\delta_v} + \frac{256}{27}\left(\frac{ \tau_v }{\delta_v}\right)^2\right)\exp\left(-\frac{16}{3}\left(\frac{ \tau_h }{\delta_h} + \frac{ \tau_v }{\delta_v}\right)\right)$	New

Squared exponential (SQX)	$\exp\left[-\pi\left(\frac{\tau_h^2}{\delta_h^2} + \frac{\tau_v^2}{\delta_v^2}\right)\right]$	15 %
Binary noise (BIN)	$\begin{cases} \left(1 - \frac{ \tau_h }{\delta_h}\right)\left(1 - \frac{ \tau_v }{\delta_v}\right) & \tau_h \leq \delta_h \text{ and } \tau_v \leq \delta_v \\ 0 & \text{otherwise} \end{cases}$	9 %
Cosine exponential (CSX)	$\cos\left(\frac{ \tau_h }{\delta_h}\right)\cos\left(\frac{ \tau_v }{\delta_v}\right)\exp\left(-\left(\frac{ \tau_h }{\delta_h} + \frac{ \tau_v }{\delta_v}\right)\right)$	10 %
Spherical (SPH)	$\begin{cases} \left[1 - \frac{9 \tau_h }{8\delta_h} + \frac{27}{128}\left(\frac{ \tau_h }{\delta_h}\right)^3\right]\left[1 - \frac{9 \tau_v }{8\delta_v} + \frac{27}{128}\left(\frac{ \tau_v }{\delta_v}\right)^3\right] & \tau_h \leq \frac{4}{3}\delta_h \text{ and } \tau_v \leq \frac{4}{3}\delta_v \\ 0 & \text{otherwise} \end{cases}$	15 %

211 **(2) Different distributions of CPT positions.** Due to constraints in equipment, site conditions,
212 and costs, the positions of CPT measurements often exhibit irregular distribution. Existing data-driven
213 models solely rely on CPT measurement data as input (Zhang et al., 2022), potentially overlooking the
214 actual distances between different CPT data points. As depicted in Fig. 4, this study enhances the
215 model's adaptability in unevenly spaced scenarios by incorporating the coordinates (L) of each CPT
216 test as an additional input.

217 **(3) Different horizontal and vertical SOFs.** The horizontal SOF range from 3 to 105 m, and the
218 vertical SOF range from 0.1 to 3.5 m (Cami et al., 2020; Zhang et al., 2021).

219 Notably, assessing the ACF type and SOFs for unevenly spaced CPT data poses a more intricate
220 challenge. This study addresses this by (1) utilizing a PCA-SCNN with robust nonlinear fitting
221 capabilities to construct a data-driven model and (2) augmenting the training data volume.
222 Consequently, in the process of constructing training samples, seven common ACFs, 2000 sets of
223 random CPT positions, and 250 sets of random horizontal and vertical SOFs combinations are
224 considered. Therefore, the total training dataset comprises $7 \times 500,000$ samples, with 500,000 samples
225 for each ACF type. Detailed model information for 3 CPTs is provided in Table 2.

226 Table 2. Overview of the Architecture of the ACF Classification Model and the SOFs Estimation Model for 3 CPTs

Model	Number of models	Inputs	Outputs	Activation function & (number of neurons) in the output layer	Sample size
Model #1 (Classification of ACFs)	1	CPT data ($\hat{\mathbf{q}}_c$)	P_i	Softmax (7)	$7 \times 500,000$
Model #2 (Estimation of SOF_x)	7	& positions ($\hat{\mathbf{L}}$)	SOF_x	ReLU (1)	500,000
Model #3 (Estimation of SOF_y)	7		SOF_y	ReLU (1)	500,000

227 Note: Model #1 is utilized to determine the probability of CPT data belonging to the 7 ACF types. Therefore, the
228 training samples for Model #1 comprise all samples of the 7 ACF types ($7 \times 500,000$). The 7 neurons in the output
229 layer of Model #1 correspond to the probabilities P_i ($i=1-7$) of CPT data belonging to each ACF type. The activation
230 function used is Softmax, ensuring $\sum P_i = 1$. Model #2 and Model #3 are used to estimate the horizontal and vertical
231 SOFs of CPT data, respectively. Therefore, both Model #2 and Model #3 consist of 7 independent sub-models
232 corresponding to the 7 ACF types. Each sub-model's training samples only include samples of one ACF type (500,000
233 samples). Each sub-model's output layer contains only 1 neuron, and the activation function used is ReLU.

234 After training, Model #1-3 can be utilized to predict the probabilities P_i ($i=1-7$) of measured CPT
235 data \mathbf{q}_{cm} belonging to the 7 ACF types, as well as the corresponding SOF_x and SOF_y for each ACF type.

236 It is noteworthy that during the training process, the input of the Model #1-3 is $[\hat{\mathbf{q}}_c, \hat{\mathbf{L}}]$, where $\hat{\mathbf{q}}_c$ is the
237 normalized \mathbf{q}_c ; $\hat{\mathbf{L}}$ is the normalized position coordinate, $\hat{\mathbf{L}} = \mathbf{L}/100$. To ensure the accuracy of the
238 prediction results, the real measurement data \mathbf{q}_{cm} should closely resemble the training data $\hat{\mathbf{q}}_c$.

239 Therefore, the following steps are needed: ① First, apply the Box-Cox method to transform \mathbf{q}_{cm} into
240 a normal distribution (Zou et al., 2017), as shown in Eq. (5); ② Calculate the mean and standard
241 deviation of the normal distribution, and transform the data into a standard normal distribution; ③
242 Apply Min-Max normalization to transform the data back into $\hat{\mathbf{q}}_c$.

$$243 \quad \mathbf{q}_c = \begin{cases} \frac{\mathbf{q}_{cm}^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln(\mathbf{q}_{cm}) & \lambda = 0 \end{cases} \quad (5)$$

244 where λ is the power parameter that needs to be estimated. The optimal λ value can be efficiently
245 determined using common statistical software such as the Scipy library implemented in Python.

246 2.3 Geotechnical Subsurface Modeling Using GCFs

247 Currently, spatial positions of sampling points are commonly represented using Euclidean
 248 distance-based coordinates. However, these coordinates have low-dimensional features and lack
 249 effective constraints from random field theory. To address this issue, this study proposes a high-
 250 dimensional GCF that conforms to random field theory to characterize the spatial positions of sampling
 251 points: (1) The GCF aligns with random field theory, where the correlation between two soil cells
 252 within the site is calculated using random field theory rather than simple Euclidean distance. (2) The
 253 GCF has higher-dimensional features. Instead of using 2D or 3D relative coordinates to represent the
 254 positions of soil cells, it employs the correlation between each soil cell and every sampled soil cell, as
 255 illustrated in Fig. 5.

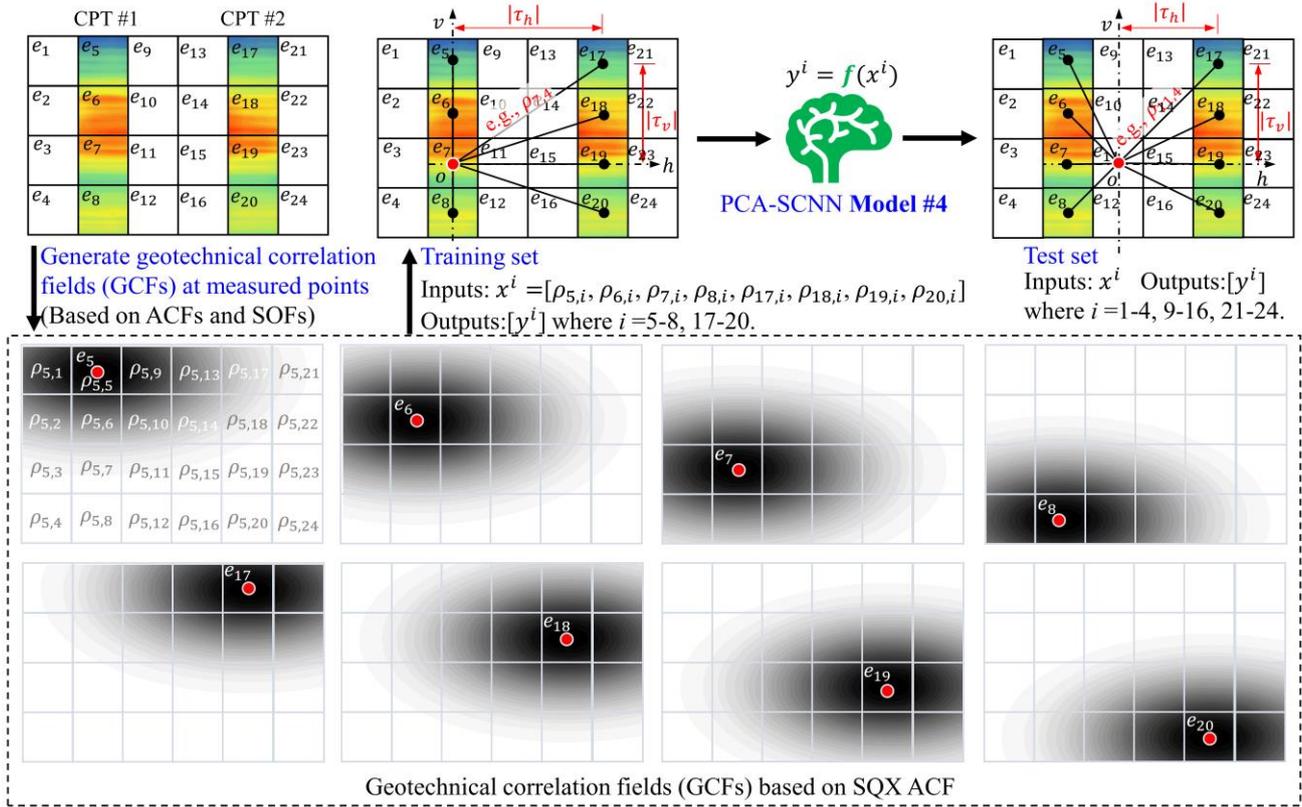
256 In this study, GCFs are employed to characterize the correlation between soil cells (calculated
 257 based on ACFs, SOFs, and relative distances), effectively integrating random field theory into the data-
 258 driven model. The process of generating GCFs is illustrated in Fig. 5. Assuming discretization of the
 259 geological site into a 4×6 grid, where the second and fifth columns represent sampled points, and the
 260 other columns represent unsampled points. Taking soil cell e_5 as an example, the GCF of e_5 is defined
 261 as the correlation matrix (GCF_{e_5}) between e_5 and other soil cells within the site:

$$262 \quad \text{GCF}_{e_5} = \begin{bmatrix} \rho_{5,1} & \rho_{5,5} & \rho_{5,9} & \rho_{5,13} & \rho_{5,17} & \rho_{5,21} \\ \rho_{5,2} & \rho_{5,6} & \rho_{5,10} & \rho_{5,14} & \rho_{5,18} & \rho_{5,22} \\ \rho_{5,3} & \rho_{5,7} & \rho_{5,11} & \rho_{5,15} & \rho_{5,19} & \rho_{5,23} \\ \rho_{5,4} & \rho_{5,8} & \rho_{5,12} & \rho_{5,16} & \rho_{5,20} & \rho_{5,24} \end{bmatrix} \quad (6)$$

263 where $\rho_{5,j}$ represents the correlation between soil cell e_5 and e_j , calculated using different ACFs from
 264 Table 1, and $\rho_{5,j} = \rho_{j,5}$. Considering that only the soil properties of cells e_5 to e_8 and e_{17} to e_{20} are
 265 known, a total of 8 GCFs can be generated: GCF_{e_5} to GCF_{e_8} and $\text{GCF}_{e_{17}}$ to $\text{GCF}_{e_{20}}$. Therefore, the
 266 spatial position of soil cell e_i ($i=1-24$) in the geological correlation field can be represented by an 8-

267 dimensional coordinate vector: $[\rho_{i,5}, \rho_{i,6}, \rho_{i,7}, \rho_{i,8}, \rho_{i,17}, \rho_{i,18}, \rho_{i,19}, \rho_{i,20}]$, which indicates the
268 correlation between e_i and the 8 sampled soil cells.

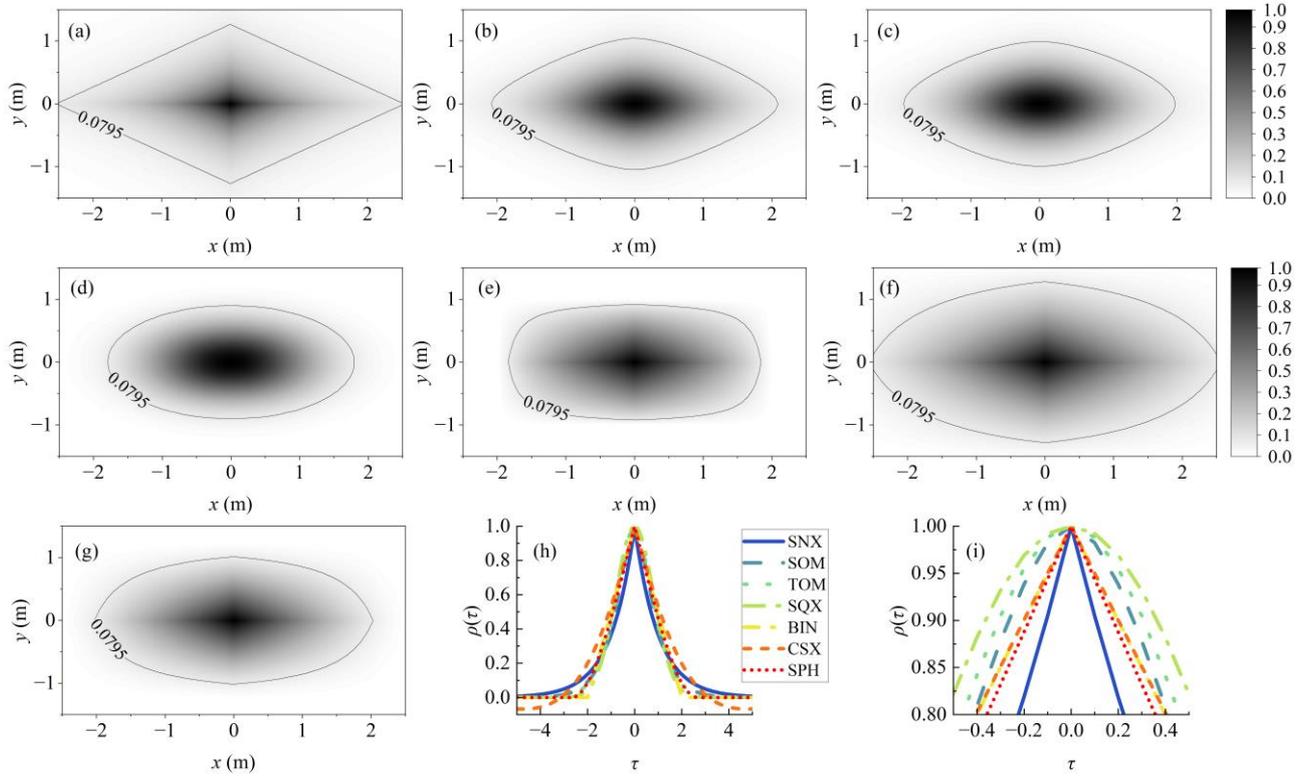
269 The essence of subsurface modeling based on GCFs is to utilize a data-driven model to learn the
270 relationship between the coordinate vectors of e_i and their corresponding geological properties,
271 represented as $\hat{q}_c^i = f([\rho_{i,5}, \rho_{i,6}, \rho_{i,7}, \rho_{i,8}, \rho_{i,17}, \rho_{i,18}, \rho_{i,19}, \rho_{i,20}])$, where f represents a complex
272 implicit function. In this study, PCA-SCNN is employed to solve f , which corresponds to Model #4 in
273 Fig. 2. The model takes an 8-dimensional coordinate vector as input and outputs the corresponding
274 geological properties of the soil cells. As shown in Fig. 5, the input and output of the training set for
275 Model #4 can be represented as $x^i = [\rho_{i,5}, \rho_{i,6}, \rho_{i,7}, \rho_{i,8}, \rho_{i,17}, \rho_{i,18}, \rho_{i,19}, \rho_{i,20}]$, $y^i = [\hat{q}_c^i]$, where $i = 5$ -
276 8 and 17-20. The test set can be represented as $x^j = [\rho_{j,5}, \rho_{j,6}, \rho_{j,7}, \rho_{j,8}, \rho_{j,17}, \rho_{j,18}, \rho_{j,19}, \rho_{j,20}]$, where
277 $j = 1$ -4, 9-16, and 21-24. Therefore, the training set consists of 8 samples, while the test set consists of
278 16 samples, all with a feature dimension of 8. If there are N discrete soil cells within the site, among
279 which M cells have been measured, then the training set will have M samples, and the test set will have
280 $N-M$ samples, all with a feature dimension of M .



281
 282 Fig. 5. Illustration of Geotechnical Correlation Fields and Associated Training and Test Set (Note: When the
 283 resolution of CPT does not align with the size of the soil cell, the average property of the soil cell is the mean value
 284 of the included measurement data.)

285 As depicted in Fig. 6 (a)-(g), using the seven different ACFs from Table 1, GCFs for soil cells at
 286 the coordinate origin (0 m, 0 m) are generated under horizontal and vertical SOFs of 2 m and 1 m,
 287 respectively. The contour lines of different GCFs exhibit distinct shapes, indicating significant
 288 variations in the influence range of the soil cells around the coordinate origin. Fig. 6 (h) and (g)
 289 illustrate one-dimensional correlation curves for different GCFs along the $y=0$ m section. While these
 290 curves demonstrate similar trends, GCFs generated by SNX, BIN, CSX, and SPH are not differentiable
 291 at zero lag, leading to lower smoothness of the corresponding sample paths (Ching et al., 2019). Real
 292 CPT data are often challenging to interpret with a single ACF. Vanmarcke (1983) proposed overlaying
 293 two or more ACFs to create a more flexible ACF. As shown in Fig. 2, this study adopts this overlay
 294 approach, generating seven sets of GCFs using seven sets of ACFs, followed by constructing seven

295 independent Model #4 for subsurface modeling. Finally, the predictions of the seven subsurface models
 296 are weighted by P_i to obtain a modeling result closely related to the measured data.



297
 298 Fig. 6. Correlations of Various GCFs Under Horizontal and Vertical Scale of Fluctuation (SOF) Set at 2 m and 1 m
 299 Respectively: (a) SNX, (b) SOM, (c) TOM, (d) SQX, (e) BIN, (f) CSX, (g) SPH; (h) Cross-Sectional Trends at $y=0$
 300 m for ACFs; (i) zoom-in view of (h)

301 2.4 Uncertainty Estimation of the Model

302 In this study, the uncertainty of the subsurface models is evaluated using the Monte Carlo Dropout
 303 structure of neural networks (P. Zhang et al., 2022). The MC dropout structure randomly deactivates a
 304 certain percentage of connections between neurons, introducing randomness into the model's
 305 architecture. Through repeated predictions (e.g., 100 times), the uncertainty of the output results can
 306 be directly obtained. As shown in Fig. 3, the models used for ACFs classification (Model #1),
 307 horizontal SOF estimation (Model #2), vertical SOF estimation (Model #3), and subsurface modeling
 308 (Model #4) all incorporate the MC dropout structure. Therefore, the output results of Model #1-4 are

309 associated with uncertainty.

310 It is noteworthy that different horizontal and vertical SOFs generate distinct GCFs, necessitating
311 the repetitive construction of different Model #4 for subsurface reconstruction. To minimize
312 computational costs, the most robust SOFs predictions are considered, obtained by running the SOFs
313 prediction model 100 times and averaging the results. Therefore, the primary source of uncertainty in
314 subsurface modeling stems from the combined contributions of ACFs classification (Model #1) and
315 subsurface modeling (Model #4). As shown in Fig. 2, the one subsurface modeling result \mathbf{Y} considers
316 the influence of 7 ACFs, represented as $\mathbf{Y} = \sum P_i \mathbf{Y}_i$, where i ranges from 1 to 7. To obtain uncertainty
317 in the prediction results, Model #1 and Model #4 need to predict 100 times each, yielding 100 sets of
318 P_i and \mathbf{Y}_i . Computing according to the aforementioned formula yields 100 subsurface modeling results.

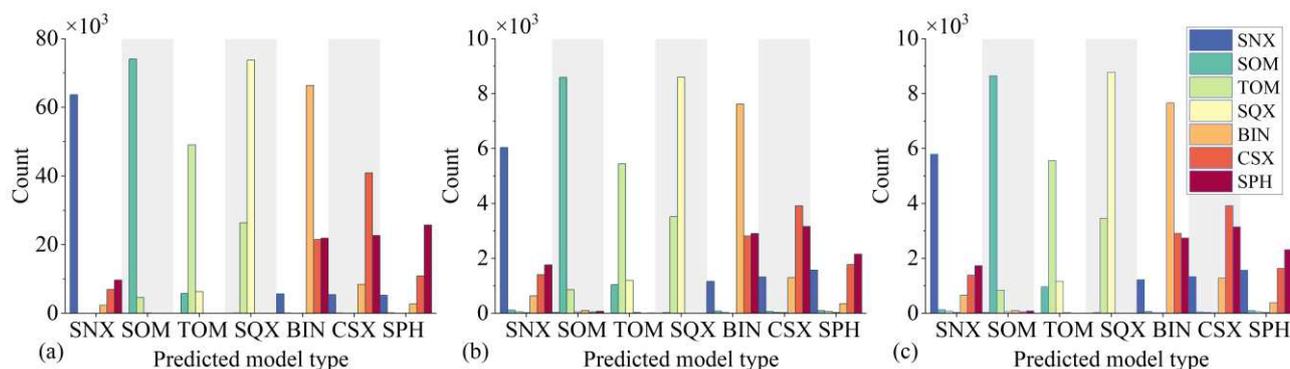
319 **3 Classification and SOF Prediction Models**

320 The PCA operation is implemented using Scipy v1.9.3, and the construction and training of the
321 SCNN are carried out using Tensorflow-GPU v2.8.0—both are open-source packages developed in
322 Python. The neural network consists of four hidden layers, each containing 512 neurons, with a dropout
323 rate set to 0.3. For SCNN, a batch size of 512 is chosen to expedite the model training process. The
324 initial learning rate is set at 0.001, with a 0.5 reduction if the loss on the validation set does not decrease
325 for 15 consecutive iterations. Early stopping is employed to control the number of model iterations,
326 terminating the training process if the model's loss on the validation set does not improve for 30
327 consecutive iterations. These hyperparameters are determined through a grid search approach. The
328 training uses the Nadam optimizer, an extension of the Adam optimizer with RMSprop and Nesterov
329 momentum.

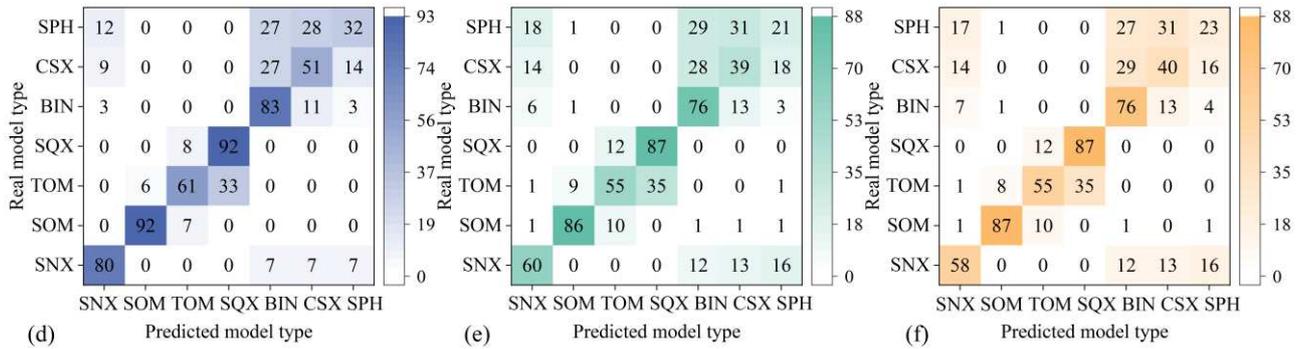
330 **3.1 Training and Validation of Classification Models**

331 The classification of ACF is a multi-class problem. Hence, the classification cross-entropy loss
 332 function is used and the activation function for the output layer is set to Softmax. One-hot encoding is
 333 employed to represent ACF categories. Each dimension of the one-hot encoding represents the
 334 probability that the measured data belongs to a specific category. For example, the SNX category can
 335 be represented as [1, 0, 0, 0, 0, 0], while the SOM category is represented as [0, 1, 0, 0, 0, 0], and so
 336 on.

337 To avoid redundancy, Fig. 7 illustrates the performance of the classification model on the training,
 338 validation, and test sets using only three CPT quantities. It can be observed that the classification model
 339 exhibits high prediction accuracy on the training set, and its performance on the validation and test
 340 sets is similar. This indicates the classification model is able to capture the complex relationship
 341 between CPT data and ACF categories and does not suffer from overfitting. Furthermore, the model
 342 can distinguish between SNX, SOM, and SQX categories. The TOM sample paths exhibit a degree of
 343 smoothness between SOM and SQX, and the model occasionally misclassifies TOM as SOM or SQX.
 344 In the case of BIN, CSX, and SPH sample paths, their smoothness is quite similar, leading to potential
 345 misclassifications among the three models. However, these misclassifications are acceptable, as even
 346 in cases of misclassification, the model provides similar predictions.



347



348

349

350

Fig. 7. Classification Model Performance with 3 Sets of CPT Measurements: Training (a, d), Validation (b, e), Test (c, f) Sets

351

352

353

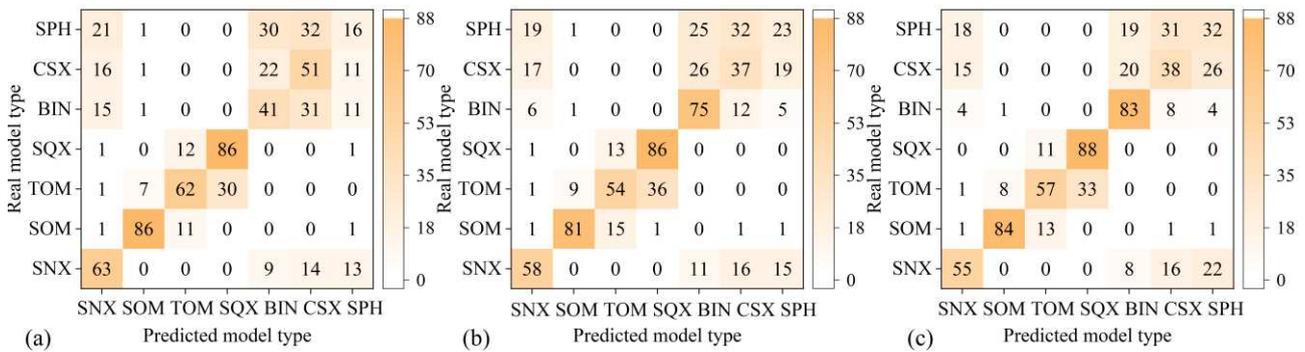
354

355

356

357

As shown in Fig. 8, the classification model exhibits strong predictive performance on the test set across different numbers of CPT data. The classification model maintains a high level of prediction accuracy even with just two CPT curves. As the measurement data increases, the model's accuracy in predicting BIN significantly improves. As depicted in Fig. 6, for a given SOF, the influence range of the BIN model varies considerably from that of the CSX and SPH models. Therefore, with an increasing amount of measurement data, the model can not only consider the smoothness of sample paths for classification but also effectively account for the influence range of SOF.



358

359

Fig. 8. Test Set Performance of the Classification Model: (a) 2CPT, (b) 4CPT, (c) 5CPT

360

3.2 Training and Validation of SOF Models

361

362

The prediction of SOF is a regression problem. The model employs mean squared error (MSE) as the loss function and applies the ReLU activation function in the output layer. All other configurations

363 remain consistent with the classification model. Root mean square error (RMSE), mean absolute
 364 percentage error (MAPE), and the coefficient of determination (R^2) are used to assess the differences
 365 between the predicted values (\hat{y}_i) and the measured values (y_i) and can be expressed as follows:

$$366 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (6)$$

$$367 \quad MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100\% \quad (7)$$

$$368 \quad R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (8)$$

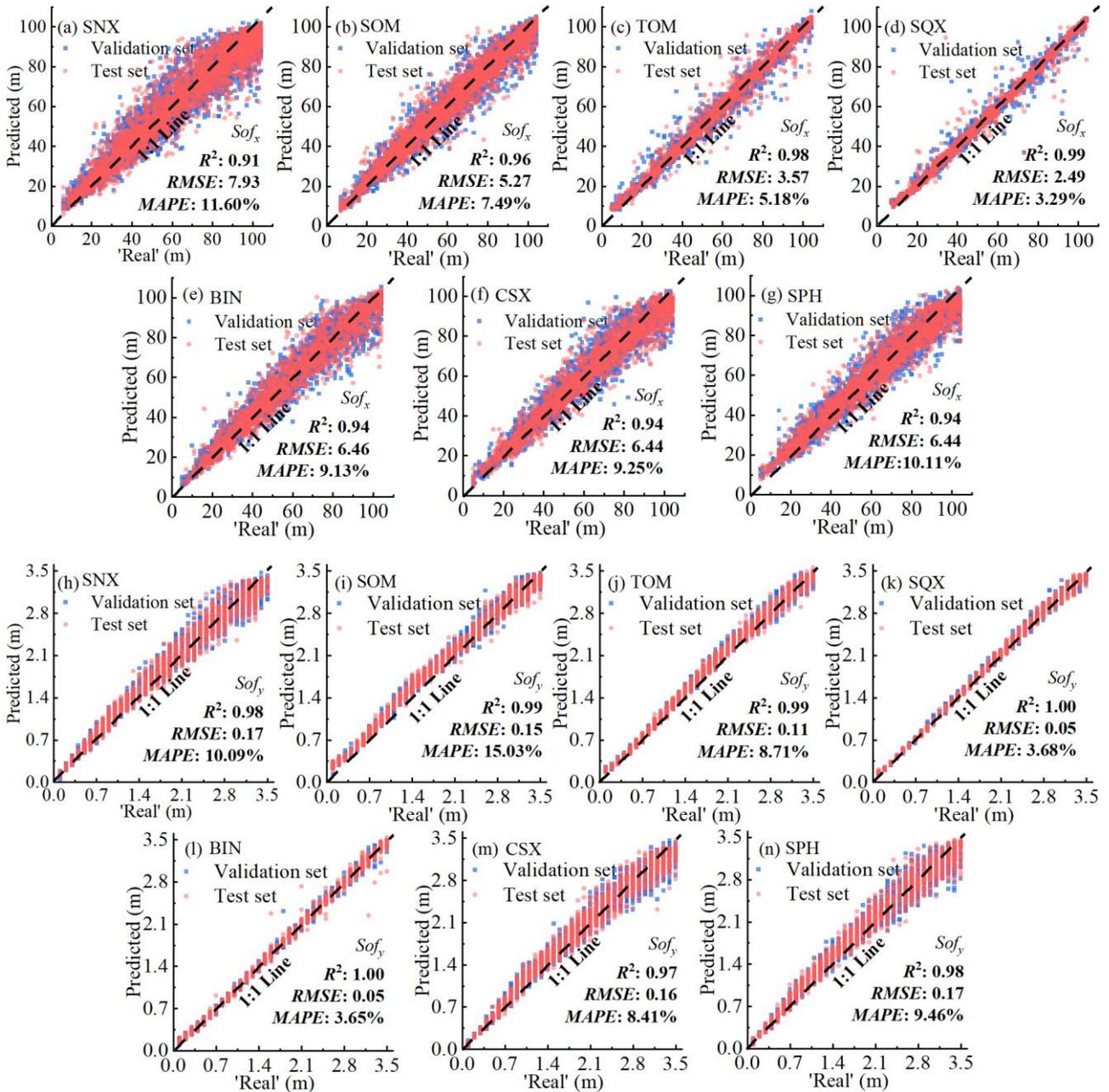
369 where \bar{y} is the mean q_c value; n is the total number of samples. When R^2 is close to 1, and RMSE and
 370 MAPE are relatively small, the predictive performance of the model is better.

371 Usually, the training set exhibits higher accuracy compared to the validation and test sets.
 372 Therefore, only the prediction results of the models for different ACFs in the validation and test sets
 373 are displayed. As shown in Fig. 9 (a)–(g), the horizontal SOF prediction models all demonstrate good
 374 performance, with the SQX model exhibiting the highest prediction accuracy. Its test set's R^2 is close
 375 to 1, with RMSE and MAPE values of only 2.49% and 3.29%, respectively. Although the SNX model
 376 has marginally lower accuracy, it still performs well on the test set, with an R^2 of 0.91 and RMSE and
 377 MAPE values of 7.93 and 11.60%, respectively.

378 As shown in Fig. 9 (h)–(n), the vertical SOF prediction models all exhibit improved performance.
 379 Among them, the SQX model shows the highest prediction accuracy. Its test set's R^2 is 1.00, with
 380 RMSE and MAPE values of 0.05 and 3.68%, respectively. The SNX model has marginally lower
 381 prediction accuracy, with an R^2 of 0.98 and RMSE and MAPE values of 0.17 and 10.09%, respectively.

382 It can be observed that the model's prediction performance is related to the smoothness of the

383 sample path. The smoother the sample path, the higher the prediction accuracy of the model. When the
 384 sample path is very rough, it becomes more challenging to distinguish whether the fluctuations in the
 385 measurement data are caused by the SOF or the roughness of the sample path itself. However, overall,
 386 the established prediction models demonstrate good accuracy.



387

388

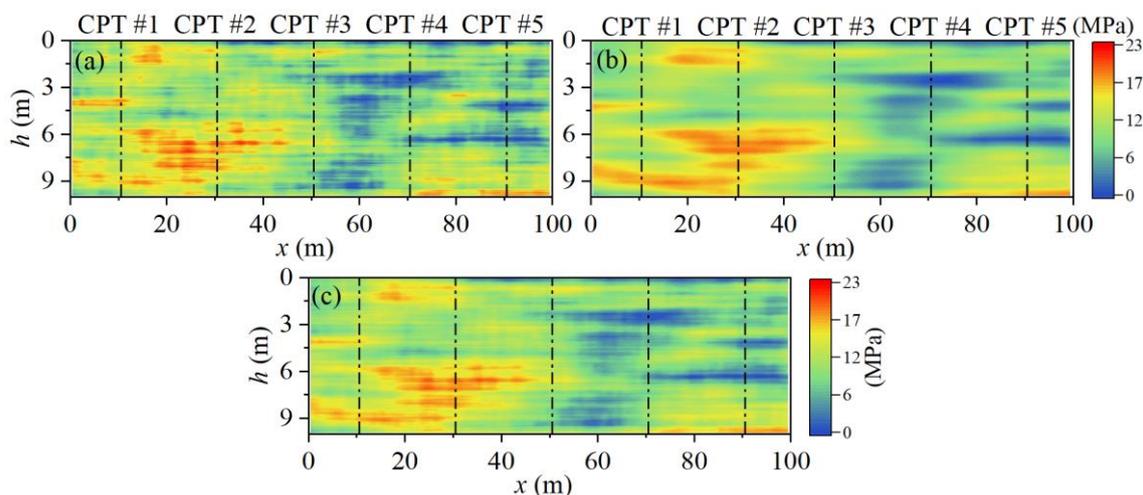
389

Fig. 9. SOF Prediction Model Performance: Horizontal SOF (a–g); Vertical SOF (h–n)

390 4 Subsurface Modeling and Validation

391 4.1 Synthetic Case Study

392 It is uncommon to record high-resolution test data along the surface of a site. Therefore, this
393 section illustrates the proposed method using a set of synthetic two-dimensional Gaussian q_c field, as
394 depicted in Fig. 10. The two-dimensional cross-section has a depth (h) of 10 m and extends 100 m
395 along the surface. The q_c field is simulated with a resolution of 0.05 m and 1.0 m along the depth and
396 horizontal directions, respectively. In this example, the mean (μ) and standard deviation (σ) are taken
397 as 10 MPa and 4 MPa, respectively. The horizontal and vertical SOF are set as 50 m and 2 m. As shown
398 in Fig. 10, under fixed random seeds, simulation results for two types of ACF, SNX, and SOM, are
399 generated using the matrix decomposition method. The sample paths of SNX exhibit more roughness
400 compared to SOM. The results of SNX and SOM are averaged to obtain synthetic data with sample
401 path smoothness between the two, as shown in Fig. 10 (c).



402

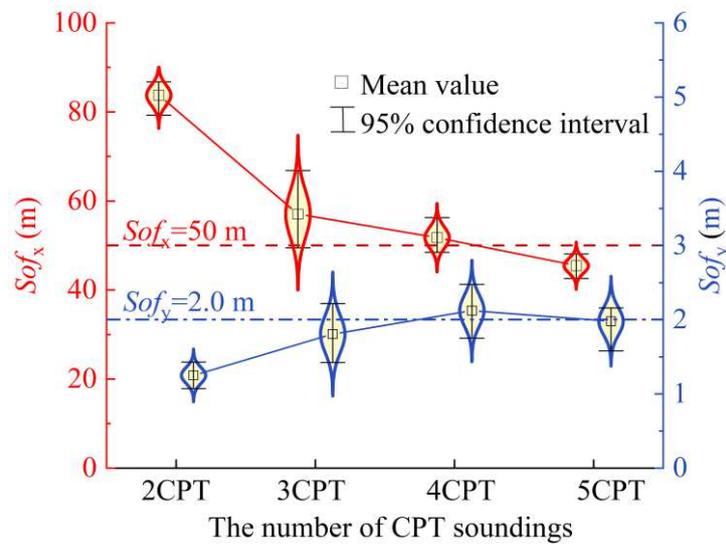
403

Fig. 10. Synthetic q_c Stratigraphy Depictions: (a) SNX, (b) SOM, (c) Averaged

404 4.2 Model Construction and Validation

405 To validate the influence of different CPT quantities on modeling accuracy, the measurement

406 locations are $x = 10.5$ m, 30.5 m, 50.5 m, 70.5 m, and 90.5 m. CPT #1 and CPT #5 are not placed at
 407 boundary positions, aiming to assess the extrapolation capability of the proposed method. As illustrated
 408 in Fig. 11, when the CPT quantity is two, the horizontal distance between CPT #1 and CPT #5 is
 409 significantly greater than the actual horizontal SOF. This leads to challenges in accurately predicting
 410 the horizontal SOF of the model. However, with more than three CPT measurements, the predicted
 411 mean values of the SOF closely align with the actual values.



412

413

Fig. 11. Predictions of Horizontal and Vertical SOFs for Different CPT Quantities

414

415

416

417

418

419

420

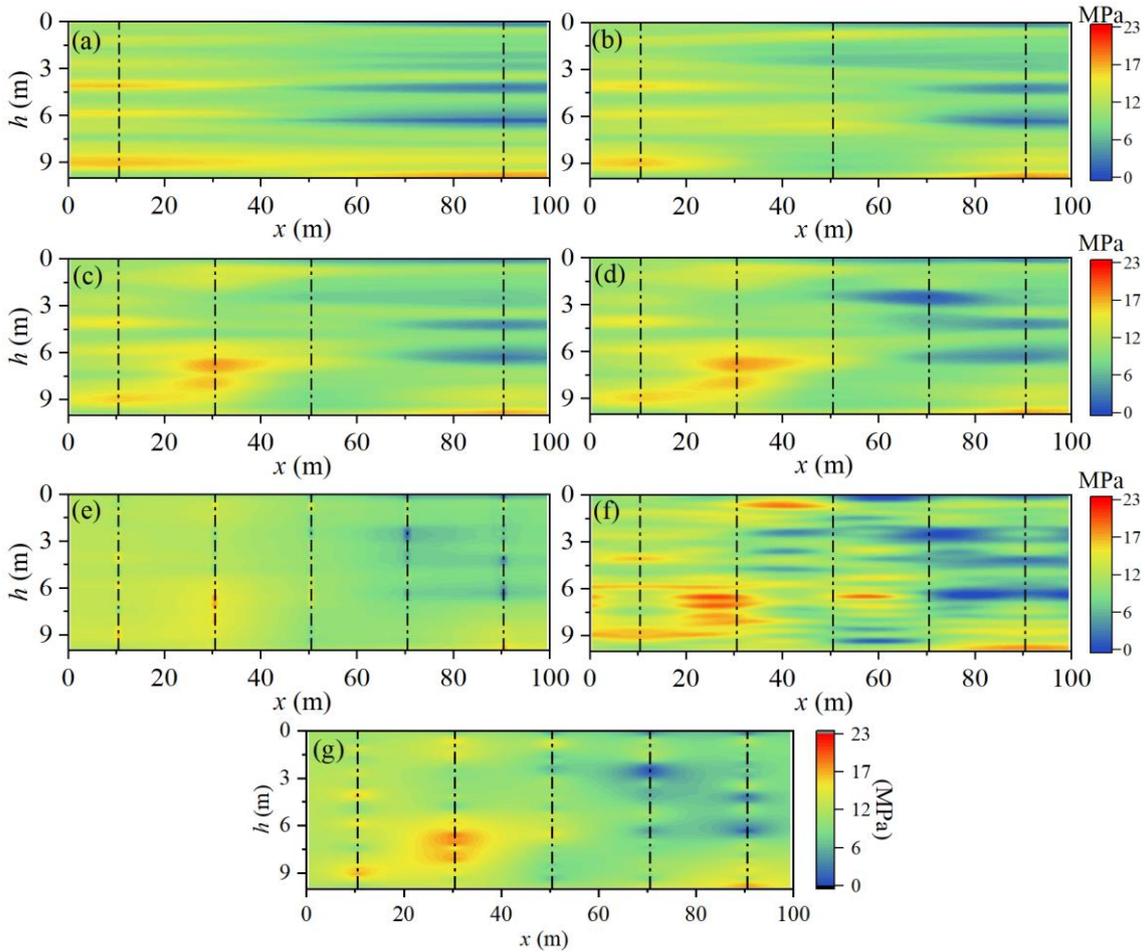
421

As shown in Fig. 12 (a)–(d), subsurface modeling is performed using the proposed method with
 2, 3, 4, and 5 CPTs, respectively. It is evident that, with an increase in the number of CPTs, the
 modeling results become progressively finer. When only 2 CPTs are used, the predicted horizontal
 SOF is relatively large, indicating the model has limited capability to predict the spatial distribution of
 soil properties near $x=50$ m. However, it accurately predicts the soil properties near the measurement
 locations. With 3 CPTs, the model's predictions near the measurement locations closely resemble those
 with more CPT records. The predictions corresponding to 4 and 5 CPTs are in close agreement,
 indicating that the information from CPT #4 is already well-predicted with just 4 CPTs.

422 As shown in Fig. 12 (d)–(f), the proposed method, GDF–ET (Xie et al., 2022b), BCS (Wang et al.,
423 2020; Zhao et al., 2020), and Kriging (Zou et al., 2017) are presented for the case of 5 CPTs. Both the
424 Kriging and the proposed methods necessitate the estimation of random field parameters. This is
425 typically accomplished through the method of moments, maximum-likelihood estimation, and
426 Bayesian analysis (Ching et al., 2020; Ching and Phoon, 2019; Liu et al., 2017; Liu and Leung, 2018;
427 Xiao et al., 2018, 2016). The maximum-likelihood estimation further facilitates the selection of the
428 optimal ACF model using either the Akaike information criterion (AIC) or the Bayesian information
429 criterion (BIC) (Chang et al., 2021). As illustrated in Fig.12 (g), the Kriging method employs the SNX
430 model, which aligns with the ACF model used for generating synthetic cases in Fig. 10, and exhibits
431 smaller AIC and BIC values.

432 As shown in Fig. 12 (d)–(f), GDF–ET excels at predicting the mean values of the soil properties
433 within the stratigraphy. BCS produces more intricate predictions compared to synthetic stratigraphy.
434 In cases where soil property variations are relatively gradual, BCS often yields accurate predictions.
435 Notably, when the subsurface modeling process lacks the constraints of geotechnical knowledge, the
436 complexity of predictions in different data-driven methods is often related to the algorithm's basis
437 functions or the complexity of input features. This can lead to over-simplification or over-complication
438 of predictions in some specific stratigraphy. Embedding geotechnical knowledge can help prevent such
439 occurrences. The Kriging method exhibits high precision in areas close to CPT boreholes, but tends to
440 estimate using the mean of nearby areas when positioned between two boreholes (Zou et al., 2017).
441 Notably, the Kriging method exhibits cubic complexity (Y. Yang et al., 2022) which may consume
442 excessive computational resources and result in slow solutions when dealing with large amounts of

443 soil cells. The method proposed in this study benefits from the control of random field theory, ensuring
 444 a good consistency between the prediction results and the actual site. Simultaneously, the proposed
 445 method is modeled based on a data-driven approach, allowing it to be applied to situations with large
 446 amounts of data at a lower computational complexity, as detailed in Section 4.4.



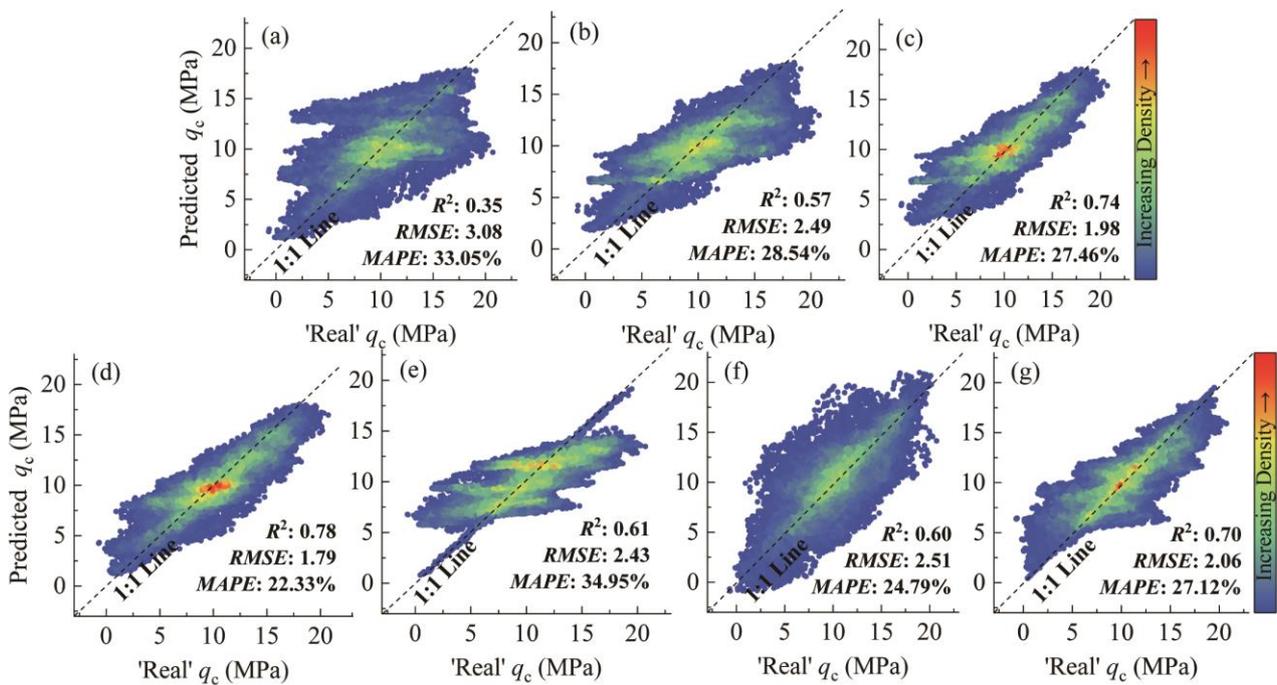
447

448

449 Fig. 12. Comparative q_c Modeling Results: Proposed Method (a–d) at 2, 3, 4, and 5 CPTs; GDF-ET Method (e),
 450 BCS Method (f) and Kriging Method (g) at 5 CPTs

451 As illustrated in Fig. 13 (a)–(d), with an increasing number of CPTs, R^2 of the predictions
 452 consistently improves, while RMSE and MAPE decrease, indicating improved prediction accuracy.
 453 Notably, since the synthetic random field follows a normal distribution, a substantial amount of data
 454 is around the mean value (10 MPa). As the volume of measurement data increases, the predictions
 455 gradually approach the 1:1 line. In Fig. 13 (d)–(f), with 5 CPTs, the R^2 values for the GDF-ET and

456 BCS models are close, with the BCS model having a slightly larger RMSE compared to the GDF-ET,
 457 and GDF-ET model exhibiting a higher MAPE than BCS. The predictions by the GDF-ET method
 458 mostly fall within the range of mean \pm one standard deviation, while the BCS method exhibits overall
 459 better consistency across the entire range. Benefiting from accurate random field parameters, the
 460 Kriging method's predictive results surpass those of the GCF-ET and BCS methods. The proposed
 461 method shows higher R^2 and lower RMSE and MAPE with better accuracy compared to the other
 462 methods.



463
 464 Fig. 13. Comparison of Modeling Results: Proposed Method (a–d) at 2, 3, 4, and 5 CPTs; GDF-ET Method (e),
 465 BCS Method (f) and Kriging Method (g) at 5 CPTs

466 As shown in Fig. 14, the predicted results and their confidence intervals at $x = 40.5$ m and 80.5 m
 467 are extracted. It is evident that the overall trend of the predicted results closely aligns with the actual
 468 values, and the majority of observed data falls within the 95% confidence interval of the predictions.
 469 Additionally, the predictions from the GDF and kriging models tend to converge toward the mean,
 470 while the BCS method exhibits unexplained fluctuations at certain locations.

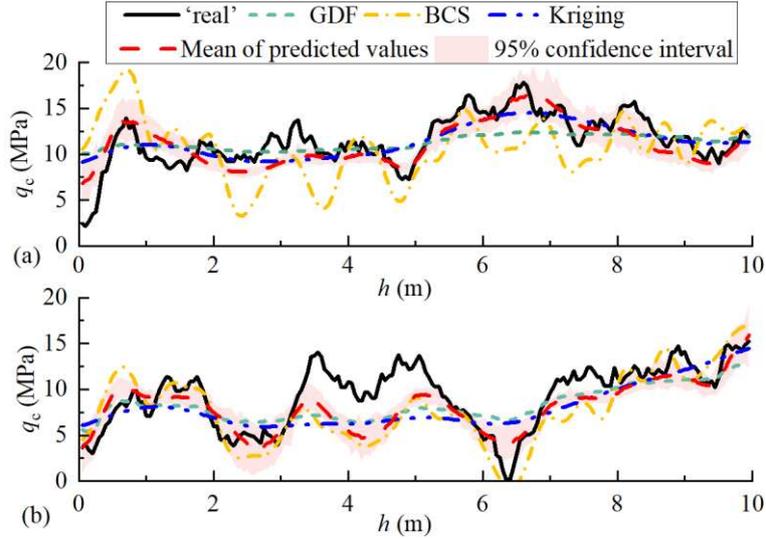
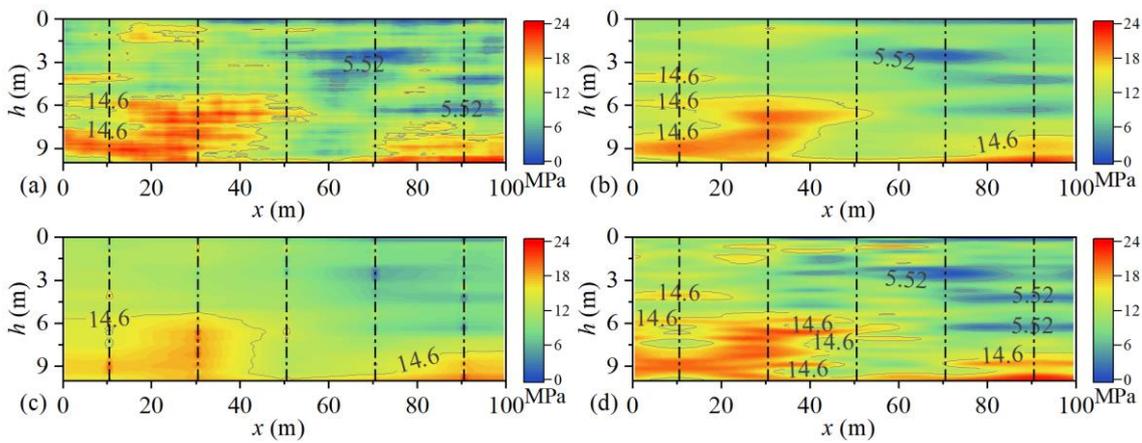


Fig. 14. Model Uncertainty Assessment at Locations: (a) $x=40.5$ m, (b) $x=80.5$ m

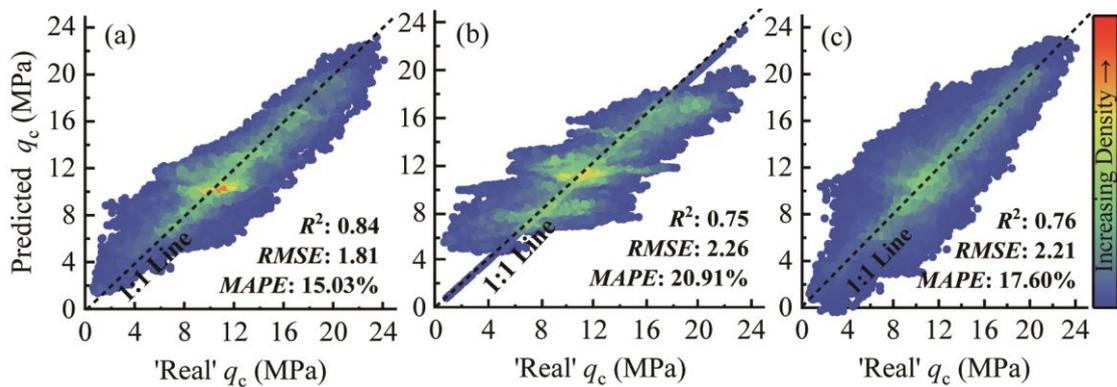
4.3 Nonstationary Synthetic Case Study

It is worth noting that when evaluating random field parameters, the input for the ACF classification model and SOF prediction model is the CPT data after detrending. This section primarily evaluates the capability of the proposed method to perform subsurface modeling directly on non-stationary data (measured CPT data) after obtaining random field parameters. This approach helps to improve modeling efficiency and reduce uncertainties caused by detrending. Therefore, a non-linear trend is introduced into the stationary random field depicted in Fig. 10(c). The synthetic data assumed a simple second-order increasing trend in the depth direction: $0.05h^2$, followed by Xie et al. (2022b). Fig. 15(a) illustrates the synthesized non-stationary random field. Fig. 15(b)-(d) show predictive results based on five CPT datasets using the proposed method, GDF, and BCS. It can be observed that all three methods capture the trends of the site well. Among them, GDF-ET and BCS show some simplification or complication compared to the synthetic stratigraphy. Although the proposed method's predictive results exhibit some simplification, due to the constraints of random field information, it effectively recovers most of the information from the synthetic stratigraphy. Combining these results

487 with Fig. 16 reveals consistently higher R^2 values for the proposed method compared to the GDF and
 488 BCS methods, with relatively minimal RMSE and MAPE. Compared to modeling results for stationary
 489 data using GCFs, as shown in Fig. 13 (d), direct modeling results for non-stationary data show higher
 490 R^2 and lower MAPE values, with RMSE values close. This indicates the proposed method achieves
 491 improved modeling accuracy and is applicable for both stationary and non-stationary data with a mild
 492 trend.



493
 494 Fig. 15. Modeling q_c Results Comparison: (a) Actual, (b) Proposed Method, (c) GDF-ET, (d) BCS



495
 496 Fig. 16. Comparative Analysis of Different Modeling Methods: (a) Proposed Method, (b) GDF-ET, (c) BCS

497 4.4 Computational Complexity of Subsurface Modelling

498 Illustrating the computational efficiency and complexity of the proposed subsurface modeling
 499 method using the non-stationary case described in Section 4.3. As shown in Fig. 15 (a), the synthetic

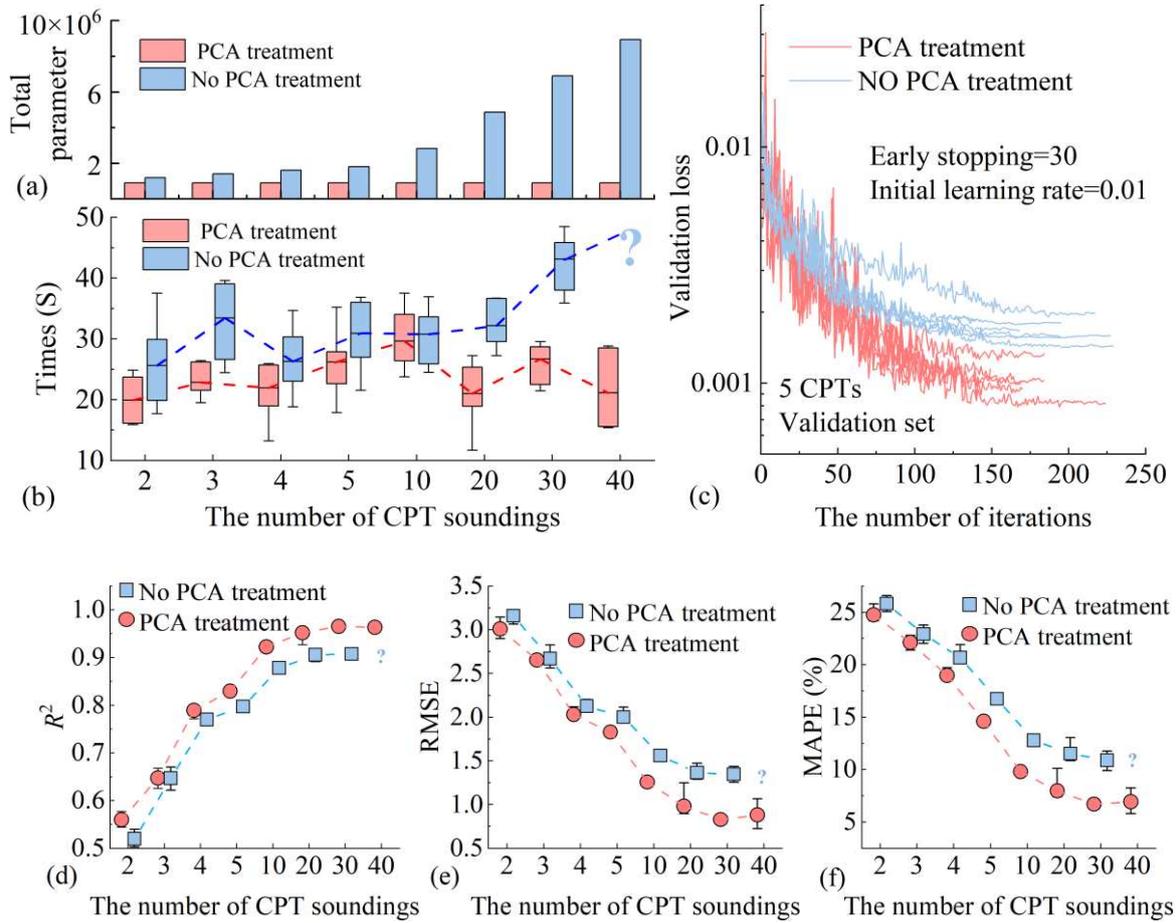
500 site consists of 20,000 discrete soil cells arranged in a 200×100 grid. Each CPT borehole includes 200
501 soil cells. Considering the case of five CPTs, a total of 1,000 soil cells are measured (5×200).
502 Consequently, the corresponding training and testing datasets consist of 1,000 and 19,000 samples,
503 respectively. Each sample has 1,000 features (denoted as M), as described in Section 2.3.

504 As depicted in Fig. 17(a), a comparison is made between the total trainable parameters in the
505 subsurface model with and without the PCA operation module. Since CPT tests provide nearly
506 continuous information in the vertical direction, the input feature dimension (M) is typically large. As
507 the number of CPTs increases, not applying PCA to the input data leads to an exponential growth in
508 the neural network's parameter count, significantly reducing its computational efficiency. The PCA-
509 SCNN model proposed in this study preprocesses the input data with PCA dimensionality reduction
510 before feeding it into the neural network. Hence, the feature dimensionality M of the input data does
511 not affect the model's complexity.

512 As depicted in Fig. 17 (b), after applying PCA processing to the input data, the training time for
513 an individual subsurface model stabilizes at around 20 seconds. The training time remains independent
514 of the number of CPTs. Furthermore, PCA processing effectively reduces the training time, especially
515 when dealing with a larger quantity of CPTs. Notably, if PCA processing is not performed, when the
516 number of CPTs reaches 40, the model consumes a significant amount of memory, increase training
517 difficulty, and may even become untrainable.

518 Fig. 17(c) illustrates the training process of the subsurface modeling model using 5 CPTs as an
519 example. Models with PCA-processed input not only converge more easily but also exhibit generally
520 lower validation losses. The PCA operation significantly reduces the training complexity of the model.

521 In Fig. 17 (d)-(f), we evaluate the model's prediction results using R^2 , RMSE, and MAPE metrics. As
 522 the number of CPTs increases, the model's prediction accuracy steadily improves, with the PCA-
 523 processed model consistently outperforming the non-PCA model.



524

525

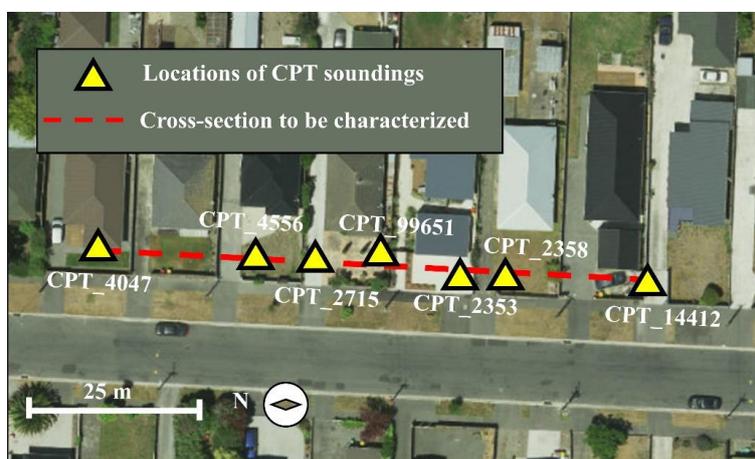
526 Fig. 17 The Impact of PCA Operations and the Number of CPTs on the Complexity of Subsurface Model (Model
 527 #4): (a) Increasing the Number of CPTs Affects the Trainable Parameter Count in Model #4. (b) The Effect of
 528 Increasing CPT Numbers on the Training Time of Model #4. (c) The Training Process of Model #4 with 5 CPTs.
 529 (d)-(f) Evaluation of the Subsurface Modeling Results Using R^2 , RMSE, and MAPE Metrics.

530 5 Real Data Case Study

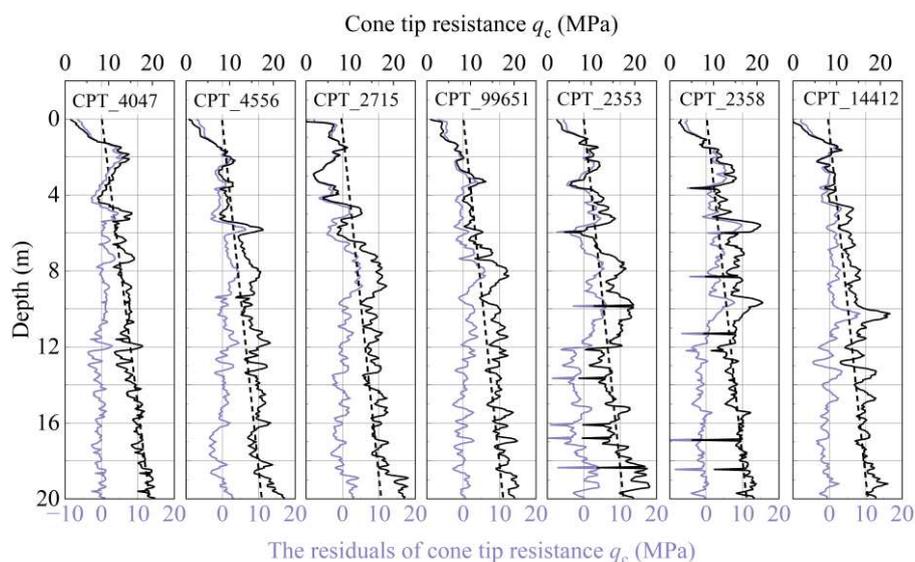
531 A set of CPT data from the Christchurch region in New Zealand is used to further demonstrate the
 532 proposed method. As shown in Fig. 18, a total of seven sets of q_c data are selected for subsurface
 533 modeling, taken from the New Zealand Geotechnical Database (NZGD) (NZGD, 2023). It is important
 534 to note that the measurement locations of these seven sets of CPT data are not strictly aligned along a

535 straight line. Therefore, CPT soundings are projected onto the two-dimensional vertical cross-section
 536 represented by the red dashed line in Fig. 18. The two-dimensional vertical cross-section extends
 537 approximately 82 m along the ground surface, with q_c data typically collected to depths of around 20
 538 m below the surface, as illustrated in Fig. 19. The codes and positions of the seven datasets in NZGD
 539 are as follows: CPT_4047 (10.5 m), CPT_4556 (33.5 m), CPT_2715 (42.5 m), CPT_99651 (47.5 m),
 540 CPT_2353 (63.5 m), CPT_2358 (70.5 m), and CPT_14412 (92.5 m). For this case, the q_c data have
 541 spatial resolutions of 1 m along the ground surface and 0.05 m in the depth direction.

542



543 Fig. 18. Geographical Layout and Cross-Section of 7 CPT Soundings in Christchurch, New Zealand (NZGD, 2023)



544

545

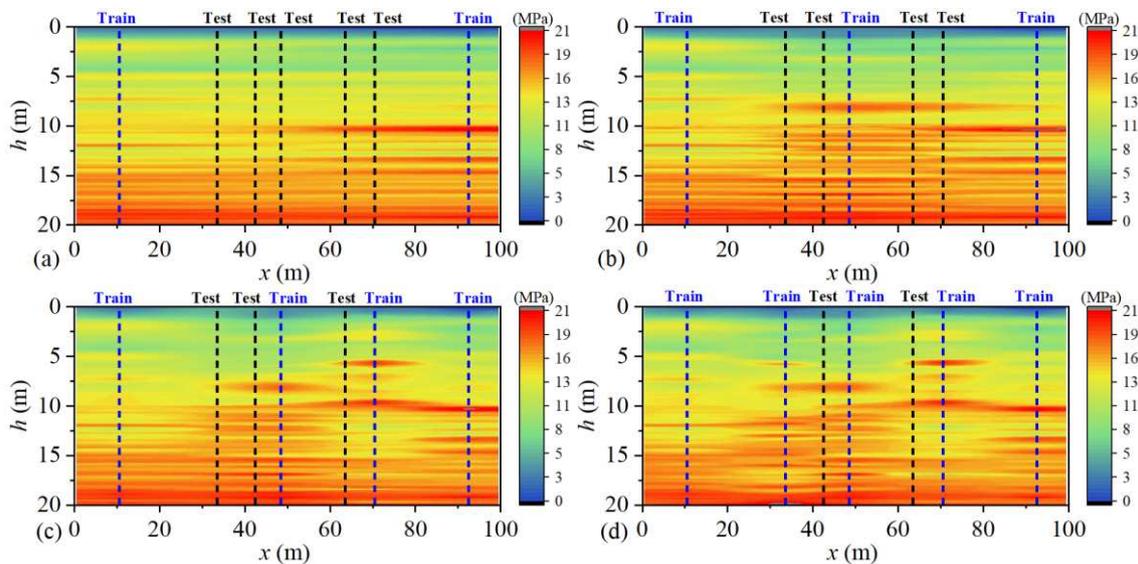
Fig. 19. Set of Seven Measured CPT Data

546 Given that actual CPT data consists of trend and residual terms, this study employs a linear
547 function for detrending, as illustrated in Fig. 19. The modified Bartlett statistical test (Phoon et al.,
548 2003) is employed to assess the stationarity of the residual term after detrending. It is worth noting that
549 the critical Bartlett peak value, B_{crit} , used in the modified Bartlett statistical test to determine the
550 stationarity of the sample sequence, depending on the ACF type. Using a classification model for
551 prediction, it is found that the actual CPT data has the highest probability of belonging to the SOM
552 and CSX models. B_{crit_SOM} and B_{crit_CSX} values of 71.85 and 48.87, detailed calculation steps are
553 available in (Phoon et al., 2003). The B_{stat} values for the residual terms from all CPT measurements
554 are mostly below 71.85, with the majority falling below 48.87. Therefore, when using the SOM ACF,
555 the null hypothesis of weak stationarity for the residual terms cannot be rejected at a significance level
556 of 5%. In the case of CSX, the assumption of weak stationarity holds for most scenarios. Furthermore,
557 it is worth noting that within the 7 sets of measured data, there are noticeable outliers, such as
558 CPT_2353 and CPT_2358. These outliers are common in real-world data, and in this study, no special
559 treatment is applied to them. This is done to further test the robustness of the proposed method.

560 Section 4.3 validates the proposed method's applicability to both stationary and non-stationary
561 data. Therefore, in this case, a detrending operation is solely applied during the assessment of random
562 field parameters. Subsequently, the GCFs, based on the obtained random field parameters, are used as
563 the model's inputs, with non-stationary observed data serving as output for subsurface modeling. As
564 shown in Fig. 20, subsurface modeling is conducted using different numbers of CPTs, with the
565 remaining CPTs serving as the test set for validation. With only 2 CPTs, the model struggles to provide
566 a detailed prediction of soil properties in the horizontal direction. As more CPT datasets are

567 incorporated, the modeling results become progressively more refined. Notably, near $x=70\text{m}$, the
568 modeling results exhibit a lens-shaped (hole effect) distribution.

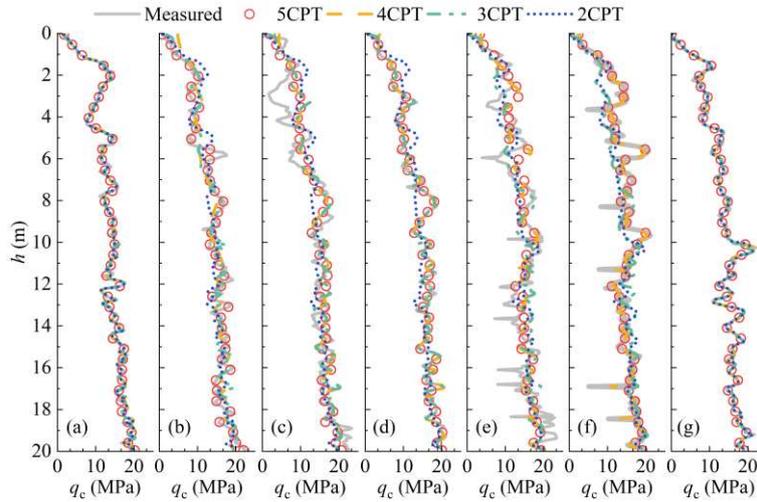
569 As depicted in Fig. 21, the model's predictions on the training set are nearly identical to the
570 observed values, indicating that using correlation as inputs for the model provides enhanced non-linear
571 expressive capabilities. On the test set, the model predictions closely match the observed values,
572 suggesting the model can predict soil properties in unknown areas based on the correlation
573 relationships established in the training set. It's worth noting that the model does not appear to be
574 affected by the outlier data from CPT_2353 and CPT_2358. This further validates the robustness of
575 the method, showing it can provide accurate predictions, even when dealing with anomalies. This
576 feature reduces the complexity of applying the method for subsurface modeling on a large scale, thus
577 making it more accessible to geotechnical engineers.



578

579

Fig. 20. Subsurface Modeling of Varied Numbers of CPT data: (a) 2; (b) 3; (c) 4; (d) 5



580

581 Fig. 21. Predictions of Measurement Locations of Varied Numbers of CPT Data: (a) CPT_4047, (b) CPT_4556, (c)
 582 CPT_2715, (d) CPT_99651, (e) CPT_2353, (f) CPT_2358, and (g) CPT_14412

583 6 Conclusions

584 This study introduces an innovative data-driven framework for soil property recovery which
 585 integrates geotechnical knowledge. This framework attempts to recover soil properties at unsampled
 586 points using sparse geotechnical measurements. Based on results from this study, the following
 587 conclusions can be drawn:

588 (1) The use of geotechnical correlation fields as inputs for the subsurface reconstruction model
 589 align with the fact that soil properties exhibit location-specific dependencies. This integration of
 590 random field theory into the data-driven model fosters enhanced collaboration between the model and
 591 geotechnical engineers. The superiority of the data-driven model has been validated through
 592 experiments with two synthetic random fields and a real-world case study.

593 (2) The subsurface modeling method proposed in this study benefits from PCA-SCNN model's
 594 dimensionality reduction of the input. The computational complexity of the model is independent of
 595 the number of soil cells, resulting in stable modeling efficiency and excellent scalability.

596 (3) A robust model for identifying types of autocorrelation functions is proposed, and this model

597 explicitly estimates the probability of observed data belonging to a specific type of autocorrelation
598 function, even in cases with limited measurement data.

599 (4) The addition of location labels to measurement data addresses the challenge of predicting SOFs
600 in irregularly spaced CPT locations, and the developed SOF prediction models exhibit superior
601 performance, as indicated by R^2 , MAPE, and RMSE.

602 The proposed geotechnical knowledge-based data-driven framework is promising for geotechnical
603 engineering applications and bridges the gap between data-driven modeling and domain-specific
604 knowledge, thereby enhancing the accuracy and reliability of estimating spatially varying geotechnical
605 properties. It should be noted that: ① The proposed framework is flexible, where each model can be
606 replaced with common methods according to the user's preference. For example, the Maximum
607 Likelihood Estimation method can be used to estimate random field parameters. Then, the estimated
608 random field parameters can be used to calculate GCFs. Subsequently, user-friendly machine learning
609 models can be employed to establish the relationship between GCFs and the measured soil properties.
610 ② The construction process of the proposed random field parameter prediction model is relatively
611 complex. Once the model is trained, it can be directly applied to the target site, making the prediction
612 process highly efficient and straightforward. Moreover, it does not require users to have a background
613 in mathematical statistics. ③ Compared to the Kriging model, which has a cubic computational
614 complexity, the computational complexity of the data-driven subsurface model proposed in this study
615 is independent of the number of soil cells, resulting in stable modeling efficiency and excellent
616 scalability.

617 Notably, when subsurface conditions involve multiple soil layers with significant variations in

618 properties, and each layer has distinct random field parameters, researchers can employ either manual
619 or IC-XGBoost methods (Shi and Wang, 2023, 2021b) to delineate the spatial distribution of
620 subsurface stratigraphic boundaries. Then, our proposed method can be used to model the spatially
621 varying soil properties within each soil layer. Furthermore, due to the sparsity of CPT data in the
622 horizontal direction, the detection data for weak thin layers often constitute only a small proportion of
623 the overall dataset. Therefore, data-driven subsurface modeling methods should further investigate
624 their applicability and improvement strategies under conditions of data imbalance.

625 Notation

626 The following terms and notations are used in this paper:

<i>Terms</i>	
CPT	Cone penetration test
CRF	Conditional random field
BCS	Bayesian compressive sensing
IDW	Inverse distance weighting
GDF	Geotechnical distance field
GCF	Geotechnical correlation field
SOF	Scale of fluctuation
ACF	The autocorrelation function
CNN	Convolutional neural network
SNX	Single exponential
SOM	Second-order Markov
TOM	Third-order Markov
SQX	Squared exponential
BIN	Binary noise
CSX	Cosine exponential
SPH	Spherical
PCA-SCNN	Principal component analysis-shortcut connection neural network
Model #1	The ACFs classification model
Model #2	The horizontal SOF estimation model
Model #3	The vertical SOF estimation model
Model #4	The subsurface modeling model
MSE	Mean squared error
RMSE	Root mean square error
MAPE	Mean absolute percentage error
R^2	The coefficient of determination
<i>Notations</i>	
C	The correlation matrix
q_c	Cone tip resistance (MPa)
L	Horizontal coordinates (m)
P	The probability of measured CPT data belonging to each ACF type.
Y	The subsurface modeling outcome (MPa)

q_{cm}	The measured q_c (MPa)
\hat{q}_c	The standardized cone tip resistance
\hat{L}	The standardized horizontal coordinates
x	The input feature vector following PCA preprocessing
h	The hidden layer feature vector
θ	The weight vectors
β	The bias vectors
o	The output of the model
τ_h and τ_v	The horizontal and vertical distances of soil properties at two discrete points (m)
δ_h and δ_v	The horizontal and vertical SOFs (m)
λ	The power parameter of the Box-Cox method
e_i	The i -th soil cell
GCF_{ei}	The GCF of e_i
$\rho_{i,j}$	The correlation between soil cell e_i and e_j
f	The complex implicit function
N	Total number of soil cells
M	The number of samples in the training set, the feature dimension of the model #4
\bar{y}	The mean value
\hat{y}_i	The predicted values
y_i	The measured values
n	Total number of samples
h	The depth of the site
μ	The mean value of the random field (MPa)
σ	The standard deviation of the random field (MPa)
B_{crit}	The critical Bartlett peak value

627 Acknowledgments

628 This study was supported by the National Natural Science Foundation of China (Grant No.
629 52378330), the Natural Science Foundation of Sichuan Province (Grant No. 2023NSFSC0391), the
630 111 Project (Grant No. B21011), and the Postgraduate Research & Practice Innovation Program of
631 Jiangsu Province (Grant No. SJCX23_0084). The authors would like to extend their most sincere
632 gratitude to the Editors and Reviewers who provided help to improve this paper.

633 CRediT authorship contribution statement

634 **Weihang Chen:** Conceptualization, Data curation, Funding acquisition, Investigation, Methodology,
635 Software, Validation, Writing – original draft, Writing – review & editing. **Jianwen Ding:**
636 Conceptualization, Funding acquisition, Supervision, Writing – review & editing. **Chao Shi:**
637 Conceptualization, Investigation, Supervision, Writing – review & editing. **Tengfei Wang:**
638 Conceptualization, Investigation, Writing – review & editing. **David P. Connolly:** Conceptualization,

639 Investigation, Writing – review & editing.

640 **References**

- 641 Bai, L., Song, C., Zhou, X., Tian, Y., Wei, L., 2023. Assessing project portfolio risk via an enhanced
642 GA-BPNN combined with PCA. *Eng. Appl. Artif. Intell.* 126, 106779.
643 <https://doi.org/10.1016/j.engappai.2023.106779>
- 644 Cami, B., Javankhoshdel, S., Phoon, K.-K., Ching, J., 2020. Scale of Fluctuation for Spatially Varying
645 Soils: Estimation Methods and Values. *ASCE-ASME J. Risk Uncertain. Eng. Syst. Part Civ.*
646 *Eng.* 6, 03120002. <https://doi.org/10.1061/AJRUA6.0001083>
- 647 Chang, Y.-C., Ching, J., Phoon, K.-K., Yue, Q., 2021. On the Hole Effect in Soil Spatial Variability.
648 *ASCE-ASME J. Risk Uncertain. Eng. Syst. Part Civ. Eng.* 7, 04021039.
649 <https://doi.org/10.1061/AJRUA6.0001168>
- 650 Chen, W., Ding, J., Wang, T., Connolly, D.P., Wan, X., 2023. Soil property recovery from incomplete
651 in-situ geotechnical test data using a hybrid deep generative framework. *Eng. Geol.* 326,
652 107332. <https://doi.org/10.1016/j.enggeo.2023.107332>
- 653 Ching, J., Huang, W.-H., Phoon, K.-K., 2020. 3D Probabilistic Site Characterization by Sparse
654 Bayesian Learning. *J. Eng. Mech.* 146, 04020134. [https://doi.org/10.1061/\(ASCE\)EM.1943-7889.0001859](https://doi.org/10.1061/(ASCE)EM.1943-7889.0001859)
- 655
- 656 Ching, J., Phoon, K.-K., 2019. Constructing Site-Specific Multivariate Probability Distribution Model
657 Using Bayesian Machine Learning. *J. Eng. Mech.* 145, 04018126.
658 [https://doi.org/10.1061/\(ASCE\)EM.1943-7889.0001537](https://doi.org/10.1061/(ASCE)EM.1943-7889.0001537)
- 659 Ching, J., Phoon, K.-K., Stuedlein, A.W., Jaksa, M., 2019. Identification of sample path smoothness
660 in soil spatial variability. *Struct. Saf.* 81, 101870.
661 <https://doi.org/10.1016/j.strusafe.2019.101870>
- 662 Ching, J., Wu, T.-J., Stuedlein, A.W., Bong, T., 2018. Estimating horizontal scale of fluctuation with
663 limited CPT soundings. *Geosci. Front., Reliability Analysis of Geotechnical Infrastructures* 9,
664 1597–1608. <https://doi.org/10.1016/j.gsf.2017.11.008>
- 665 Collico, S., Arroyo, M., Devincenzi, M., 2024. A simple approach to probabilistic CPTu-based
666 geotechnical stratigraphic profiling. *Comput. Geotech.* 165, 105905.
667 <https://doi.org/10.1016/j.compgeo.2023.105905>
- 668 Dasaka, S. m., Zhang, L. m., 2012. Spatial variability of in situ weathered soil. *Géotechnique* 62, 375–
669 384. <https://doi.org/10.1680/geot.8.P.151.3786>
- 670 Gu, X., Wang, L., Ou, Q., Zhang, W., 2023. Efficient stochastic analysis of unsaturated slopes
671 subjected to various rainfall intensities and patterns. *Geosci. Front.* 14, 101490.
672 <https://doi.org/10.1016/j.gsf.2022.101490>
- 673 Guan, Z., Wang, Y., Cao, Z., Hong, Y., 2020. Smart sampling strategy for investigating spatial
674 distribution of subsurface shallow gas pressure in Hangzhou Bay area of China. *Eng. Geol.*
675 274, 105711. <https://doi.org/10.1016/j.enggeo.2020.105711>
- 676 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep Residual Learning for Image Recognition. Presented
677 at the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–
678 778. <https://doi.org/10.1109/CVPR.2016.90>

- 679 Hong, L., Wang, X., Zhang, W., Li, Y., Zhang, R., Chen, C., 2024. System reliability-based robust
680 design of deep foundation pit considering multiple failure modes. *Geosci. Front.* 15, 101761.
681 <https://doi.org/10.1016/j.gsf.2023.101761>
- 682 Hu, Y., Wang, Y., 2024. Evaluating statistical homogeneity of cone penetration test (CPT) data profile
683 using auto-correlation function. *Comput. Geotech.* 165, 105852.
684 <https://doi.org/10.1016/j.compgeo.2023.105852>
- 685 Li, H., Xu, Z., Taylor, G., Studer, C., Goldstein, T., 2018. Visualizing the Loss Landscape of Neural
686 Nets.
- 687 Liu, W.F., Leung, Y.F., 2018. Characterising three-dimensional anisotropic spatial correlation of soil
688 properties through in situ test results. *Géotechnique* 68, 805–819.
689 <https://doi.org/10.1680/jgeot.16.P.336>
- 690 Liu, W.F., Leung, Y.F., Lo, M.K., 2017. Integrated framework for characterization of spatial variability
691 of geological profiles. *Can. Geotech. J.* 54, 47–58. <https://doi.org/10.1139/cgj-2016-0189>
- 692 Lloret-Cabot, M., Fenton, G.A., Hicks, M.A., 2014. On the estimation of scale of fluctuation in
693 geostatistics. *Georisk Assess. Manag. Risk Eng. Syst. Geohazards* 8, 129–140.
694 <https://doi.org/10.1080/17499518.2013.871189>
- 695 Nag, P., Sun, Y., Reich, B.J., 2023. Spatio-temporal DeepKriging for interpolation and probabilistic
696 forecasting. *Spat. Stat.* 57, 100773. <https://doi.org/10.1016/j.spasta.2023.100773>
- 697 NZGD, 2023. World Wide Web Address [WWW Document]. URL <https://www.nzgd.org.nz> (accessed
698 10.18.23).
- 699 Onyejekwe, S., Kang, X., Ge, L., 2016. Evaluation of the scale of fluctuation of geotechnical
700 parameters by autocorrelation function and semivariogram function. *Eng. Geol.* 214, 43–49.
701 <https://doi.org/10.1016/j.enggeo.2016.09.014>
- 702 Phoon, K.-K., Cao, Z.-J., Ji, J., Leung, Y.F., Najjar, S., Shuku, T., Tang, C., Yin, Z.-Y., Ikumasa, Y.,
703 Ching, J., 2022. Geotechnical uncertainty, modeling, and decision making. *Soils Found.* 62,
704 101189. <https://doi.org/10.1016/j.sandf.2022.101189>
- 705 Phoon, K.-K., Quek, S.-T., An, P., 2003. Identification of Statistically Homogeneous Soil Layers Using
706 Modified Bartlett Statistics. *J. Geotech. Geoenvironmental Eng.* 129, 649–659.
707 [https://doi.org/10.1061/\(ASCE\)1090-0241\(2003\)129:7\(649\)](https://doi.org/10.1061/(ASCE)1090-0241(2003)129:7(649))
- 708 Phoon, K.-K., Wang, Y., 2019. Chicken (method) and egg (data)— Which comes first? Presented at
709 the Int. Symp. on Reliability of Multidisciplinary Engineering Systems under Uncertainty
710 (ISMES2019), Ministry of Education and Ministry of Science and Technology, Da’an, Taipei.
- 711 Qi, X.-H., Liu, H.-X., 2019. Estimation of autocorrelation distances for in-situ geotechnical properties
712 using limited data. *Struct. Saf.* 79, 26–38. <https://doi.org/10.1016/j.strusafe.2019.02.003>
- 713 Shi, C., Wang, Y., 2023. Data-driven spatio-temporal analysis of consolidation for rapid reclamation.
714 *Géotechnique* 1–21. <https://doi.org/10.1680/jgeot.22.00016>
- 715 Shi, C., Wang, Y., 2021a. Non-parametric machine learning methods for interpolation of spatially
716 varying non-stationary and non-Gaussian geotechnical properties. *Geosci. Front.* 12, 339–350.
717 <https://doi.org/10.1016/j.gsf.2020.01.011>
- 718 Shi, C., Wang, Y., 2021b. Development of Subsurface Geological Cross-Section from Limited Site-
719 Specific Boreholes and Prior Geological Knowledge Using Iterative Convolution XGBoost. *J.*
720 *Geotech. Geoenvironmental Eng.* 147, 04021082. [42](https://doi.org/10.1061/(ASCE)GT.1943-</p></div><div data-bbox=)

721 5606.0002583

722 Uzielli, M., Vannucchi, G., Phoon, K.K., 2005. Random field characterisation of stress-normalised cone
723 penetration testing parameters. *Géotechnique* 55, 3–20.
724 <https://doi.org/10.1680/geot.2005.55.1.3>

725 Vanmarcke, E.H., 1983. *Random Fields: Analysis and Synthesis*. MIT Press.

726 Wang, F., Chen, J.E., 2023. Efficient modeling of random fields by using Gaussian process inducing-
727 point approximations. *Comput. Geotech.* 157, 105304.
728 <https://doi.org/10.1016/j.compgeo.2023.105304>

729 Wang, H., Wang, X., Wellmann, J.F., Liang, R.Y., 2018. Bayesian Stochastic Soil Modeling
730 Framework Using Gaussian Markov Random Fields. *ASCE-ASME J. Risk Uncertain. Eng.*
731 *Syst. Part Civ. Eng.* 4, 04018014. <https://doi.org/10.1061/AJRUA6.0000965>

732 Wang, Y., Hu, Y., Zhao, T., 2020. Cone penetration test (CPT)-based subsurface soil classification and
733 zonation in two-dimensional vertical cross section using Bayesian compressive sampling. *Can.*
734 *Geotech. J.* 57, 947–958. <https://doi.org/10.1139/cgj-2019-0131>

735 Wang, Y., Shi, C., 2023. Data-driven analysis of soil consolidation with prefabricated vertical drains
736 considering stratigraphic variation. *Comput. Geotech.* 161, 105569.
737 <https://doi.org/10.1016/j.compgeo.2023.105569>

738 Wang, Y., Zhao, T., 2017. Statistical interpretation of soil property profiles from sparse data using
739 Bayesian compressive sampling. *Géotechnique* 67, 523–536.
740 <https://doi.org/10.1680/jgeot.16.P.143>

741 Xiao, T., Li, D.-Q., Cao, Z.-J., Au, S.-K., Phoon, K.-K., 2016. Three-dimensional slope reliability and
742 risk assessment using auxiliary random finite element method. *Comput. Geotech.* 79, 146–158.
743 <https://doi.org/10.1016/j.compgeo.2016.05.024>

744 Xiao, T., Li, D.-Q., Cao, Z.-J., Zhang, L.-M., 2018. CPT-Based Probabilistic Characterization of Three-
745 Dimensional Spatial Variability Using MLE. *J. Geotech. Geoenvironmental Eng.* 144,
746 04018023. [https://doi.org/10.1061/\(ASCE\)GT.1943-5606.0001875](https://doi.org/10.1061/(ASCE)GT.1943-5606.0001875)

747 Xie, J., Huang, J., Lu, J., Burton, G.J., Zeng, C., Wang, Y., 2022a. Development of two-dimensional
748 ground models by combining geotechnical and geophysical data. *Eng. Geol.* 300, 106579.
749 <https://doi.org/10.1016/j.enggeo.2022.106579>

750 Xie, J., Huang, J., Zeng, C., Huang, S., Burton, G.J., 2022b. A generic framework for geotechnical
751 subsurface modeling with machine learning. *J. Rock Mech. Geotech. Eng.* 14, 1366–1379.
752 <https://doi.org/10.1016/j.jrmge.2022.08.001>

753 Xie, J., Zeng, C., Huang, J., Zhang, Y., Lu, J., 2024. A back analysis scheme for refined soil
754 stratification based on integrating borehole and CPT data. *Geosci. Front.* 15, 101688.
755 <https://doi.org/10.1016/j.gsf.2023.101688>

756 Yan, W., Shen, P., Zhou, W.-H., Ma, G., 2023. A rigorous random field-based framework for 3D
757 stratigraphic uncertainty modelling. *Eng. Geol.* 323, 107235.
758 <https://doi.org/10.1016/j.enggeo.2023.107235>

759 Yang, Y., Wang, P., Brandenberg, S.J., 2022. An algorithm for generating spatially correlated random
760 fields using Cholesky decomposition and ordinary kriging. *Comput. Geotech.* 147, 104783.
761 <https://doi.org/10.1016/j.compgeo.2022.104783>

762 Yang, Z., Ching, J., 2021. Simulation of three-dimensional random field conditioning on incomplete

763 site data. *Eng. Geol.* 281, 105987. <https://doi.org/10.1016/j.enggeo.2020.105987>

764 Yang, Z., Li, X., Qi, X., 2022. Efficient simulation of multivariate three-dimensional cross-correlated
765 random fields conditioning on non-lattice measurement data. *Comput. Methods Appl. Mech.*
766 *Eng.* 388, 114208. <https://doi.org/10.1016/j.cma.2021.114208>

767 Yang, Z., Nie, J., Peng, X., Tang, D., Li, X., 2021. Effect of random field element size on reliability
768 and risk assessment of soil slopes. *Bull. Eng. Geol. Environ.* 80, 7423–7439.
769 <https://doi.org/10.1007/s10064-021-02422-z>

770 Zhang, J.-Z., Phoon, K.K., Zhang, D.-M., Huang, H.-W., Tang, C., 2021. Novel approach to estimate
771 vertical scale of fluctuation based on CPT data using convolutional neural networks. *Eng. Geol.*
772 294, 106342. <https://doi.org/10.1016/j.enggeo.2021.106342>

773 Zhang, J.-Z., Zhang, D.-M., Huang, H.-W., Phoon, K.K., Tang, C., Li, G., 2022. Hybrid machine
774 learning model with random field and limited CPT data to quantify horizontal scale of
775 fluctuation of soil spatial variability. *Acta Geotech.* 17, 1129–1145.
776 <https://doi.org/10.1007/s11440-021-01360-0>

777 Zhang, P., Yin, Z.-Y., Jin, Y.-F., 2022. Bayesian neural network-based uncertainty modelling:
778 application to soil compressibility and undrained shear strength prediction. *Can. Geotech. J.*
779 59, 546–557. <https://doi.org/10.1139/cgj-2020-0751>

780 Zhang, W., Gu, X., Tang, L., Yin, Y., Liu, D., Zhang, Y., 2022. Application of machine learning, deep
781 learning and optimization algorithms in geoenvironment and geoscience: Comprehensive
782 review and future challenge. *Gondwana Res.* 109, 1–17.
783 <https://doi.org/10.1016/j.gr.2022.03.015>

784 Zhang, W., Wu, C., Tang, L., Gu, X., Wang, L., 2023. Efficient time-variant reliability analysis of
785 Bazimen landslide in the Three Gorges Reservoir Area using XGBoost and LightGBM
786 algorithms. *Gondwana Res., Data driven models* 123, 41–53.
787 <https://doi.org/10.1016/j.gr.2022.10.004>

788 Zhang, W., Zhang, R., Wu, C., Goh, A.T.C., Lacasse, S., Liu, Z., Liu, H., 2020. State-of-the-art review
789 of soft computing applications in underground excavations. *Geosci. Front.* 11, 1095–1106.
790 <https://doi.org/10.1016/j.gsf.2019.12.003>

791 Zhao, T., Hu, Y., Wang, Y., 2018. Statistical interpretation of spatially varying 2D geo-data from sparse
792 measurements using Bayesian compressive sampling. *Eng. Geol.* 246, 162–175.
793 <https://doi.org/10.1016/j.enggeo.2018.09.022>

794 Zhao, T., Xu, L., Wang, Y., 2020. Fast non-parametric simulation of 2D multi-layer cone penetration
795 test (CPT) data without pre-stratification using Markov Chain Monte Carlo simulation. *Eng.*
796 *Geol.* 273, 105670. <https://doi.org/10.1016/j.enggeo.2020.105670>

797 Zou, H., Liu, S., Cai, G., Bheemasetti, T.V., Puppala, A.J., 2017. Mapping probability of liquefaction
798 using geostatistics and first order reliability method based on CPTU measurements. *Eng. Geol.*
799 218, 197–212. <https://doi.org/10.1016/j.enggeo.2017.01.021>

800