# Similarity quantification of soil spatial variability between two cross-sections using auto-correlation functions

Yue Hu[1], Yu Wang[2,*], Kok-Kwang Phoon[3], and Michael Beer[4,5,6]

[1] Research Fellow, Institute for Risk and Reliability, Leibniz Universität Hannover, Hannover, Germany. (Email): yue.hu@irz.uni-hannover.de

[2] Professor, Department of Architecture and Civil Engineering, City University of Hong Kong, Kowloon, Hong Kong SAR, China. (Tel): 852-3442-7605 (Fax): 852-3442-0427 (Email): yuwang@cityu.edu.hk ([*]Corresponding author)

[3] Cheng Tsang Man Chair Professor and Provost, Singapore University of Technology and Design, Singapore. (Email): kkphoon@sutd.edu.sg

[4] Professor and Head, Institute for Risk and Reliability, Leibniz Universität Hannover, Hannover, Germany. (Email): beer@irz.uni-hannover.de

[5] Professor, Institute for Risk and Uncertainty, University of Liverpool, United Kingdom.

[6] Guest Professor, International Joint Research Center for Resilient Infrastructure & International Joint Research Center for Engineering Reliability and Stochastic Mechanics, Tongji University, Shanghai, China.

**Abstract**

In geotechnical engineering, an appreciation of local geological conditions from similar sites is beneficial and can support informed decision-making during site characterization. This practice is known as "site recognition", which necessitates a rational quantification of site similarity. This paper proposes a data-driven method to quantify the similarity between two cross-sections based on the spatial variability of one soil property from a spectral perspective. Bayesian compressive sensing (BCS) is first used to obtain the discrete cosine transform (DCT) spectrum for a cross-section. Then DCT-based auto-correlation function (ACF) is calculated based on the obtained DCT spectrum using a set of newly derived ACF calculation equations. The cross-sectional similarity is subsequently reformulated as the cosine similarity of DCT-based ACFs between cross-sections. In contrast to the existing methods, the proposed method explicitly takes soil property spatial variability into account in an innovative way. The challenges of sparse investigation data, non-stationary and anisotropic spatial variability, and inconsistent spatial dimensions of different cross-sections are tackled effectively. Both numerical examples and real data examples from New Zealand are provided for illustration. Results show that the proposed method can rationally quantify cross-sectional similarity and associated statistical uncertainty from sparse investigation data. The proposed method advances data-driven site characterization, a core application area in data-centric geotechnics.

Keywords: Geotechnical site investigation; Site similarity; Auto-correlation; Bayesian compressive sensing

## 1. Introduction

Reliable site characterization is a cornerstone of effective geotechnical designs and construction safety. However, it is often subject to significant uncertainty due to the spatially variable geological conditions and a sparsity of investigation points (e.g., boreholes, in-situ tests). In practice, to supplement limited knowledge at a target site and to mitigate the resultant uncertainty, engineers usually attempt to refer to and review available information (e.g., interpreted soil cross-sections) of previous construction sites in the neighborhood where geological conditions are expected to be similar to the target project site. This practice is also known as "site recognition" which helps engineers to better understand the site-specific features at the target site and plays an important role in data-driven site characterization and informed decision-making in the presence of inter-site variabilities (e.g., Fenton, 1999b; Phoon et al., 2022; Yang et al., 2022; Phoon and Zhang, 2023; Shi et al., 2023; Zhao et al., 2023).

Consider, for example, a two-dimensional (2D) soil property cross-section, which has been commonly adopted in practical engineering designs and analyses for explicitly representing site geological conditions along both depth and horizontal directions. Interpretation of a target 2D cross-section may be underpinned and supplemented by referring to 2D cross-section interpretations available from other pre-existing and documented construction sites. To this end, before introducing knowledge from other 2D cross-sections to inform decision making at a target site, it is desirable to assess the similarity between the target 2D cross-section and the other available 2D cross-sections. Geotechnical engineers primarily do this *qualitatively*, because there are no quantitative methods that are tractable/effective in the presence of sparse and incomplete data to name a few data attributes. From a geotechnical engineering viewpoint, 2D cross-sectional similarity shall be closely related to the similarity of corresponding 2D soil property spatial variability, which is a natural product of complicated geological formation processes (e.g., erosion, weathering, deposition) undergone by the

corresponding sites (e.g., Fenton, 1999a; Phoon and Kulhawy, 1999; Baecher and Christian, 2003; Juang et al., 2019; Wang et al., 2022). However, a direct and quantitative comparison of spatial variability in different 2D cross-sections is challenging because of the following issues: (1) the available site investigation data (e.g., borehole and in-situ test data) in both a target 2D cross-section and existing 2D cross-sections are usually sparse and are often not measured over an identical sampling grid (e.g., Xu et al., 2021; Guan and Wang, 2023). It is unlikely that the site investigation plans for two sites are identical (e.g., boreholes or CPT soundings layout). Therefore, classical statistical correlation analysis may not be applicable to quantify the similarity between two such cross-sections; (2) spatial variability in 2D cross-sections may exhibit non-stationary trends and spatial variability anisotropy. Accurate identification of the underlying trend and spatial variability anisotropy for different cross-sections is critical for cross-sectional similarity quantification, but challenging in the presence of sparse data (e.g., Ching and Phoon, 2017; Ching et al., 2017; Hu et al., 2019; Wang et al., 2019; Ching et al., 2020; Shuku et al., 2020; Yoshida et al., 2021; Ching et al., 2022; Katsman and Painuly, 2022); and (3) the dimensions of different 2D cross-sections along depth and horizontal directions are often different due to projects occupying different footprints and extending to different depths. Directly comparing 2D cross-sections with different spatial dimensions is often a tricky task (e.g., Shechtman and Irani, 2007; Simakov et al., 2008; Shi and Wang, 2021b). Therefore, how to quantitatively evaluate the similarity between two given 2D cross-sections from their sparse site investigation data measured over different grids and covering different spatial dimensions remains unsolved.

Recently, the topic of site similarity, or site retrieval, has been investigated from different perspectives. For example, Ching and Phoon (2020) proposed a Bayesian method for measuring similarity between data records (e.g., two or more soil parameter values at a location and or depth) at a target site and data records from other sites. Sharma et al. (2022) further

2

developed a novel hierarchical Bayesian model for measuring similarity between the target site

and database sites, achieving site similarity quantification beyond solely data record similarity.

Phoon and Ching (2022) presented a summary of different methods for similarity measures.

Their frameworks focused on the MUSIC data attributes framework (Multivariate, Uncertain

and Unique, Sparse, Incomplete, and potentially Corrupted) (e.g., Phoon et al., 2022) and

treated the likelihood function of past data records given site-specific data records as an index

of the similarity. Only cross correlations between different soil parameters are considered. The

spatial variability of a soil parameter was not considered in these studies, although it was

recognized as a critical aspect. In addition, Han et al. (2022) used confidence ellipses to

quantify the similarity of soil parametric data using existing databases. Their framework

required abundant data over identical depth ranges to be compiled at every site. The

performance was highly dependent on the specific volume of available data at different depths.

More importantly, the geotechnical spatial variability might not be fully preserved after

preprocessing of the data. Shi and Wang (2021a, 2021b) proposed to use training images to

incorporate and summarize past geological knowledge on stratigraphy and to quantify the site

similarity by measuring the similarity of edge orientation statistics of soil layer boundaries

between site-specific borehole data and geological training images. The spatial variability of

soil property was not considered. Currently, there is no rational method available for

quantifying similarity between 2D cross-sections of soil property from sparse site investigation

data with explicit consideration of spatial variability.

       This paper proposes a novel method for data-driven quantification of 2D cross-sectional

similarity from a spectral perspective. This paper attempts to fill an important gap in the site

recognition challenge (e.g., Phoon et al., 2022). A non-parametric method called Bayesian

compressive sensing (BCS) is used to directly approximate the sparse spectrum of 2D cross-

section. A new efficient and robust formulation of 2D auto-correlation function (ACF) is

115    derived for a unified representation of 2D spatial variability based on a sparse spectrum

116    approximated by BCS. The 2D cross-sectional similarity is then quantified by the similarity

117    between ACFs of the corresponding 2D cross-sections. In contrast to current methods in the

118    literature, cross-sectional similarity quantification in this study deals explicitly with soil

119    property spatial variability. Theoretical derivation suggests that the spatial variability patterns

120    in a 2D cross-section, either stationary or non-stationary, spatially isotropic or anisotropic, can

121    be quantified concisely by 2D ACF. The three challenges highlighted above are solved by the

122    proposed method. The proposed method also has significant practical relevance in geotechnical

123    site recognition. For example, given a global geotechnical database (e.g., Ching et al., 2023)

124    containing a wealth of information from different sites, the proposed method can efficiently

125    pick up a limited number of similar and informative records for a target site, which is also

126    referred to as a "quasi-regional clustering" strategy (e.g., Phoon and Ching, 2022; Guan et al.,

127    2023b).

128          The rest of this paper is organized as follows. Section 2 briefly illustrates the

129    background and practical significance of 2D cross-sectional similarity using real examples.

130    The proposed method for data-driven quantification of 2D cross-sectional similarity from

131    sparse site investigation data is described in Section 3. The implementation procedure is

132    provided in Section 4, followed by illustrative examples in Section 5. In Section 6, a real case

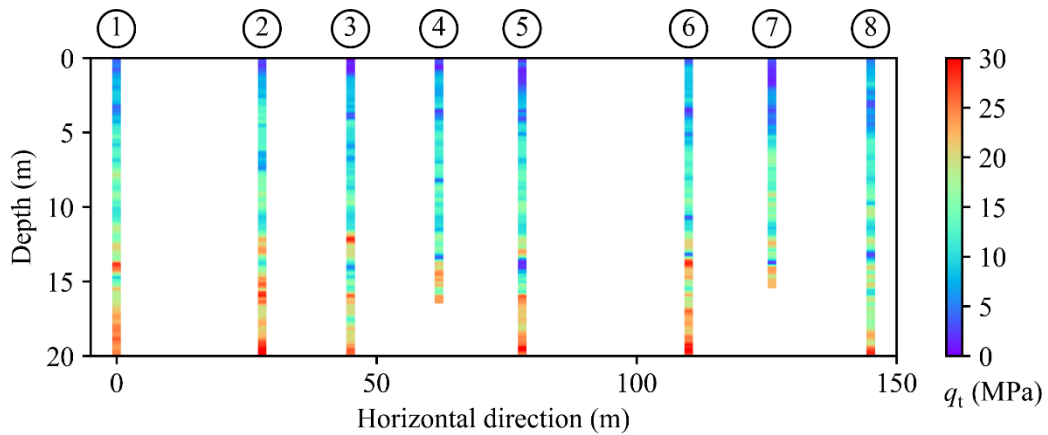133    study is used to demonstrate the application of the proposed method.

134

Figure 1. A layout of cone penetration tests (CPTs) performed in two cross-sections in

Christchurch, New Zealand (NZGD, 2023)

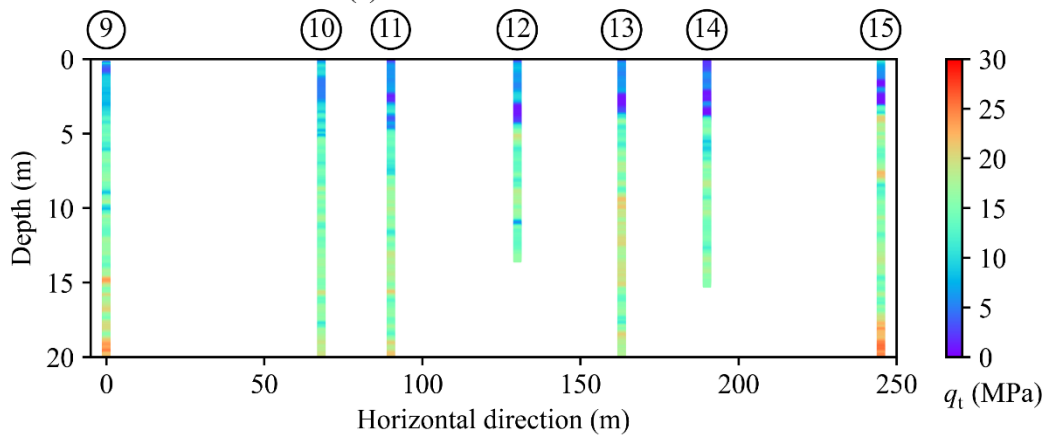## 2. Similarity between two 2D cross-sections of a soil property

A 2D cross-section in this study refers to 2D spatial variability of one soil property within a

single soil layer. To illustrate the 2D soil property cross-sectional similarity, Figure 1 shows a

map with a layout of 15 cone penetration tests (CPTs) performed in Christchurch, New Zealand.

The CPTs data are obtained from the New Zealand Geotechnical Database (e.g., NZGD, 2023).

In Figure 1, the CPTs are denoted by yellow triangles and numbered from #1 to #15. It is seen

that these CPTs were performed at two separate sites, i.e., Site 1 and Site 2, respectively. Eight

CPTs (CPT #1 to CPT #8) were performed in Site 1 (see the left-hand side of Figure 1), while

seven CPTs (CPT #9 to CPT #15) were carried out in Site 2 (see the right-hand side of Figure

1). The actual IDs of these CPTs used in the NZGD database are also provided in the map.

Note that the CPTs at these two sites are roughly laid alone a straight line, leading to two cross-

sections denoted by two red dashed lines. Engineers may assess the cross-sections at Site 1 and

Site 2 to be similar since the distance between the two cross-sections is only roughly 700m.

Figures 2a and 2b show the corrected cone resistance ($q_t$) data of available CPTs by color-

coded columns in the two cross-sections, respectively. It appears that the general patterns of $q_t$
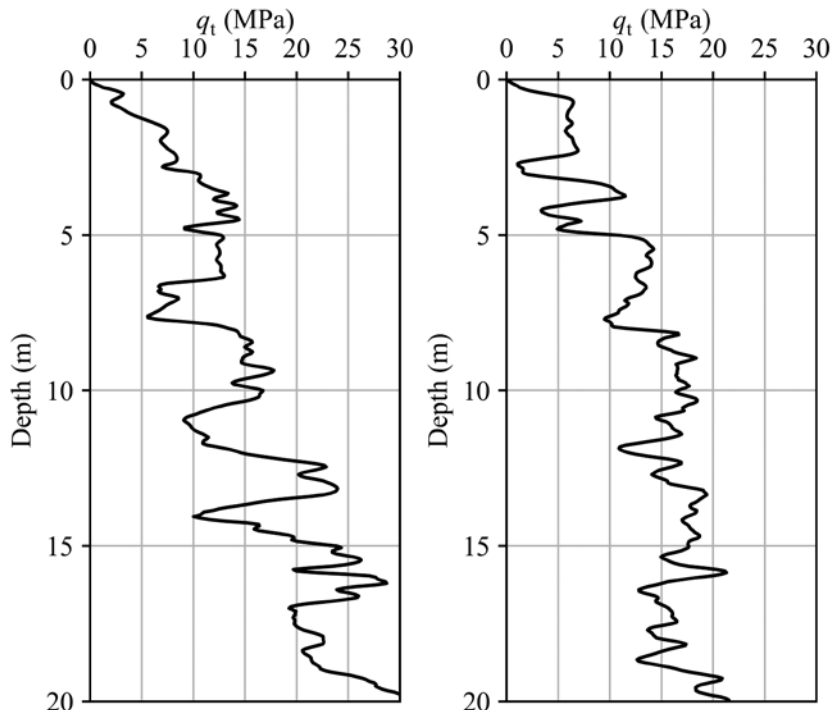
5

152    data at Sites 1 and 2 are comparable, both with $q_t$ values fluctuating but generally increasing

153    with depth. For example, as shown in Figures 2c and 2d, the $q_t$ profiles of CPT #2 (CPT_50725)

154    from Site 1 and CPT #11 (CPT_35561) from Site 2 exhibit comparable variation patterns.

155    Between these two specific 1D $q_t$ profiles, similarity quantification can be conducted directly

156    and readily. Mathematically, this may be routinely achieved by calculating the cross-

157    correlation between these two 1D $q_t$ profiles after re-configuring the data with an identical

158    sampling interval, or by comparing the corresponding estimated spectrum of these two $q_t$

159    profiles (e.g., Priestley, 1981; Dai et al., 2022; Guan and Wang, 2023). However, it is very

160    challenging to quantify the $q_t$ data cross-sectional similarity between Site 1 and Site 2, as shown

161    in Figures 2a and 2b. This real example clearly demonstrates the three challenges for direct

162    quantification of 2D cross-sectional similarity mentioned above in concrete terms. First, in

163    Figures 2a and 2b, the available CPT soundings are sparse within these two cross-sections with

164    non-uniform horizontal spacing and sounding depths. Second, the $q_t$ data shown in these two

165    cross-sections exhibit evidence of non-stationarity and spatial variability anisotropy. Third,

166    these two cross-sections have different spatial dimensions. The cross-section of Site 1 has a

167    length of around 145m, while the cross-section of Site 2 has a length of 245m. The challenges

168    highlighted above cannot be addressed by existing methods in literature (e.g., Ching and Phoon,

169    2020; Han et al., 2022). This next section addresses these challenges by proposing a novel data-

170    driven approach.

171

172
173　　Figure 2. Illustration of CPT example: (a) Cross-section of Site 1; (b) Cross-section of Site 2;
174　　　　(c) $q_t$ data profile of CPT #2 from Site 1; and (d) $q_t$ data profile of CPT #11 from Site 2
175

7

## 3. Proposed method for quantification of cross-sectional similarity

The concept for quantification of cross-sectional similarity is to treat soil property cross-sections as images and then compare them from a spectral perspective. Note that image spectral analysis is able to identify non-stationary patterns, spatial variability anisotropy, and spatial shift-invariant patterns (e.g., Shalvi and Weinstein, 1996; Wen and Gu, 2004; Blumensath and Davies, 2006). The results of image spectral analysis are also independent of image dimension. In the proposed method, the following two steps shall be performed before similarity quantification. First, Bayesian compressive sensing (BCS) is adopted to obtain the discrete cosine transform (DCT) spectrum of a 2D soil property cross-section directly from sparse data. It has been shown in past research that BCS can deal with non-stationarity (e.g., Wang et al., 2019; Zhao and Wang, 2020), spatial variability anisotropy of soil property (e.g., Hu et al., 2019), and the associated statistical uncertainty quantification (e.g., Wang et al., 2022). Second, to tackle the difficulty in comparing target 2D cross-sections with different dimensions, a novel and efficient 2D DCT-based ACF is developed to facilitate a unified representation of 2D soil property spatial variability. The DCT-based ACF is utilized in this study as a data-driven surrogate to represent 2D cross-sections and enables direct pattern comparison between cross-sections with different spatial dimensions. Subsequently, cross-sectional similarity is quantified by DCT-based ACF similarity between two cross-sections. Details of the proposed method are elaborated in the following three subsections. The approximation of sparse DCT spectrum by BCS will be described in Subsection 3.1. A unified representation of 2D spatial variability using DCT-based ACF is then derived in Subsection 3.2. Quantification of DCT-based ACF similarity is established in Subsection 3.3.

### 3.1. Approximation of DCT spectrum from sparse data using BCS

Compressive sensing (CS) is a technique for efficiently acquiring and reconstructing signals or images (e.g., Candès et al., 2006; Donoho, 2006; Candès and Wakin, 2008). Utilizing the

201    sparsity featured by many signals or images after adopting appropriate basis functions, CS is

202    able to reconstruct a signal or image from far fewer measurement data points than the number

203    indicated by conventional Nyquist sampling theorem (e.g., Shannon, 1948; Candès et al., 2006).

204    From a spectral perspective, complicated soil property spatial variability in terms of a 1D

205    profile or 2D cross-section (or image) can be sparsely represented after transformation using

206    basis functions. For example, the DCT functions, which have been widely used in digital signal

207    processing and data compression (e.g., Rao and Yip, 1990; Wallace, 1992), are used to

208    construct basis functions in this study. The commonly used type-II 1D DCT basis function is

209    defined as:

$$
B_t(x) = \begin{cases} \dfrac{1}{\sqrt{N}} & for\ t = 1;\ x = 1,2,\cdots,N \\[4mm] \sqrt{\dfrac{2}{N}}\ cos\pi\dfrac{(t-1)(2x-1)}{2N} & for\ t = 2,\cdots,N;\ x = 1,2,\cdots,N \end{cases} \tag{1}
$$

210    in which $x$ represents the 1D index ($x$=1, 2, …, $N$); $t$ indicates the order of $B_t(x)$. In Figure 3,

211    the first five DCT basis functions (i.e., $t$ = 1, 2, 3, 4, 5) with $N$ = 200 are illustrated by colored

212    lines with different styles. The frequency of these DCT basis function $B_t(x)$ is controlled by $t$

213    and increases with $t$. Based on the 1D DCT basis functions in Equation (1), 2D DCT basis

214    functions may be constructed by a tensor product of two 1D DCT basis functions (e.g., Itskov,
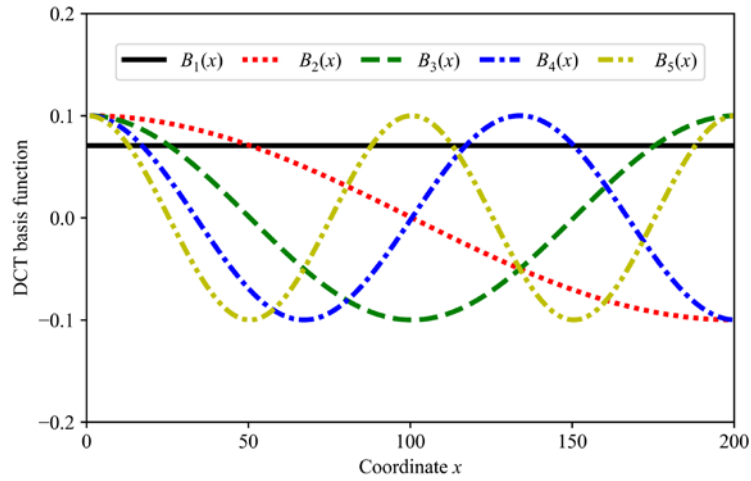
215    2007):

$$
B_{t,s}(x_1, x_2) = B_t(x_1) \times B_s(x_2) \tag{2}
$$

216    in which $B_t(x_1)$ and $B_s(x_2)$ are two basis functions along two directions, respectively; $t$ and $s$

217    indicate the corresponding orders ($t$=1, 2, …, $N_1$; $s$=1, 2, …, $N_2$). For example, Figure 4

218    illustrates the construction process of 25 2D DCT basis functions $B_{t,s}(x_1, x_2)$. Each 2D DCT

219    basis function is constructed by a tensor product of two 1D DCT basis functions of the same

220    length at each frequency (e.g., $t$, $s$=1, 2, 3, 4, 5). Using the 2D DCT basis function $B_{t,s}(x_1, x_2)$,
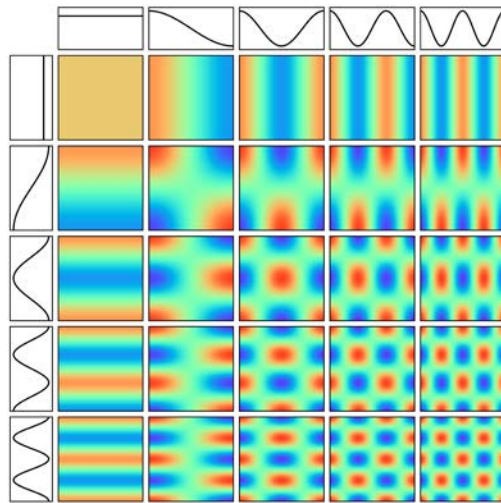
221 soil property spatial variability in a cross-section can be regarded as an image F with size

222 $N_1 \times N_2$, which is formulated as (e.g., Tipping, 2001; Candès and Wakin, 2008):

$$F(x_1, x_2) = \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D} B_{t,s}(x_1, x_2) \qquad (3)$$

223 in which $F(x_1, x_2)$ is the 2D spatial variability in the cross-section; $x_1$, $x_2$ are indexes along

224 two directions, respectively ($x_1$=1, 2, …, $N_1$; $x_2$=1, 2, …, $N_2$); $\omega_{t,s}^{2D}$ is the weight coefficient of

225 $B_{t,s}(x_1, x_2)$. The weight coefficients and their corresponding frequencies collectively form the

226 DCT spectrum of $F(x_1, x_2)$.



227

228 Figure 3. Illustration of five 1D discrete cosine transform (DCT) basis functions



229

230 Figure 4. Construction of 2D DCT basis functions from 1D DCT basis functions

231         The DCT spectrum enables an effective representation of variability patterns at various

232     frequencies along two directions. Note that both non-stationarity and spatial variability

233     anisotropy may be preserved by a combination of various $B_{t,s}(x_1, x_2)$ with large weight

234     coefficients and specific patterns (e.g., see various 2D DCT basis functions in Figure 4). In

235     addition, due to the spatial correlation contained in soil property spatial variability, $F(x_1, x_2)$,

236     usually leads to a sparse representation in its DCT spectrum. In other words, most weight

237     coefficients in Equation (3) have negligible magnitudes and only limited weight coefficients

238     are significant, or non-trivial, after adopting proper basis functions (e.g., Candès and Wakin,

239     2008; Zhao et al., 2018; Hu et al., 2019). It is therefore feasible to approximate those limited

240     non-trivial weight coefficients in DCT spectrum using sparse data Y, which is expressed as:

$$Y(x_1, x_2) = \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D} A_{t,s}(x_1, x_2) \qquad (4)$$

241     in which $Y(x_1, x_2)$ is the measured data points from $F(x_1, x_2)$, and $A_{t,s}(x_1, x_2)$ are

242     corresponding values extracted from the 2D basis function $B_{t,s}(x_1, x_2)$ based on the

243     measurement locations. In the context of geotechnical site investigation, generally the number,

244     $M$, of available measurement $Y(x_1, x_2)$, is much smaller than the size of $F(x_1, x_2)$ (i.e.,

245     $N_1 \times N_2$), and this leads to an underdetermined system in Equation (4), which cannot be solved

246     directly. Numerical algorithms, e.g., orthogonal matching pursuit (OMP), may be used to

247     approximate the solution $\hat{\omega}_{t,s}^{2D}$ in Equation (4) by minimizing the error between the measured

248     data $Y(x_1, x_2)$ and the estimated values at measured locations (e.g., Pati et al., 1993; Wang and

249     Zhao, 2016). The idea of the OMP algorithm is to iteratively find out $B_{t,s}(x_1, x_2)$ that can well

250     match the $Y(x_1, x_2)$. However, since the available site investigation data $Y(x_1, x_2)$ are sparse,

251     CS may not produce a perfect reconstruction of the spatial variability $F(x_1, x_2)$, and the

252     approximated $\hat{\omega}_{t,s}^{2D}$ contain significant statistical uncertainty. To quantify the associated

253 statistical uncertainty, CS may be integrated with the Bayesian framework (i.e., Bayesian

254 compressive sensing, BCS) to estimate the non-trivial weight coefficients. In BCS, the prior of

255 non-trivial weight coefficients is formulated as independent normal random variables with

256 relatively large variance to achieve an uninformative prior. The posterior distribution of $\widehat{\omega}_{t,s}^{2D}$

257 also follows a normal distribution and can be solved efficiently by a Markov chain Monte Carlo

258 (MCMC) simulation, leading to a series of random samples of $\widehat{\omega}_{t,s}^{2D}$ (e.g., Zhao and Wang, 2020;

259 Wang et al., 2022; Lyu et al., 2023). After repeatedly generating random samples of $\widehat{\omega}_{t,s}^{2D}$ $N_B$

260 times, the best estimate of DCT spectrum can be approximated by taking the mean of $N_B$

261 random samples of $\widehat{\omega}_{t,s}^{2D}$. In a MCMC simulation, the statistical independence of $N_B$ random

262 samples of $\widehat{\omega}_{t,s}^{2D}$ is guaranteed by taking only one sample in every larger number (e.g., 20 or 50)

263 Markov chain samples as the random sample. The associated statistical uncertainty can be

264 quantified using the standard deviation (SD) of $N_B$ samples. Another noteworthy advantage of

265 BCS is that it is applicable to a non-uniform measurement grid, a scenario commonly

266 encountered in site investigation (e.g., Zhao and Wang 2020; Guan et al., 2023a). Both the CS

267 and BCS algorithms have been compiled into a user-friendly free download software which is

268 available        from        the        corresponding        author's        website

269 (https://sites.google.com/site/yuwangcityu/software-download/bayesian-compressive-

270 samplingsensing-bcs). The approximated DCT spectrum for the cross-section allows

271 subsequent development of a unified representation of 2D spatial variability and quantification

272 of a cross-sectional similarity.

273 *3.2. Unified representation of 2D spatial variability using DCT-based ACF*

274 Note that each weight coefficient $\omega_{t,s}^{2D}$ in DCT spectrum corresponds to a specific basis

275 function $B_{t,s}(x_1, x_2)$, which is defined over specific $N_1$ and $N_2$ (see Equations (1) and (2)). In

276 other words, the DCT spectrum is relative to the dimension $N_1$ and $N_2$, which are essentially

277    determined by both the cross-section dimension and discretization resolution. This indicates

278    that direct comparison of the DCT spectrums of spatial variability $F(x_1, x_2)$ obtained in

279    different 2D cross-sections may not be feasible. To this end, a unified representation of 2D

280    spatial variability using DCT-based ACF is developed, which enables direct comparison

281    between 2D cross-sections with different spatial dimensions.

282          In signal processing, ACF is an effective tool for evaluating the correlation structure of

283    signals (e.g., Vanmarcke, 2010; Onyejekwe et al., 2016). ACF essentially measures the

284    correlation of a signal with a shifted version of itself. Mathematically, it describes how the

285    correlation between two points varies as the lag distance between the two points changes. In

286    the context of 2D cross-sectional spatial variability, ACF of $F(x_1, x_2)$ can be calculated as (e.g.,

287    Webster and Oliver, 2007; Vanmarcke, 2010):

$$
\text{ACF}[F(x_1, x_2), \tau_1, \tau_2] = \frac{E\left\{\left[F(x_1, x_2) - \mu_{F(x_1,x_2)}\right]\left[F(x_1+\tau_1, x_2+\tau_2) - \mu_{F(x_1,x_2)}\right]\right\}}{\sigma_{F(x_1,x_2)}^2} \tag{5}
$$

288    in which $\tau_1, \tau_2$ are the lag distances along $x_1$ and $x_2$ directions, respectively; $\mu_{F(x_1,x_2)}$ and

289    $\sigma_{F(x_1,x_2)}$ are the mean value and SD of all data points in $F(x_1, x_2)$, respectively. Equation (5)

290    does not assume stationarity. The ACF reflects the auto-correlation structure of 2D spatial

291    variability with respect to its mean value. It is a normalized and non-parametric measure

292    because it is normalized by the variance value and it is not fitted to any parametric function

293    form. In many fields, ACF has been widely used to identify predominant patterns/frequencies

294    embedded in the signals or images of interest (e.g., Priestley, 1981; Rafiee and Tse, 2009;

295    Zhang et al., 2021). Mathematically, ACF is closely related to the spectrum, and they form a

296    Wiener–Khinchin transform pair (e.g., Priestley, 1981). Although ACF may be used as a

297    surrogate to represent patterns of 2D cross-sectional spatial variability, it is worth noting that

298    calculating 2D ACF accurately and efficiently is usually difficult. Conventional approach of

299    calculating 2D ACF using Equation (5) often yields instable ACF values at large lag distances

13

300    and may be subject to significant computational efforts when dealing with high-resolution

301    images/matrices (e.g., Phoon and Fenton, 2004). To tackle this issue, this subsection derives a

302    new efficient formulation of ACF based on the DCT spectrum obtained from BCS:

$$
\begin{aligned}
&\text{ACF}[F(x_1, x_2), \tau_1, \tau_2] \\
&= \frac{1}{\sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D\,2}} \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D\,2} \, \text{ACF}[B_{t,s}(x_1, x_2), \tau_1, \tau_2]
\end{aligned}
\tag{6}
$$

$$(t, s \neq 1 \; concurrently)$$

303    in which $\text{ACF}[B_{t,s}(x_1, x_2), \tau_1, \tau_2]$ is the ACF of $B_{t,s}(x_1, x_2)$. Step-by-step derivation of

304    Equation (6) is provided in Appendix. Equation (6) shows that the ACF of 2D spatial variability

305    $F(x_1, x_2)$ is a weighted summation of the ACFs of 2D DCT basis functions, which are functions

306    of lag distances $\tau_1$ and $\tau_2$. The weight is the corresponding squared weight coefficient in DCT

307    spectrum. Note that Equation (6) establishes a theoretical basis for the unified representation

308    of 2D cross-sectional spatial variability, including non-stationarity and spatial variability

309    anisotropy. Equation (6) can also be interpreted as a special case of covariance decomposition

310    in traditional Karhunen-Loève expansion (e.g., Huang et al., 2001). Moreover, DCT-based

311    ACF enables a direct and convenient comparison between different 2D cross-sections. Using

312    $N_B$ random samples of DCT spectrum, $N_B$ DCT-based ACFs are obtained by substituting $N_B$

313    random samples of $\widehat{\omega}_{t,s}^{2D}$ into Equation (6). The best estimate DCT-based ACF is calculated as

314    the mean of $N_B$ DCT-based ACFs. Statistical uncertainty of approximated DCT spectrum also

315    propagates to DCT-based ACFs and can be quantified using SD of the $N_B$ DCT-based ACFs.

### 3.3. Quantification of DCT-based ACF similarity between cross-sections

317    With the DCT-based ACFs of two cross-sections determined in Section 3.2, cross-sectional

318    similarity can be quantified by the similarity between the corresponding DCT-based ACFs.

319    Consider, for example, two cross-sections A and B. Note that the actual dimensions of cross-

320      sections A and B may be different, leading to the different dimensions and patterns of the

321      corresponding DCT-based ACFs. To enable fair and effective comparison between two cross-

322      sections, only the largest overlapped sections with common lag distances of DCT-based ACFs

323      of cross-sections A and B are used accordingly. For example, if the sizes of spatial variability

324      matrices for cross-sections A and B are (200, 300) and (300, 200), respectively, only the

325      overlapped DCT-based ACFs with an identical range of lag distances, i.e., $\tau_1 = 0, 1, 2, \ldots, 199$

326      and $\tau_2 = 0, 1, 2, \ldots, 199$, are considered for similarity quantification. The overlapped ACFs

327      called effective ACFs offer a benchmark for comparison of two different cross-sections in a

328      statistical manner. Mathematically, a generalized cosine similarity between the effective DCT-

329      based ACFs of cross-sections A and B is calculated as (e.g., Dong et al., 2006; Nguyen and

330      Bai, 2011; Hu and Wang 2024):

$$\rho_{AB} = \frac{tr\left(\mathbf{ACF}_A \cdot \mathbf{ACF}_B^T\right)}{\sqrt{tr\left(\mathbf{ACF}_A \cdot \mathbf{ACF}_A^T\right)}\sqrt{tr\left(\mathbf{ACF}_B \cdot \mathbf{ACF}_B^T\right)}} \qquad (7)$$

331      in which $\rho_{AB}$ is defined as the similarity value between cross-sections A and B; $\mathbf{ACF}_A$ and

332      $\mathbf{ACF}_B$ are matrix representations of the effective 2D DCT-based ACFs of cross-sections A and

333      B, respectively; "T" is a transpose operation of a matrix; "$tr$" is the trace operation of a matrix.

334      This formula is equivalent to calculating the sum of element-wise product of $\mathbf{ACF}_A$ and $\mathbf{ACF}_B$,

335      divided by the product of the Frobenius norms of $\mathbf{ACF}_A$ and $\mathbf{ACF}_B$. $\rho_{AB}$ is therefore defined

336      over the range of [-1, 1]. Equation (7) essentially treats $\mathbf{ACF}_A$ and $\mathbf{ACF}_B$ as high-dimensional

337      vectors and measures the cosine value of the angle between the two vectors. High $\rho_{AB}$ indicates

338      closeness between the two vectors and hence high similarity between cross-sections A and B,

339      and vice versa.

340      Note that a deterministic $\rho_{AB}$ is obtained when substituting the corresponding best

341      estimate DCT-based ACFs of cross-sections A and B into Equation (7). To consider the

342    associated statistical uncertainty in both cross-sections A and B simultaneously, Equation (7)

343    may be used in a probabilistic manner. A random sample of $\rho_{AB}$ is calculated using Equation

344    (7) by substituting a pair of random samples of DCT-based ACFs of cross-sections A and B

345    respectively into Equation (7). After repeating the $\rho_{AB}$ calculation for all $N_B$ pairs of DCT-

346    based ACFs of two cross-sections, statistical analysis is performed on the obtained $N_B$ $\rho_{AB}$

347    values. The statistical uncertainty associated with the cross-sectional similarity quantification

348    is expressed by the SD of the $\rho_{AB}$ samples. The SD reflects the variability of the cross-sectional

349    similarity quantification in the presence of uncertainties in both cross-sections. Note that in

350    engineering practice, the required number of $N_B$ depends on the characteristics of the spatial

351    variability in the two cross-sections. The optimum value of $N_B$ may be identified by examining

352    the convergence behavior of the obtained similarity values.


353    **4. Implementation procedures**

354    To facilitate its applicability in engineering practice, this section summarizes the

355    implementation procedure of the proposed method for cross-sectional similarity quantification.

356    For example, two cross-sections, e.g., A and B, are to be evaluated. Five steps are involved in

357    implementing the cross-sectional similarity quantification between A and B, as described

358    below:

359        **Step 1**: Obtain the actual spatial dimensions (i.e., depths and horizontal lengths), and

360            determine the corresponding spatial resolutions $N_1 \times N_2$ of 2D spatial variability

361            $F(x_1, x_2)$ for cross-sections A and B, respectively. For example, if a cross-section has

362            a depth of 20m and a horizontal distance of 30m, spatial resolutions of 0.1m along both

363            directions will lead to a discretized 2D cross-section of shape 200×300.

364    **Step 2**: Compile the available soil property data within the cross-sections A and B as

365         measurement data $Y(x_1, x_2)$, which is a subset of $F(x_1, x_2)$. This step leads to two sets

366         of $Y(x_1, x_2)$ for cross-sections A and B, respectively.

367    **Step 3**: Perform BCS simulation to generate $N_B$ (e.g., $N_B$=500) random samples of

368         DCT spectrum from the corresponding measurement data $Y(x_1, x_2)$ in two cross-

369         sections, respectively.

370    **Step 4**: Calculate the $N_B$ DCT-based ACFs using Equation (6) and $N_B$ random samples

371         of DCT spectrum for cross-sections A and B, respectively.

372    **Step 5**: Perform probabilistic cross-sectional similarity quantification. $N_B$ DCT-based

373         ACFs of cross-section A are randomly paired with the $N_B$ DCT-based ACFs of cross-

374         section B. One similarity value is then obtained using Equation (7) for each pair of

375         DCT-based ACFs. $N_B$ pairs of DCT-based ACFs lead to $N_B$ similarity values.

376         Statistical analysis is then performed on the obtained $N_B$ similarity values.

377    With the mean of these $N_B$ similarity values, the similarity between cross-sections A and B can

378    be evaluated based on a pre-specified threshold. The threshold is purpose-dependent and

379    problem-specific. The question of how to select an optimal threshold is an interesting research

380    topic and will be investigated in a future study. One possible approach is to develop

381    characteristic values of similarity for a specific geotechnical problem based on many case

382    studies performed in similar geological settings. The statistical uncertainty of cross-sectional

383    similarity is quantified using SD of $N_B$ similarity values. Note that geotechnical analysis is

384    purpose-dependent and problem-specific. Different engineering projects might be sensitive to

385    the spatial variability of soil properties to different extents. This study only considers the

386    statistical similarity (e.g., the ACF similarity) of 2D cross-sectional spatial variability of one

387    soil property. This study does not consider the response of structures installed in the soil as a

388    result of spatial variability. Geotechnical analyses of different projects still need to resort to

389    specific domain knowledge from geotechnical engineers, even for the same or highly similar

390    site but with different purposes (e.g., deep foundation design versus liquefaction assessment)

391    (e.g., Leung, 2023). In the following section, the proposed method is illustrated using numerical

392    examples.


## 393    5. Illustrative examples

### 394    *5.1. Numerical examples of soil property cross-sections*

395    In this section, three cross-sections, i.e., namely A, B, and C, are simulated for illustration of

396    the proposed method. The configurations of these three cross-sections are summarized in Table

397    1. These three cross-sections have different spatial dimensions. Cross-section A has a depth of

398    20m and a width of 20m, i.e., a 20m×20m cross-section. Cross-sections B and C are 30m×20m

399    and 20m×30m cross-sections, respectively. A discretization resolution of 0.1m is adopted for

400    these three cross-sections, leading to discretized cross-sections with spatial resolution $N_1 \times N_2$

401    = 200×200, 300×200, and 200×300, respectively. Non-stationary undrained shear strength $s_u$

402    data are simulated using random field (RF) models for these three discretized cross-sections.

403    The non-stationary $s_u$ random fields are realized by adding a non-stationary trend function to a

404    2D zero-mean random field. As summarized in Table 1, in these three cross-sections, the $s_u$

405    data have different trend functions, which are formulated as the sum of 50 kPa and a scaled

406    cosine function term (in kPa) in different frequencies or phases. The cosine trend functions

407    adopted herein are to model the periodic property of geological depositional conditions (e.g.,

408    Einsele et al., 1996). Note that the trend functions of cross-section A and cross-section B

409    exhibit the same frequency, i.e., 0.5, which is equivalent to a period of around 12.5m, while

410    cross-section B incorporates an additional spatial shift of 10m. Cross-section B may be

411    interpreted as a cross-section exhibiting similar geological depositional conditions to A, but

412    occurring at another elevation level, as a scenario commonly encountered in engineering

practice. Cross-section C contains a frequency of 1, which is higher than cross-sections A and

B, and does not incorporate any spatial shift. The trend functions of the three cross-sections are

illustrated in Figures 5a-c, respectively. For each cross-section, an SD of 5 kPa and a Gaussian

auto-correlation structure are adopted for the 2D zero-mean random field. Different correlation

lengths along vertical and horizontal directions are configured for each cross-section. As shown

in Table 1, the vertical correlation lengths for three cross-sections are 1m, 1.5m, and 1m,

respectively; the horizontal correlation lengths for three cross-sections are 3.5m, 3m, and 2.5m,

respectively, leading to different spatial variability anisotropy structures for $s_u$ data. For each

cross-section, one $s_u$ data cross-section is realized and used for illustration, as shown in Figures

5d-f, respectively. Each $s_u$ data cross-section is a realization of a random field simulated by the

spectral representation method (e.g., Shinozuka and Deodatis 1991; Müller et al., 2022).

Table 1. Configurations of simulated undrained shear strength $s_u$ data for three cross-sections

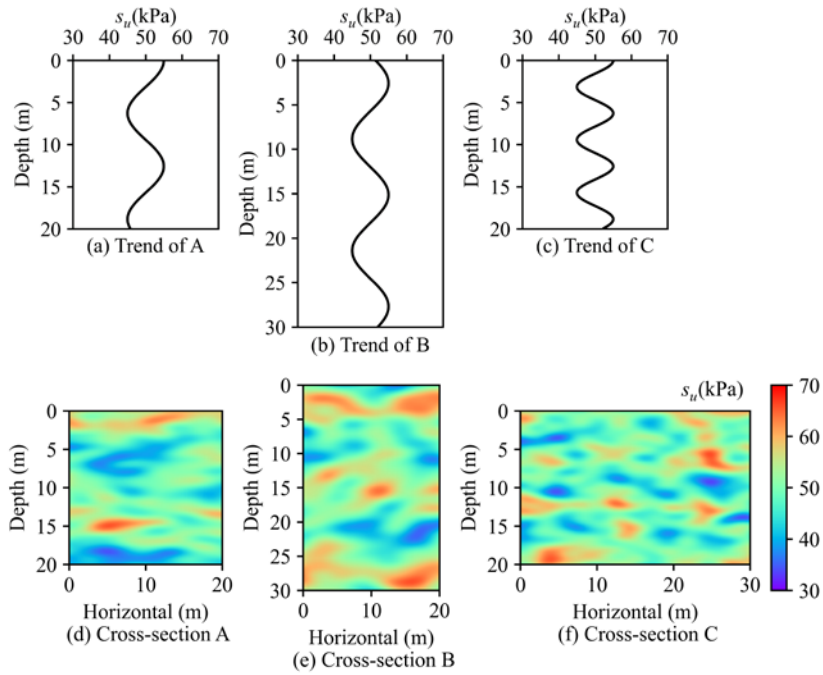| Parameters | Cross-section A | Cross-section B | Cross-section C |
|---|---|---|---|
| Depth (m) | 20 | 30 | 20 |
| Width (m) | 20 | 20 | 30 |
| Trend $s_u(z)$ versus depth (kPa) | $50+5\times\cos(0.5\times z)$ | $50+5\times\cos(0.5\times(z+10))$ | $50+5\times\cos(z)$ |
| Standard deviation (kPa) | 5 | 5 | 5 |
| Correlation function | Gaussian | Gaussian | Gaussian |
| Vertical correlation length (m) | 1 | 1.5 | 1 |
| Horizontal correlation length (m) | 3.5 | 3 | 2.5 |

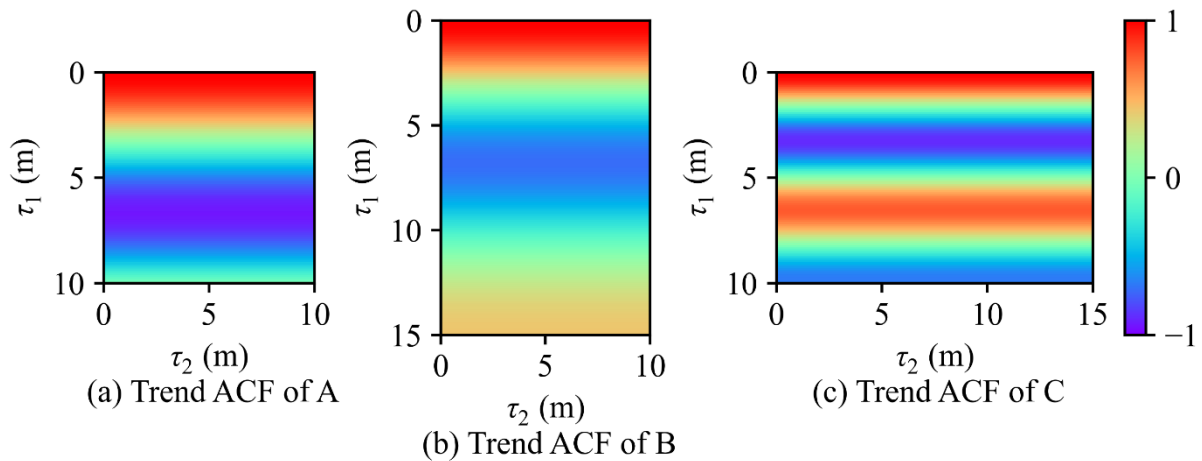Figure 5. Simulated undrained shear strength ($s_u$ in kPa) data cross-sections A, B, and C

Note that the three cross-sections are configured to illustrate the challenges of cross-sectional similarity quantification. In Figure 5, it is seen that cross-sections A, B, and C have different spatial dimensions and show different non-stationary and spatial variability anisotropy patterns. It is very challenging to rationally quantify the similarity among these cross-sections using conventional statistical methods. In this study, the derived DCT-based ACF tackles this challenge and offers an effective way to quantify the cross-sectional similarity among these three cross-sections. For each of the three cross-sections, the associated DCT spectrum can be readily obtained using Equation (3), and subsequently, the associated DCT-based ACF can be calculated using Equation (6). Note that the three cross-sections are respectively synthesized by adding up a non-stationary trend function and a 2D zero-mean random field. Therefore, the ACFs of $s_u$ cross-sections are controlled by both the underlying trends and zero-mean RFs. The ACFs of the trend functions for the three cross-sections are shown in Figures 6a-6c, respectively, while the corresponding theoretical RF ACFs are shown in Figures 6d-6f, respectively. Figures 6g-6i, respectively, show the ACFs of the three
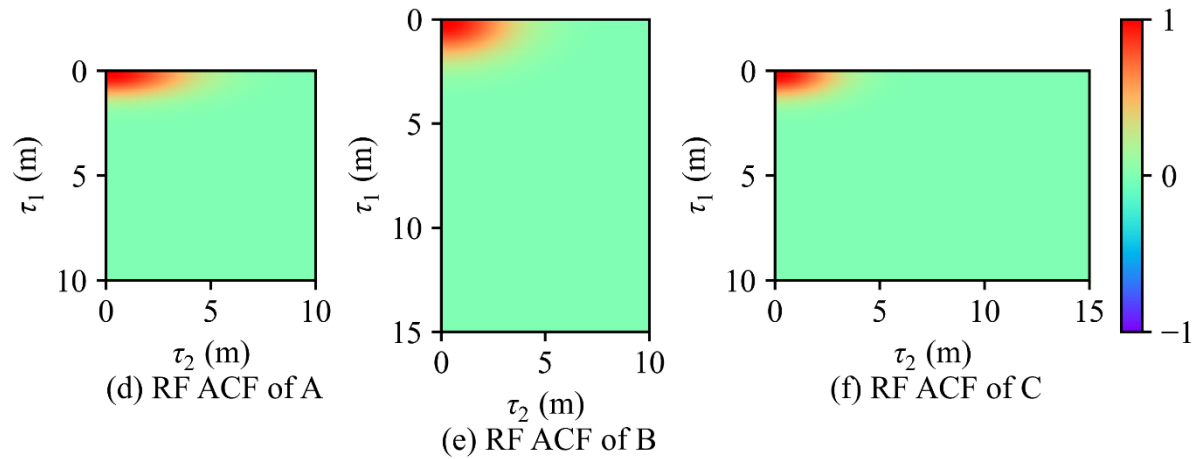
443  synthesized $s_u$ cross-sections. The 2D ACFs are plotted as a colormap versus varying lag

444  distances along two directions. Since the auto-correlation decreases as the lag distance

445  increases, only half of the maximum lag distances along both directions are considered as

446  shown in the figure. In Figures 6a-6c, it is shown that the DCT-based ACFs of the trend

447  functions behave also like cosine functions, with ACF values fluctuating at corresponding

448  frequencies along the depth direction. In Figures 6d-6f, theoretical RF ACFs decay along both

449  directions accordingly to the corresponding RF parameters. Note that in Figures 6g-6i, the

450  DCT-based ACFs of the three $s_u$ cross-sections show combined patterns exhibiting features of

451  the ACFs from the trend functions and the ACFs from the RFs. This indicates that ACF not

452  only may be used to characterize a zero-mean RF, but also simultaneously characterize the

453  underlying deterministic trend function (e.g., Brockwell and Davis, 1991).

454      Note that the three DCT-based ACFs have different shapes. To fairly compare these

455  DCT-based ACFs, the largest overlapped sections between any two cross-sections are selected,

456  as delineated by red dashed lines in Figures 6g-6i. For any two cross-sections, a similarity value

457  is calculated using Equation (7) and the corresponding overlapped DCT-based ACFs. The

458  similarity values for different pairs of cross-sections are calculated as $\rho_{AB} = 0.97$, $\rho_{AC} = 0.27$,

459  and $\rho_{BC} = 0.36$. The similarity values are consistent with the theoretical configurations of the

460  three cross-sections. It has been indicated in Table 1 that cross-sections A and B demonstrate

461  better spectral coherence since their non-stationary trend functions have an identical frequency,

462  although the trend function of cross-section B has a spatial shift. However, the trend function

463  of cross-section C contains a higher frequency than A and B. From a spectral perspective, cross-

464  section C may not be similar to cross-sections A and B. Note that the similarity quantification

465  using DCT-based ACF is invariant to the spatial offset values of spatial variability. Similar

466  examples with different spatial offset values had been analyzed, and consistent results were

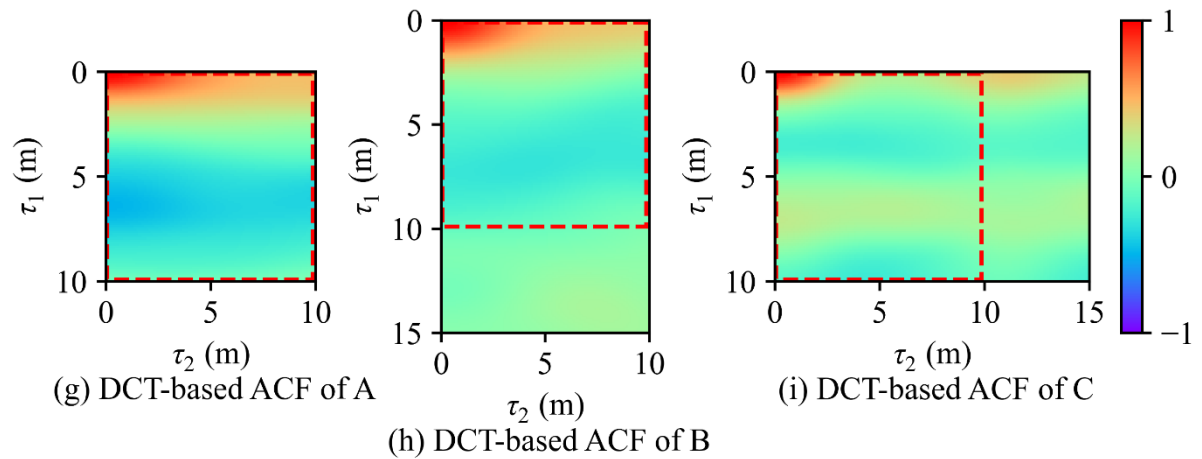467  obtained. In this paper, only the cross-sections configured in Table 1 are presented for brevity.

(a) Trend ACF of A

(b) Trend ACF of B

(c) Trend ACF of C

(d) RF ACF of A

(e) RF ACF of B

(f) RF ACF of C

(g) DCT-based ACF of A

(h) DCT-based ACF of B

(i) DCT-based ACF of C
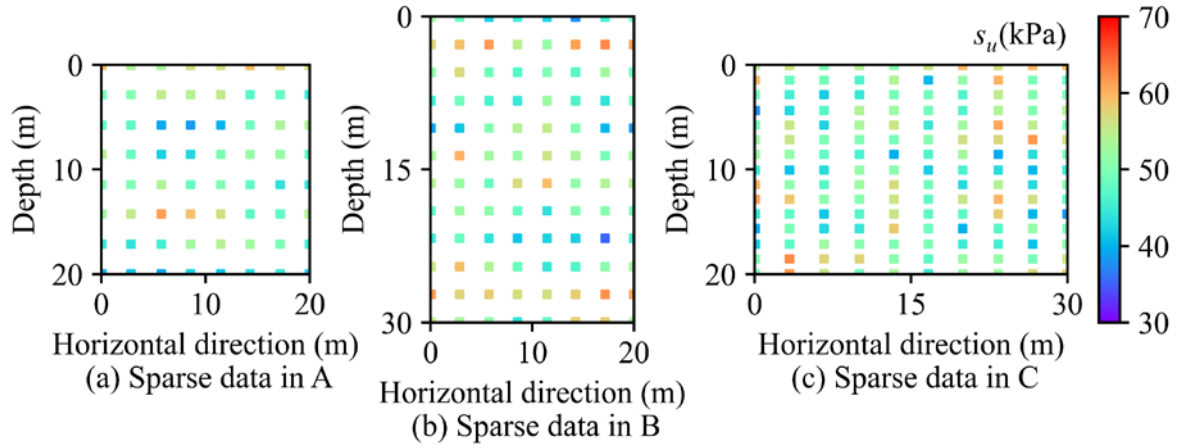
472     Figure 6. DCT-based ACF in cross-sections A, B, and C: (a) Trend ACF of A; (b) Trend ACF
473     of B; (c) Trend ACF of C; (d) RF ACF of A; (e) RF ACF of B; (f) RF ACF of C; (g) ACF of
474     synthetic $s_u$ data in A; (h) ACF of synthetic $s_u$ data in B; (i) ACF of synthetic $s_u$ data in C

475

476

477     Note that the above cross-sectional similarity quantification is based on the three

478     simulated cross-sections with complete $s_u$ data. To illustrate the typical scenario of sparse

479     investigation data, the following subsection performs cross-sectional similarity quantification

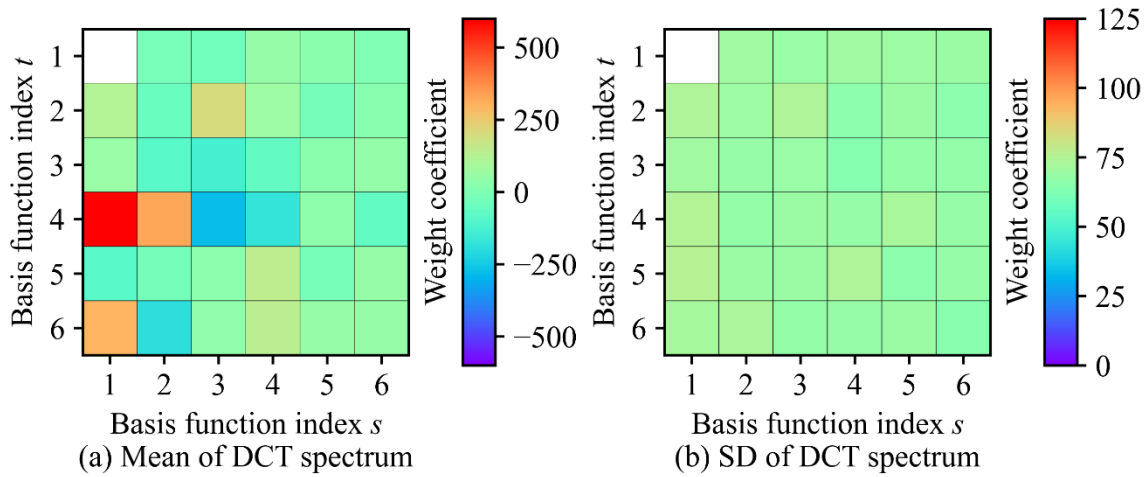480     using the same examples, but with sparsely measured data.



481

482                     Figure 7. Sparse data measured from cross-sections A, B, and C

483     *5.2. Similarity quantification between simulated cross-sections using sparse data*

484     To illustrate the challenge of sparse investigation data in cross-sections, selective $s_u$ data are

485     sampled from the simulated cross-sections as measurement data, which are subsequently used

486     for probabilistic quantification of similarity between cross-sections with consideration of

487     statistical uncertainty. As shown in Figure 7, uniform grid sampling is implemented in the three

488     cross-sections. In cross-sections A, B, and C, 8×8, 12×8, and 15×10 $s_u$ data are measured,

489     respectively. The measured data account for around 0.16%, 0.16%, and 0.25%, respectively,

490     of the corresponding discretized cross-section. The sampling ratios adopted are generally

491     comparable to the engineering practice of site investigation, where the ratio of the volume of

492     sampled soils over the volume of soils loaded/affected is normally around or less than 0.1%,

493     depending on the project requirements, site complexity, and the level of details needed (e.g.,

494     Look 2014; Guan and Wang, 2020). For the three cross-sections, the spatial resolutions

495     $N_1 \times N_2$ are set as the original simulated cross-section, i.e., 200×200 for A, 300×200 for B,

23

496    and 200×300 for C. The corresponding measurement data are then adopted as Y($x_1, x_2$), and

497    subsequently, BCS simulation is performed to generate $N_B$=500 samples of DCT spectrum for
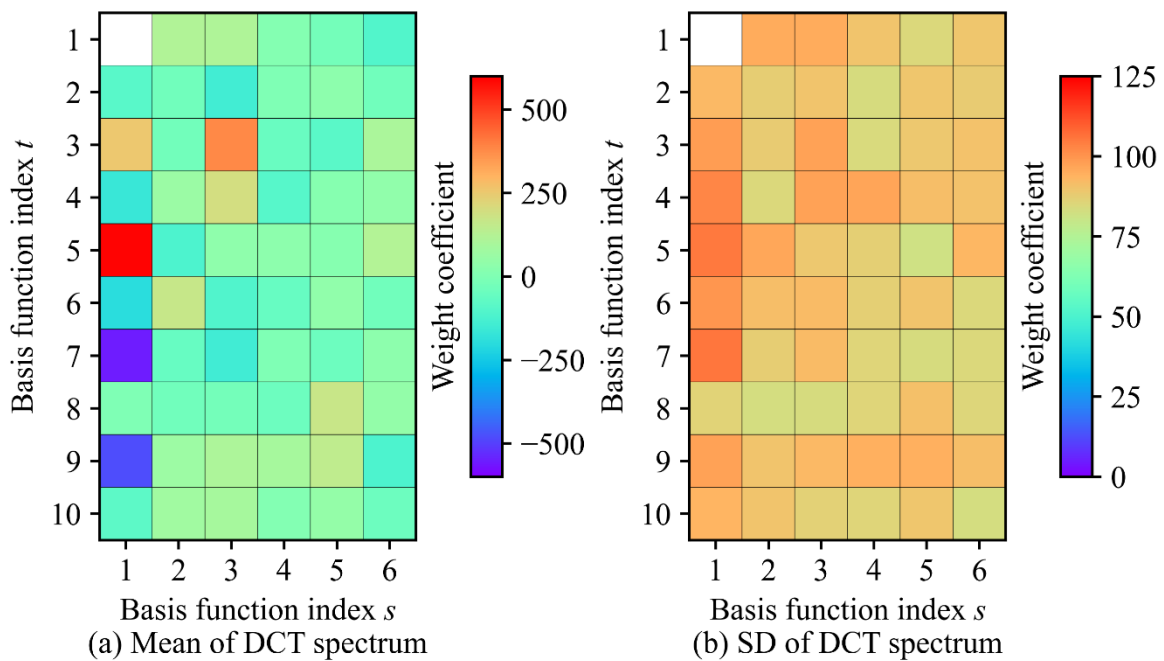
498    each cross-section.

499



(a) Mean of DCT spectrum        (b) SD of DCT spectrum

Figure 8. Statistics of DCT spectrum in cross-section A

502



(a) Mean of DCT spectrum        (b) SD of DCT spectrum

Figure 9. Statistics of DCT spectrum in cross-section B

505

24

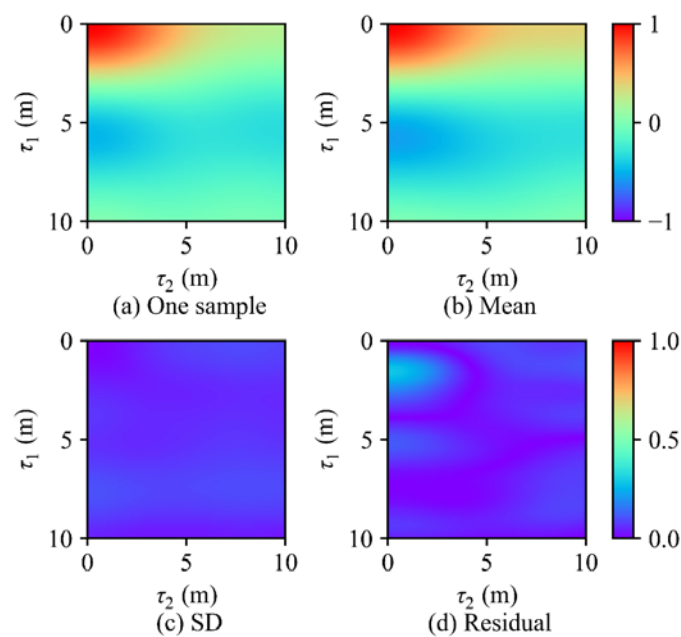Figure 10. Statistics of DCT spectrum in cross-section C

Figures 8a and 8b show the statistics of DCT spectrum for cross-sections A in the form of colored meshes. As shown in Figure 8a, 36 weight coefficients are identified from sparse data using BCS, and they are characterized by a 6×6 matrix of DCT spectrum. Each mesh is color-coded using the mean value of the weight coefficient with indexes $t$ and $s$, which are indexes (or frequencies) of the corresponding 2D basis functions, as indicated in Equation (2). Note that the $\omega_{1,1}^{2D}$ corresponding to the contribution of the constant basis function $B_{1,1}(x_1, x_2)$ is much greater than the other coefficients and is not required for the calculation of DCT-based ACF (see Equation (6)). Therefore, the upper left cell in Figure 8 is shown as empty for visualization clarity. It is observed that among these 35 weight coefficients, only a few of them are significant, with the remaining ones close to zeros. Figure 8b shows the SD values of the corresponding weight coefficients in Figure 8a, which reflect the statistical uncertainty of the DCT spectrum. Figures 9 and 10 show the statistics of $N_B$ samples of DCT spectrum for cross-sections B and C, respectively. Figures 8-10 display that the number of weight coefficients, as well as the indexes of significant coefficients, in the approximated DCT spectrum are different

522    in different cross-sections. The DCT spectrums are not directly comparable, since they

523    correspond to cross-sections with different spatial dimensions. DCT-based ACF is

524    subsequently calculated for a unified representation of 2D spatial variability in different cross-
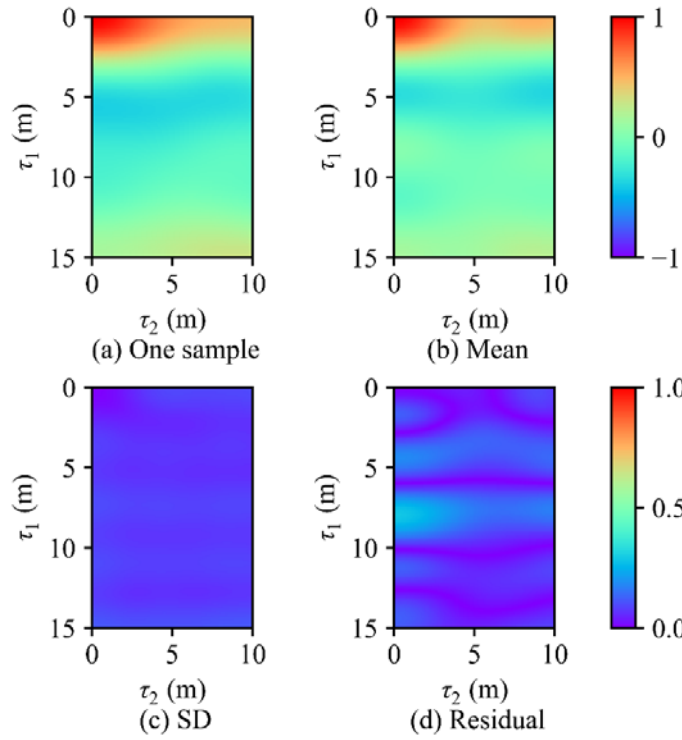
525    sections.

526            Using Equation (6), $N_B$ samples of DCT-based ACFs are obtained for each cross-

527    section. Figure 11 shows the statistics of the obtained DCT-based ACFs in cross-section A by

528    colormaps. Subplot (a) demonstrates one sample of DCT-based ACF, which portrays a possible

529    representation of the 2D spatial variability patterns of $s_u$ data in cross-section A. The associated

530    non-stationarity and spatial variability anisotropy are represented in DCT-based ACF in a

531    unified manner. Subplots (b) and (c) show the mean and the SD of $N_B$ samples of DCT-based

532    ACFs. Subplot (d) reveals the absolute residual between the mean in Subplot (b) and the

533    original DCT-based ACF in Figure 6g. The mean of DCT-based ACFs in Subplot (b) is

534    interpreted as the best estimate for the cross-section A in the presence of sparse investigation

535    data (e.g., see Figure 7a). Figures 12 and 13 show the statistics of the obtained DCT-based

536    ACFs in cross-sections B and C, respectively, following the same presentation format. It is

537    observed from Figures 11-13 that, for all three cross-sections, the colormaps in Subplots (c)

538    and (d) are generally comparable, indicating the quantified statistical uncertainty of DCT-based

539    ACF is rational.

540            To quantify the similarity between any two cross-sections with consideration of

541    statistical uncertainty, $N_B$ DCT-based ACF samples from one cross-section are randomly

542    paired with $N_B$ DCT-based ACF samples from another cross-section. For each pair of DCT-

543    based ACFs, a similarity value is calculated using Equation (7), leading to $N_B$ similarity values.

544    Probabilistic cross-sectional similarity quantification is then performed by statistical analysis

545    of these $N_B$ similarity values. Figure 14 shows the obtained similarity values by blue

546    histograms. Figure 14a presents the similarity values between cross-sections A and B. It shows

547 that for cross-sections A and B, the associated histogram of similarity values peaks at a value

548 approaching 1. The mean of the similarity values is calculated as 0.88 and is close to the true

549 similarity value 0.97, which is denoted by a vertical red dashed line in Figure 14a. High cross-

550 sectional similarity between A and B is reasonably quantified from sparse data. Figure 14b

551 presents the similarity values between cross-sections A and C. The mean of the similarity

552 values is calculated as 0.26 which is close to the true value 0.27. Although the correlation

553 lengths of cross-sections A and B are slightly different as indicated in Table 1, the similarity

554 quantification may be dominated by the respective trend functions that have the same frequency.

555 The low cross-sectional similarity between cross-sections A and C is also quantified accurately.

556 In Figure 14, the SD of similarity values can be interpreted as the statistical uncertainty of

557 quantified cross-sectional similarity, which integrates the statistical uncertainty of DCT

558 spectrum from both concerned cross-sections. The SD values are calculated as 0.06 for cross-

559 sections A and B, and 0.08 for cross-sections A and C. It suggests that for the above two

560 comparisons, the associated statistical uncertainty is relatively small. In other words, the three

561 cross-sections, particularly the associated trend functions, may be characterized well by sparse

562 data.


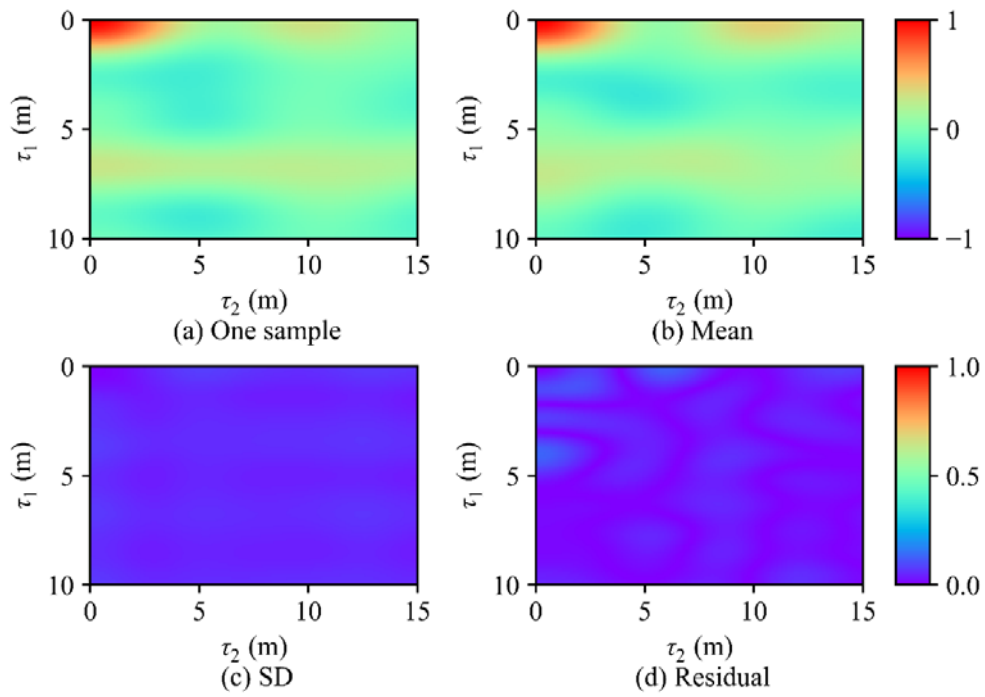
Figure 11. Statistics of DCT-based ACFs in cross-section A

565

Figure 12. Statistics of DCT-based ACFs in cross-section B
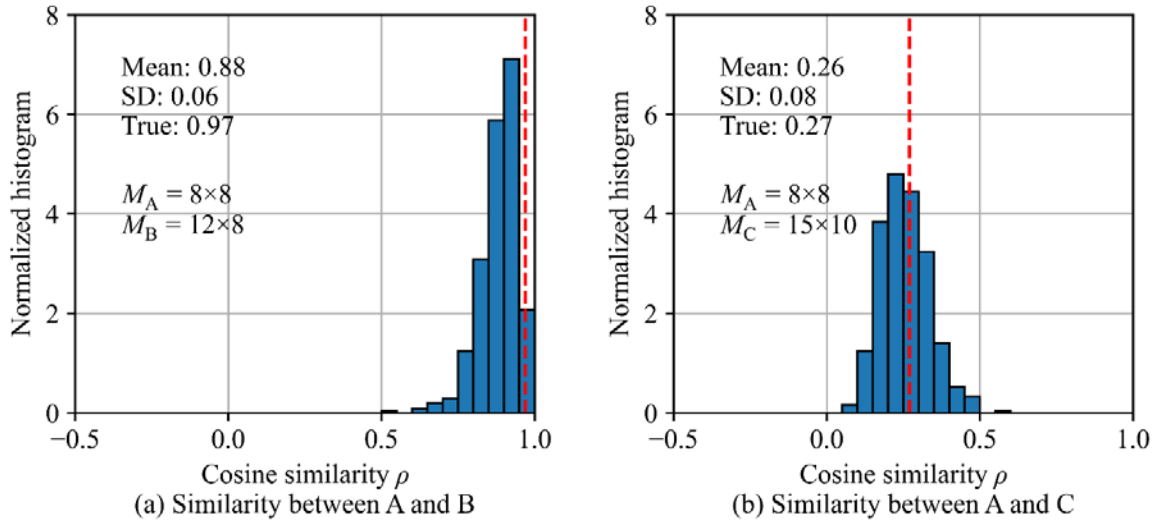
567



568

Figure 13. Statistics of DCT-based ACFs in cross-section C

570

571

572    Figure 14. Normalized histogram of cosine similarity values between cross-sections: (a) A
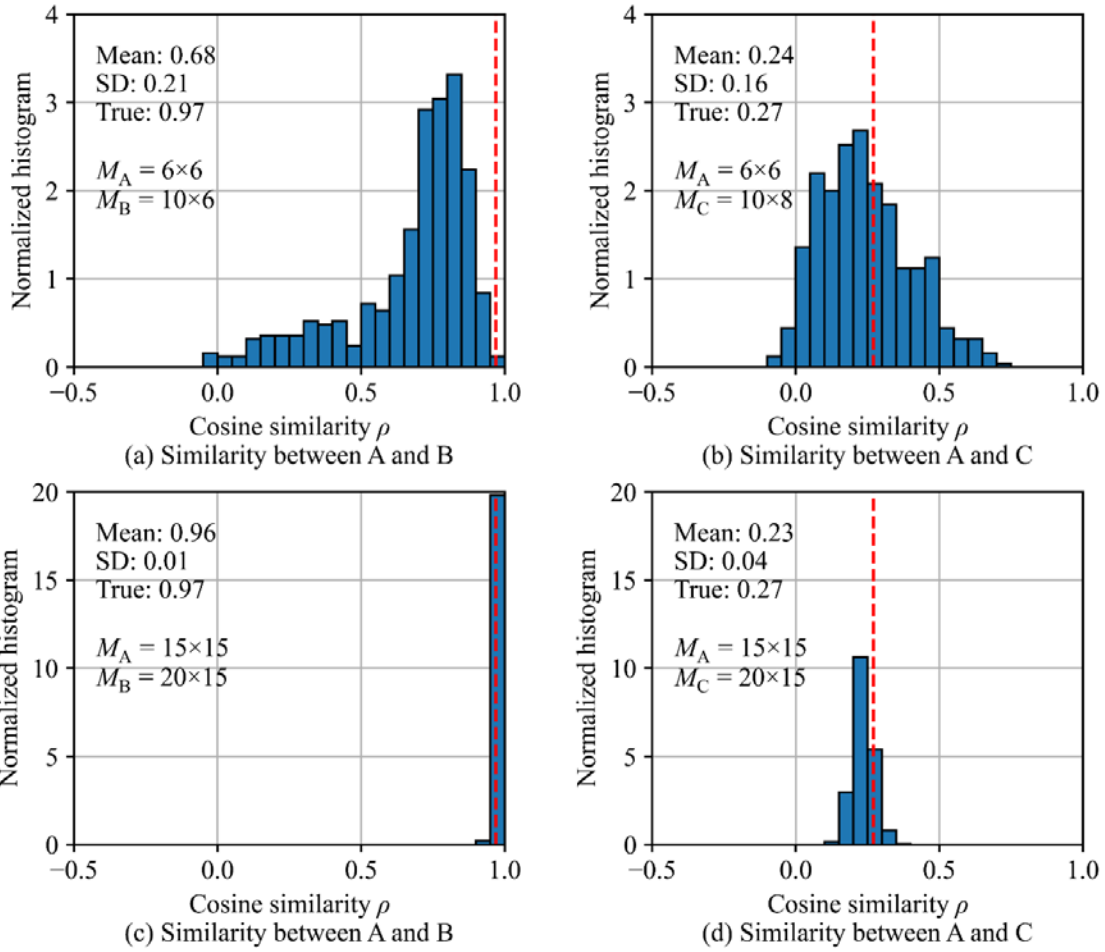
573    and B; (b) A and C

574

575    Table 2. Statistics of cross-sectional similarity values obtained from measurement data of

576    cross-sections A, B, and C

| $M$ scenario | Cross-sections | A | B | C |
|---|---|---|---|---|
| $M_A=6\times6$ | A | 1 | 0.68(0.21) | 0.24(0.16) |
| $M_B=10\times6$ | B | - | 1 | 0.35(0.14) |
| $M_C=10\times8$ | C | - | - | 1 |
| $M_A=8\times8$ | A | 1 | 0.88(0.06) | 0.26(0.08) |
| $M_B=12\times8$ | B | - | 1 | 0.47(0.09) |
| $M_C=15\times10$ | C | - | - | 1 |
| $M_A=15\times15$ | A | 1 | 0.96(0.01) | 0.23(0.04) |
| $M_B=20\times15$ | B | - | 1 | 0.36(0.03) |
| $M_C=20\times15$ | C | - | - | 1 |

Data format: Mean (Standard deviation)

577

Figure 15. Normalized histogram of cosine similarity values between cross-sections under different measurement scenarios ($M$)

### 5.3. Effect of the number of measured data points

This subsection investigates the effect of the number, $M$, of measured data points on the performance of cross-sectional similarity quantification. Two more measurement scenarios for the three cross-sections are added. One added scenario has a smaller number of measured $s_u$ data, i.e., $M$=6×6, 10×6, and 10×8 $s_u$ data with a uniform grid sampling are measured in cross-sections A, B, and C, respectively. Another added scenario has a larger number of measured $s_u$ data, i.e., $M$=15×15, 20×15, and 20×15 $s_u$ data are measured in the three cross-sections respectively. For each added scenario, cross-sectional similarity quantifications are performed, following the implementation procedures described in Section 4.

590     Figure 15 shows the histograms of cross-sectional similarity values for two added

591     scenarios, between A and B, and between A and C, respectively. Figure 15a shows the

592     similarity between cross-sections A and B, when the number of measurement data is relatively

593     small. In comparison to Figure 14a, it is observed that the mean of similarity values decreases

594     significantly from 0.88 to 0.68. In addition, the SD of similarity values increases significantly

595     from 0.08 to 0.21. Figure 15c corresponds to the similarity between cross-sections A and B,

596     when the number of measurement data is relatively large. It shows that the histogram is

597     narrowed down significantly and almost overlaps with the true value. Similar observations are

598     also obtained for similarity values between cross-sections A and C, where the corresponding

599     small and large measurement data number scenarios are shown in Figures 15b and 15d,

600     respectively. For cross-sectional similarity between A and C, it appears that the true similarity

601     value, which is as low as 0.27, can be accurately identified using extremely sparse data. Both

602     high cross-sectional similarity between A and B and low cross-sectional similarity between A

603     and C can be quantified effectively using sparse data. The results of this sensitivity study, as

604     well as the cross-sectional similarity between B and C, are summarized in Table 2. The results

605     indicate that the performance of the proposed method for quantifying cross-sectional similarity

606     depends on the number $M$ of available measured data in corresponding cross-sections. When

607     $M$ is low in the two cross-sections to be compared, the associated statistical uncertainty might

608     become dominant in the subsequent cross-sectional similarity quantification. In this case,

609     additional site investigation might be required to get more measurements and insights into the

610     spatial variability in concerned cross-sections. As $M$ increases, the quantified cross-sectional

611     similarity converges to the true value. Moreover, the proposed method also applies to scenarios

612     where the amount of measured data differs significantly in the two cross-sections, e.g., one

613     cross-section characterized with limited data while another characterized with much more data.

614     This enables the proposed method to be performed in a data-driven manner in practical site

615  investigation. Note that the relationship between the number $M$ of available measurement data

616  and the cross-sectional similarity is problem-specific and might not necessarily be a general

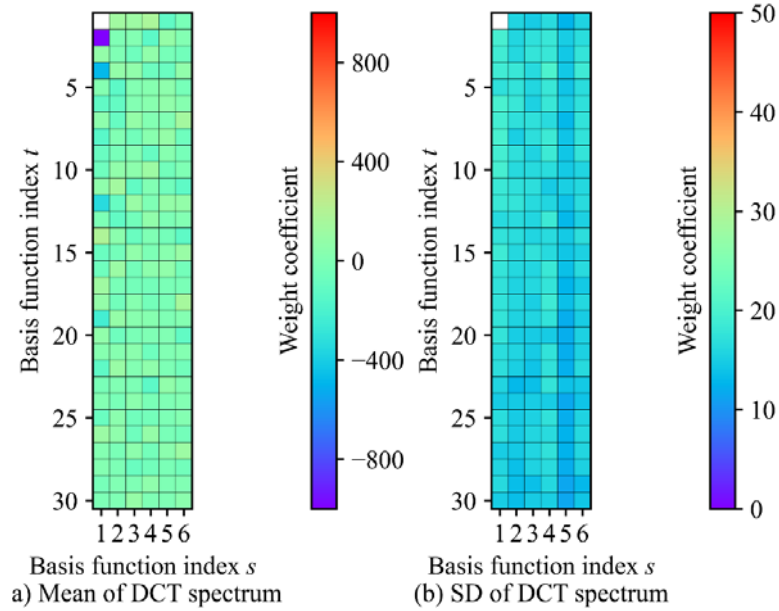617  one that can possibly be applied to other cross-sections.

618



619

620  Figure 16. Statistics of DCT spectrum at Site 1

621



622

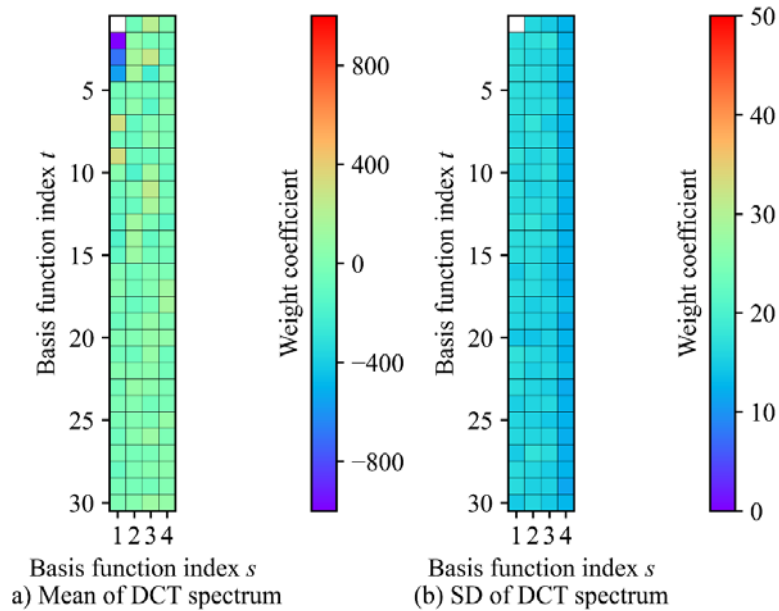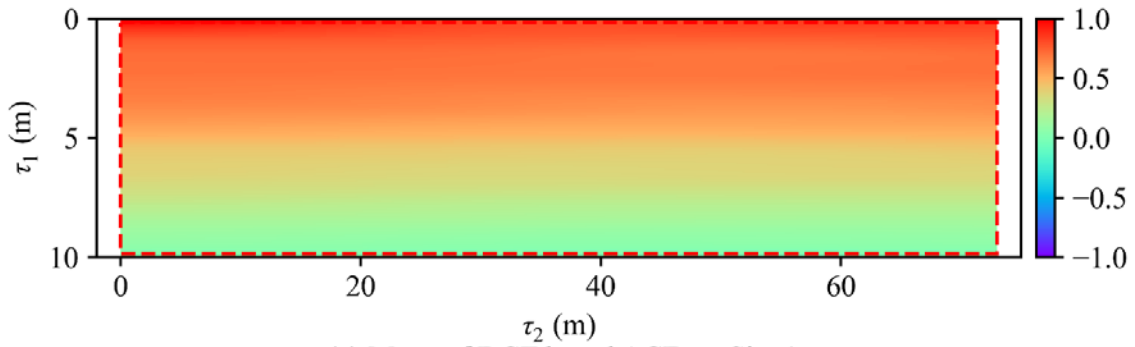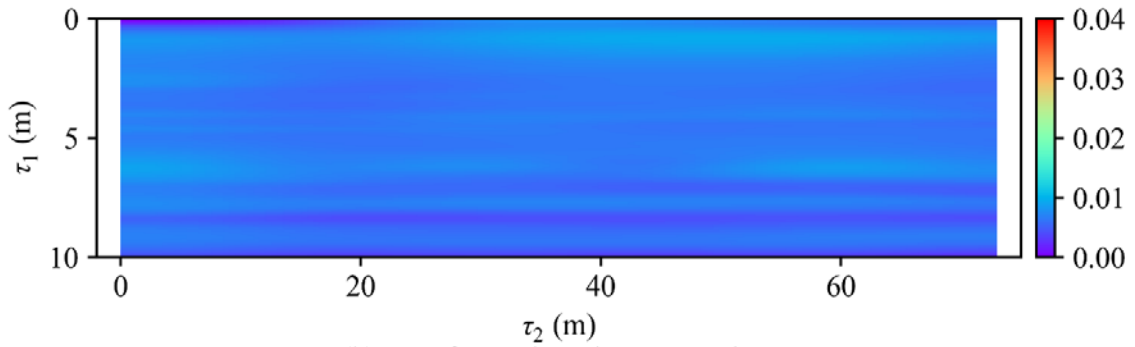623  Figure 17. Statistics of DCT spectrum at Site 2

## 6. Real examples

This section demonstrates an application of the proposed method to the real examples shown in Figure 1. Probabilistic quantification of cross-sectional similarity between Site 1 and Site 2 is performed, following the implementation procedures in Section 4. Sites 1 and 2 have a length of 145m and 245m, respectively, and they are both 20m deep. In step 1, a vertical resolution of 0.05m and a horizontal resolution of 0.5m are adopted to discretize the cross-sections at the two sites, leading to a cross-section image with a size of 400×290 for Site 1, and a cross-section image with a size of 400×490 for Site 2. In step 2, the CPT data (e.g., corrected cone resistance $q_t$ in this example) within these two cross-sections are obtained. As shown in Figures 2a and 2b, eight $q_t$ data profiles are within cross-section at Site 1 and seven profiles are within cross-section at Site 2. In step 3, $N_B$=500 samples of DCT spectrum are generated from $q_t$ data for each site using BCS. Figures 16 and 17 show the statistics of DCT spectrum at Site 1 and Site 2, respectively. It is seen that for both sites, the numbers of identified weight coefficients are different, with 179 coefficients for Site 1 and 119 coefficients for Site 2. The indexes of significant coefficients for the two sites are also apparently different. In step 4, $N_B$ DCT-based ACFs are calculated based on $N_B$ samples of DCT spectrum for both sites. Figures 18a and 18b show the mean and SD of $N_B$ samples of DCT-based ACFs at Site 1. Figures 19a and 19b show the corresponding results at Site 2. It is evident that the means of DCT-based ACFs in Figures 18a and 19a have generally consistent patterns, i.e., predominant spatial variability patterns along vertical directions and relatively minor variability patterns along horizontal directions. The SD maps in Figures 18b and 19b show similar magnitudes, suggesting the statistical uncertainty for these two sites is comparable. In step 5, to perform a probabilistic quantification of cross-sectional similarity, $N_B$ DCT-based ACFs at Site 1 are randomly paired with $N_B$ DCT-based ACFs at Site 2. Since the horizontal lengths are different for these two sites (i.e., 145m for Site 1 and 245m for Site 2), the associated DCT-based ACFs of these two sites have

33

649  different dimensions, as shown in Figures 18 and 19. Therefore, the overlapped sections of

650  Sites 1 and 2, denoted by red dashed lines in Figures 18a and 19a, are used for cross-sectional

651  similarity quantification. Using $N_B$ pairs of DCT-based ACFs and Equation (6), $N_B$ similarity

652  values are obtained, which are presented by a histogram in Figure 20. Note that the histogram

653  of similarity values is narrow and mainly located at the $\rho$ range of about 0.97 to 0.98. The mean

654  and SD of the $N_B$ similarity values are calculated as 0.977 and 0.002, respectively. According

655  to the results, the proposed method suggests that Sites 1 and 2 are highly similar, and the

656  associated statistical uncertainty is insignificant. Although the numbers of CPTs soundings are

657  sparse from both sites, the spatial variability patterns of $q_t$ data are prominent and consistent.

658  The increasing trends of $q_t$ data profiles at both sites, as shown in Figure 2 are comparable and

659  properly identified. A systematic study on spatial variability with increasing trend functions is

660  worth exploring in a future study to clearly demonstrate generalizability of the proposed

661  method. In addition, the results indicate that the proposed method performs well even for cross-

662  sections with non-uniformly measured data. This scenario can be regarded as incomplete data,

663  because any non-uniform measurement grid can be derived from a uniform measurement grid

664  by removing measurements from selected points. Hence, this scenario refers to one aspect of

665  MUSIC, which is "I" for incomplete data.

(a) Mean of DCT-based ACFs at Site 1
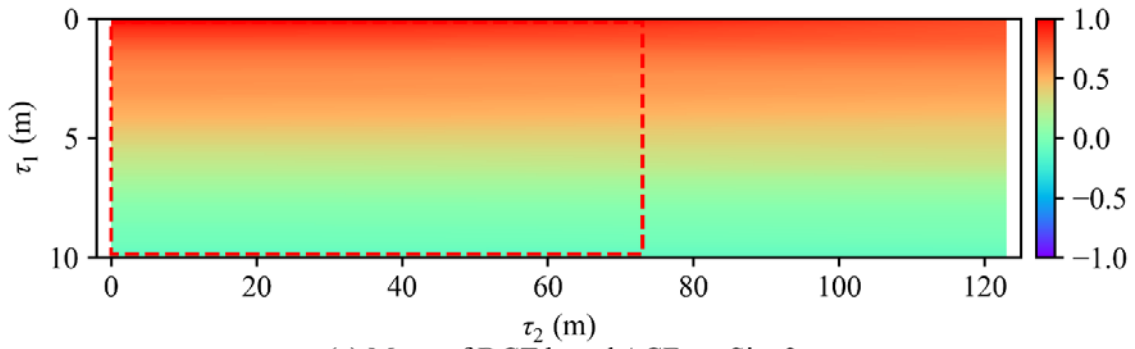

(b) SD of DCT-based ACFs at Site 1
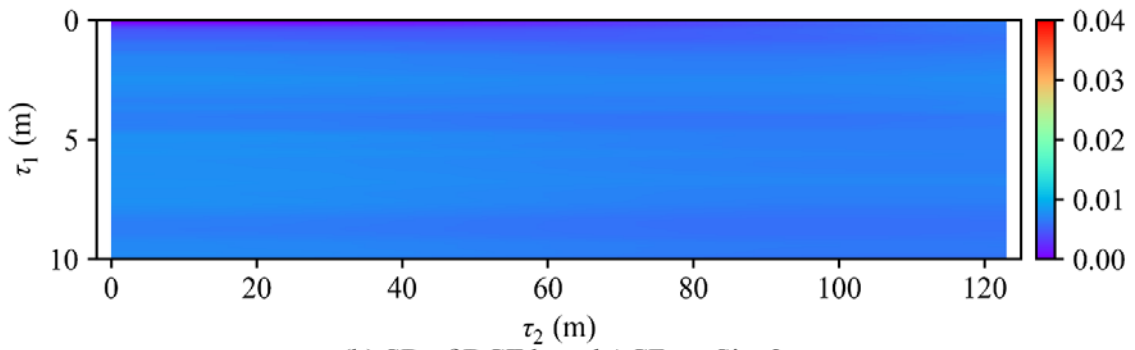
666

667                    Figure 18. Statistics of DCT-based ACFs at Site 1


(a) Mean of DCT-based ACFs at Site 2


(b) SD of DCT-based ACFs at Site 2

668

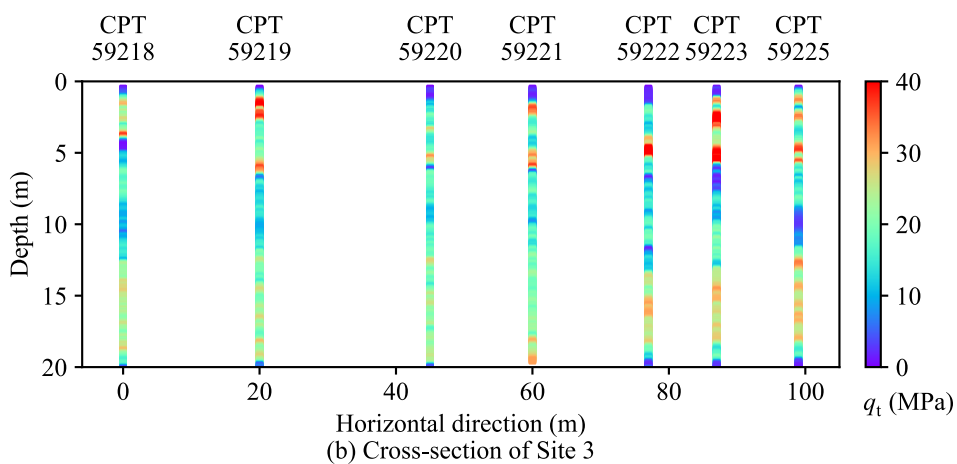669                    Figure 19. Statistics of DCT-based ACFs at Site 2

670

671

Figure 20. Normalized histogram of cosine similarity values between Sites 1 and 2 in the real data example
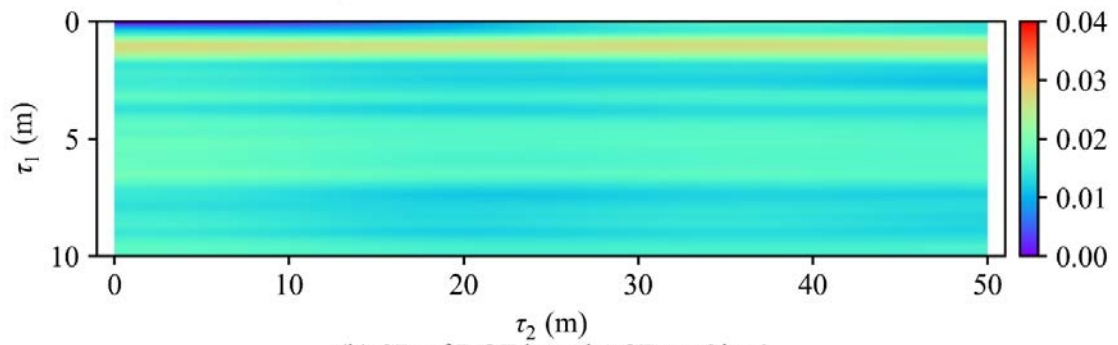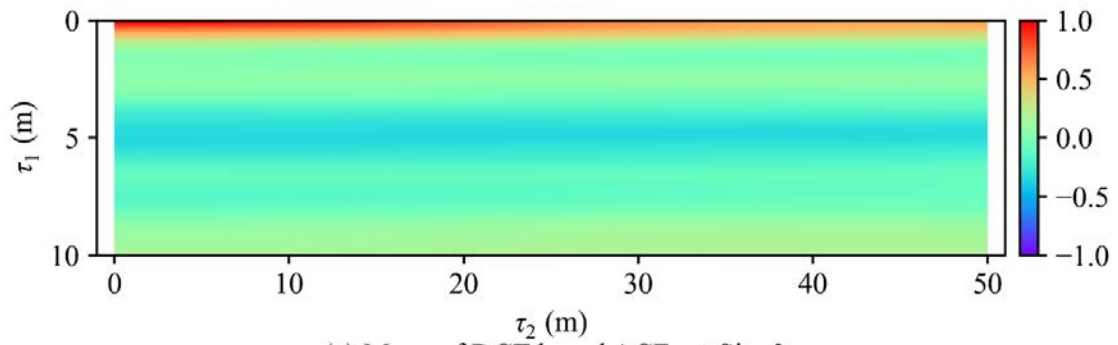


(a) Layout of seven CPTs in Site 3

675



(b) Cross-section of Site 3

676

Figure 21. Cone penetration tests (CPTs) performed in Site 3: (a) layout map of seven CPTs; (b) Cross-section of seven CPTs
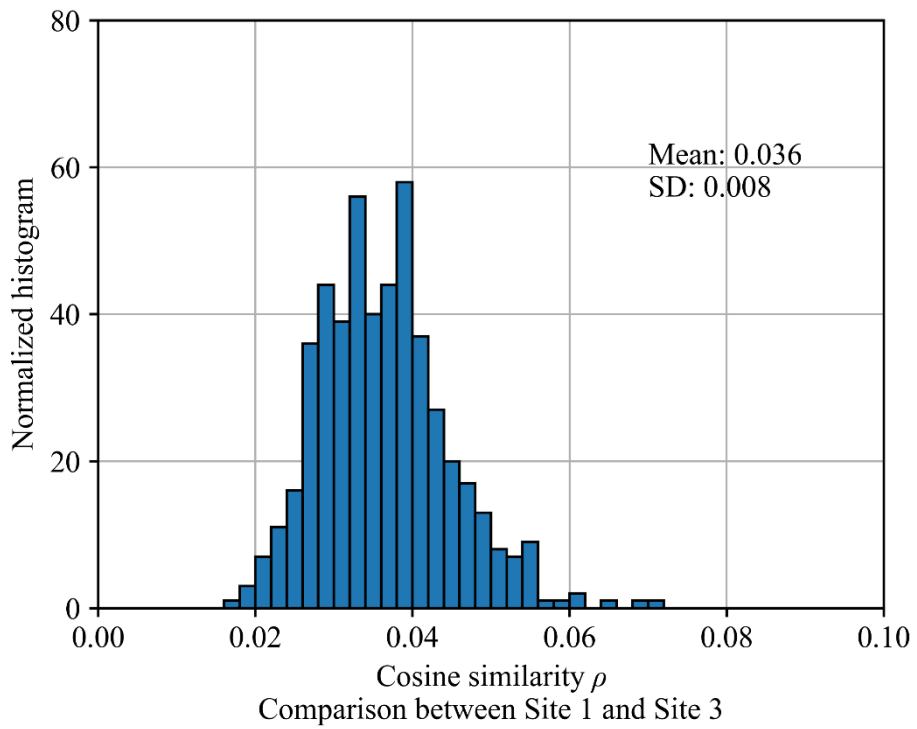
36

679        To further demonstrate the performance of the proposed method, another cross-section

680    example at Site 3 is compared with the cross-section at Site 1. As shown in Figure 21a, Site 3

681    is a cross-section with a horizontal length of around 100m, and seven CPTs were performed in

682    this cross-section. In contrast to Site 2, which is only approximately 700m away from Site 1,

683    Site 3 is relatively far from Site 1, and they are roughly 5km apart. In view of this spatial

684    distance, CPT data at Site 3 might exhibit different spatial variability patterns from that of Site

685    1. Figure 21b shows the $q_t$ data profiles in this cross-section. Figure 2a and Figure 21b are

686    visually different. Probabilistic quantification of cross-sectional similarity between Sites 1 and

687    3 is performed, following the implementation procedures described in Section 4. After

688    configuring the same spatial resolutions for Site 3 as Site 1 (i.e., vertical resolution of 0.05m

689    and a horizontal resolution of 0.5m), $N_B$=500 samples of DCT spectrum are generated for Site

690    3 using BCS. Subsequently, DCT-based ACFs are calculated for Site 3. Figures 22a and 22b

691    show the mean and SD of DCT-based ACFs at Site 3, respectively. It shows that the mean of

692    DCT-based ACFs in Figure 22a differs significantly from Figure 18a. The ACF at Site 3 decays

693    and fluctuates faster than the ACF at Site 1 along the depth direction. This suggests that the

694    correlation length along the depth direction at Site 3 is much smaller than that at Site 1. In

695    addition, the SD map in Figure 22b shows a higher magnitude than the SD map in Figure 18b,

696    suggesting the statistical uncertainty of $q_t$ data at Site 3 is greater than Site 1. Then, $N_B$ DCT-

697    based ACFs at Site 3 are randomly paired with $N_B$ DCT-based ACFs at Site 1. The resulting

698    $N_B$ similarity values are calculated and shown in Figure 23. The mean and SD of the $N_B$

699    similarity values are calculated as 0.036 and 0.008, respectively. The difference between Figure

700    20 (similarity between Sites 1 and 2) and Figure 23 (similarity between Sites 1 and 3) is stark.

701    According to the results, the proposed method suggests that Sites 3 and 1 are not similar.

(a) Mean of DCT-based ACFs at Site 3



(b) SD of DCT-based ACFs at Site 3

702

703                    Figure 22. Statistics of DCT-based ACFs at Site 3



Mean: 0.036
SD: 0.008

Cosine similarity $\rho$
Comparison between Site 1 and Site 3

704

705      Figure 23. Normalized histogram of cosine similarity values between Site 1 and Site 3

## 7. Conclusions

In this study, a novel data-driven method was proposed to quantify 2D cross-sectional similarity based on the spatial variability of one soil property. To the best of the authors' knowledge, the proposed method is the first method to quantify geotechnical site similarity with explicit consideration of 2D spatial variability. It tackled the challenges of sparse investigation data, non-stationary spatial variability, and inconsistent spatial dimensions of different 2D cross-sections. A unified representation framework of 2D spatial variability using DCT-based ACF was developed. Cross-sectional similarity was quantified by the similarity of DCT-based ACFs between cross-sections. For a given 2D cross-section, BCS was adopted to approximate the DCT spectrum directly from sparse investigation data. The associated statistical uncertainty was also quantified by simulation of many random samples of the DCT spectrum. Samples of DCT-based ACF were then calculated using random samples of DCT spectrum and the newly derived equations. Then, cross-sectional similarity was quantified by the cosine similarity of DCT-based ACFs between two cross-sections. Theoretical derivation in the Appendix suggested that DCT-based ACF is an effective surrogate to represent 2D cross-sectional spatial variability from a spectral perspective and enables direct pattern comparison between cross-sections with different spatial dimensions.

Numerical examples of three soil property cross-sections were provided to illustrate the performance of the proposed method. The similarity between any two of the three cross-sections was quantified. Results indicated that the proposed method rationally quantifies the cross-sectional similarity and associated statistical uncertainty from sparse data in a data-driven manner. The quantified similarity values converged to the true value when the number of measured data increases. Real data examples from New Zealand were also used to demonstrate the application of the proposed method. High cross-sectional similarity was obtained between two sites which are approximately 700m apart. The proposed method also suggested low

731  similarity between two sites which are 5km apart. The similarity quantification developed in

732  this study assists engineering geologists and geotechnical engineers with an efficient

733  identification of similar project sites and informed decision-making in site characterization.

734  Geotechnical experiences may be effectively shared between the identified similar sites. In

735  practice, the proposed method might also be implemented with masked geographical

736  coordinate information, since such information may be restricted for confidentiality reasons.

737  The proposed method can identify a quasi-regional cluster of CPT soundings in a global

738  database that is more relevant to understanding the spatial variability of a soil property at a

739  target site. It relieves the engineer from sole reliance on subjective judgment and tedious

740  manual visual inspection to complete the same task.

## Appendix: Derivation of Equation (6)

742  Equation (3) suggests that the patterns of basis function $B_{t,s}(x_1, x_2)$ and corresponding weight

743  coefficient $\omega_{t,s}^{2D}$ jointly control the patterns of $F(x_1, x_2)$. To investigate the effect of basis

744  function $B_{t,s}(x_1, x_2)$ patterns on the $F(x_1, x_2)$ patterns when DCT basis function is adopted,

745  the ACF of $B_{t,s}(x_1, x_2)$ is firstly derived based on Equations (1), (2), and (5) (e.g., Shinozuka

746  and Deodatis, 1991; Vanmarcke, 2010):

$$
\begin{aligned}
&\mathrm{ACF}\big[B_{t,s}(x_1, x_2), \tau_1, \tau_2\big] \\
&= \frac{E\left\{\left[B_{t,s}^{2D}(x_1, x_2) - \mu_{B_{t,s}^{2D}(x_1,x_2)}\right]\left[B_{t,s}^{2D}(x_1+\tau_1, x_2+\tau_2) - \mu_{B_{t,s}^{2D}(x_1,x_2)}\right]\right\}}{\sigma^2_{B_{t,s}^{2D}(x_1,x_2)}}
\end{aligned}
\tag{A.1}
$$

747  in which $\mu_{B_{t,s}^{2D}(x_1,x_2)}$ and $\sigma_{B_{t,s}^{2D}(x_1,x_2)}$ are the mean value and standard deviation of $B_{t,s}(x_1, x_2)$,

748  respectively. It is seen from Equation (1) and Figure 4 that, when $t = s = 1$, $\mu_{B_{1,1}^{2D}(x_1,x_2)}$ is a

749  constant function, for which ACF is undefined. For $t \geq 2$, the mean value $\mu_{B_t(x)}$ is zero

750  because of the nature of cosine function:

$$\int_0^1 cos(t\pi)idi = 0 \quad (for\ t \in \mathbb{Z}, t \geq 2, and\ i \in (0,1)) \tag{A.2}$$

751   Therefore, for $t, s \neq 1$ concurrently, the $\mu_{B_{t,s}^{2D}(x_1,x_2)}$ is reduced to zero according to the

752   definition of $B_{t,s}(x_1, x_2)$ in Equation (2). Equation (A.1) is then rewritten as:

$$ACF[B_{t,s}(x_1,x_2), \tau_1, \tau_2]$$
$$= \frac{1}{\sigma_{B_{t,s}^{2D}(x_1,x_2)}^2} E\{[B_{t,s}^{2D}(x_1,x_2)][B_{t,s}^{2D}(x_1+\tau_1, x_2+\tau_2)]\} \tag{A.3}$$

753   Since the 2D basis function $B_{t,s}^{2D}(x_1, x_2)$ is constructed by a tensor product of two 1D DCT

754   basis functions along two orthonormal directions, respectively (see Equation (2)), two 1D DCT

755   basis functions are independent of each other. $B_{t,s}^{2D}(x_1, x_2)$ and $B_{t,s}^{2D}(x_1+\tau_1, x_2+\tau_2)$ hence can

756   be factorized, and Equation (A.3) is rewritten as:

$$ACF[B_{t,s}(x_1,x_2), \tau_1, \tau_2]$$
$$= \frac{1}{\sigma_{B_{t,s}(x_1,x_2)}^2} E[B_t(x_1)B_t(x_1+\tau_1)]E[B_s(x_2)B_s(x_2+\tau_2)] \tag{A.4}$$

757   The first expectation term in Equation (A.4) for $x_1$ direction can be derived, after substituting

758   1D DCT basis function from Equation (1) to Equation (A.4):

$$E[B_t(x_1)B_t(x_1+\tau_1)]$$
$$= E\left[\sqrt{\frac{2}{N_1}}cos\pi(t-1)\frac{(x_1-0.5)}{N_1}\sqrt{\frac{2}{N_1}}cos\pi(t-1)\frac{(x_1+\tau_1-0.5)}{N_1}\right] \tag{A.5}$$

759   Using product-to-sum identity, the product of two cosine functions can be rewritten as:

$$E[B_t(x_1)B_t(x_1 + \tau_1)]$$
$$= \frac{1}{N_1}\left\{E\left[cos\pi(t-1)\frac{(2x_1 + \tau_1 - 1)}{N_1}\right] + E\left[cos\pi(t-1)\frac{\tau_1}{N_1}\right]\right\} \tag{A.6}$$

760     The first expectation term with index $x_1$ is reduced to zero after expectation operation for $x_1$.

761     Therefore, Equation (A.6) is derived as a cosine function of $\tau_1$ at frequency $\pi(t-1)$:

$$E[B_t(x_1)B_t(x_1 + \tau_1)] = \frac{1}{N_1}cos\pi(t-1)\frac{\tau_1}{N_1} \tag{A.7}$$

762     In a similar fashion, the second expectation in Equation (A.4) for $x_2$ direction can be derived

763     as:

$$E[B_s(x_2)B_s(x_2 + \tau_2)] = \frac{1}{N_2}cos\pi(s-1)\frac{\tau_2}{N_2} \tag{A.8}$$

764     Therefore, substituting Equations (A.7) and (A.8) into Equation (A.4) leads to:

$$ACF\left[B_{t,s}(x_1, x_2), \tau_1, \tau_2\right]$$
$$= \frac{1}{\sigma_{B_{t,s}(x_1,x_2)}^2}E[B_t(x_1)B_t(x_1 + \tau_1)]E[B_s(x_2)B_s(x_2 + \tau_2)] \tag{A.9}$$
$$= \frac{1}{\sigma_{B_{t,s}(x_1,x_2)}^2}\frac{1}{N_1 N_2}cos\pi(t-1)\frac{\tau_1}{N_1}cos\pi(s-1)\frac{\tau_2}{N_2}$$

765     Note that the variance term $\sigma_{B_{t,s}(x_1,x_2)}^2$ can be derived based on the fact that the 2D DCT basis

766     function is zero-mean and orthonormal (e.g., Rao and Yip, 1990). The Frobenius norm of 2D

767     basis function $B_{t,s}(x_1, x_2)$ is unity:

$$\left\|B_{t,s}(x_1, x_2)\right\| = \sqrt{B_{t,s}(1,1)^2 + B_{t,s}(2,1)^2 + \cdots + B_{t,s}(N_1, N_2)^2} = 1 \tag{A.10}$$

768     According to the definition of variance, the $\sigma_{B_{t,s}(x_1,x_2)}^2$ is derived as:

$$\sigma^2_{B_{t,s}(x_1,x_2)} = \frac{1}{N_1 N_2}\left[B_{t,s}(1,1)^2 + B_{t,s}(2,1)^2 + \cdots + B_{t,s}(N_1,N_2)^2\right] = \frac{1}{N_1 N_2} \tag{A.11}$$

769 Therefore, substituting Equation (A.11) into Equation (A.9) leads to the normalized ACF of

770 2D basis function:

$$
\begin{aligned}
&ACF\left[B_{t,s}(x_1,x_2),\tau_1,\tau_2\right]\\
&= \frac{1}{\sigma^2_{B_{t,s}(x_1,x_2)}} E[B_t(x_1)B_t(x_1+\tau_1)]E[B_s(x_2)B_s(x_2+\tau_2)]\\
&= cos\pi(t-1)\frac{\tau_1}{N_1}cos\pi(s-1)\frac{\tau_2}{N_2} \qquad (t,s \neq 1\ concurrently)
\end{aligned} \tag{A.12}
$$

771 Equation (A.12) shows that the ACF of a 2D DCT basis function $B_{t,s}(x_1,x_2)$ is derived as the

772 product of two cosine functions along $x_1$ and $x_2$ directions, respectively. These two cosine

773 functions are of lag distances $\tau_1$ and $\tau_2$, respectively. Note that the indexes $x_1$ and $x_2$ are

774 eliminated in Equation (A.12).

775       Based on Equations (2), (3), and (5), the ACF of 2D spatial variability $F(x_1,x_2)$ can

776 also be derived under DCT framework:

$$
\begin{aligned}
&ACF[F(x_1,x_2),\tau_1,\tau_2]\\
&= \frac{E\left\{\left[F(x_1,x_2)-\mu_{F(x_1,x_2)}\right]\left[F(x_1+\tau_1,x_2+\tau_2)-\mu_{F(x_1,x_2)}\right]\right\}}{\sigma^2_{F(x_1,x_2)}}
\end{aligned} \tag{A.13}
$$

777 in which $\mu_{F(x_1,x_2)}$ and $\sigma_{F(x_1,x_2)}$ are mean value and variance of $F(x_1,x_2)$. Note that $\mu_{F(x_1,x_2)}$

778 can be replaced by the contribution of $B_{1,1}(x_1,x_2)$, since $B_{1,1}(x_1,x_2)$ is the only 2D DCT basis

779 function with a non-zero mean value (see upper left 2D basis function in Figure 4). Substituting

780 Equation (3) into Equation (5) leads to:

$$ACF[F(x_1, x_2), \tau_1, \tau_2]$$

$$= \frac{E\left\{\left[\sum_{t=1}^{N_1}\sum_{s=1}^{N_2}\omega_{t,s}^{2D}B_{t,s}(x_1,x_2)\right]\left[\sum_{t'=1}^{N_1}\sum_{s'=1}^{N_2}\omega_{t',s'}^{2D}B_{t',s'}(x_1+\tau_1,x_2+\tau_2)\right]\right\}}{\sigma_{F(x_1,x_2)}^2} \tag{A.14}$$

$$(t, s \neq 1\ concurrently;\ t', s' \neq 1\ concurrently)$$

781   The prime symbols are used to distinguish two factors along an individual direction. By

782   utilizing the orthonormal property of $B_{t,s}(x_1, x_2)$, Equation (A.14) can be expanded and

783   rearranged as a multiple summation:

$$ACF[F(x_1, x_2), \tau_1, \tau_2]$$

$$= \frac{1}{\sigma_{F(x_1,x_2)}^2}\sum_{t=1}^{N_1}\sum_{s=1}^{N_2}\sum_{t'=1}^{N_1}\sum_{s'=1}^{N_2}\omega_{t,s}^{2D}\omega_{t',s'}^{2D}E[B_t(x_1)B_{t'}(x_1+\tau_1)]E[B_s(x_2)B_{s'}(x_2+\tau_2)] \tag{A.15}$$

$$(t, s \neq 1\ concurrently;\ t', s' \neq 1\ concurrently)$$

784   The expectation operations are performed for $x_1$ and $x_2$ directions, separately. For $x_1$ direction,

785   after substituting 1D DCT basis function in Equation (1), the expectation operation

786   $E[B_t(x_1)B_{t'}(x_1+\tau_1)]$ can be further expressed as:

$$E[B_t(x_1)B_{t'}(x_1+\tau_1)]$$

$$= \begin{cases} E\left[\sqrt{\frac{2}{N_1}}\sqrt{\frac{2}{N_1}}\cos\pi\dfrac{(t'-1)(x_1+\tau_1-0.5)}{N_1}\right] = 0 & when\ t = 1, t' \neq 1 \\[3mm] E\left[\sqrt{\frac{2}{N_1}}\sqrt{\frac{2}{N_1}}\cos\pi\dfrac{(t-1)(x_1-0.5)}{N_1}\right] = 0 & when\ t \neq 1, t' = 1 \\[3mm] E\left[\sqrt{\frac{2}{N_1}}\sqrt{\frac{2}{N_1}}\cos\pi\dfrac{(t-1)(x_1-0.5)}{N_1}\cos\pi\dfrac{(t'-1)(x_1+\tau_1-0.5)}{N_1}\right] & when\ t \neq 1,\ t' \neq 1 \end{cases} \tag{A.16}$$

787   Equation (A.16) consists of three scenarios, i.e., when $t = 1, t' \neq 1$; $t \neq 1, t' = 1$; and $t \neq$

788   $1,\ t' \neq 1$. Since the first two scenarios lead to single cosine functions of $x_1$, the associated

789   terms are reduced to zeros after expectation operation on $x_1$. When $t \neq 1,\ t' \neq 1$, the

790     expectation of product of two cosine functions yields zero when the frequencies of two basis

791     function are not equal, i.e., $t \neq t'$. Only the product of two cosine functions with equal

792     frequencies remain, i.e., $t = t'$:

$$E[B_t(x_1)B_{t'}(x_1+\tau_1)] =$$
$$= \begin{cases} 0 & when\ t \neq t' \\ \dfrac{2}{N_1}E\left[cos\pi\dfrac{(t-1)(x_1-0.5)}{N_1}cos\pi\dfrac{(t-1)(x_1+\tau_1-0.5)}{N_1}\right] & when\ t = t' \end{cases} \tag{A.17}$$

793     When $t = t'$, the product of two cosine functions can be rewritten using product-to-sum

794     identity as:

$$E[B_t(x_1)B_{t'}(x_1+\tau_1)]$$
$$= \frac{1}{N_1}E\left[cos\pi\frac{2(t-1)(x_1-0.5)+(t-1)\tau_1}{N_1} + cos\pi(t-1)\frac{\tau_1}{N_1}\right] \tag{A.18}$$

795     Similarly, in Equation (A.18), the first cosine function term reduces to zeros after expectation

796     operation on $x_1$. Therefore, Equation (A.18) is then derived as:

$$E[B_t(x_1)B_{t'}(x_1+\tau_1)] = \frac{1}{N_1}cos\pi(t-1)\frac{\tau_1}{N_1} \qquad when\ t = t' \tag{A.19}$$

797     In a similar fashion, the second expectation term in Equation (A.15) for $x_2$ direction is derived

798     as:

$$E[B_s(x_2)B_{s'}(x_2+\tau_2)] = \frac{1}{N_2}cos\pi(s-1)\frac{\tau_2}{N_2} \qquad when\ s = s' \tag{A.20}$$

799     After substituting Equations (A.19) and (A.20) into Equation (A.15), Equation (A.15) can be

800     rewritten as:

$$\text{ACF}[F(x_1, x_2), \tau_1, \tau_2]$$

$$= \frac{1}{N_1 N_2} \frac{1}{\sigma^2_{F(x_1,x_2)}} \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D^2} \cos\pi(t-1)\frac{\tau_1}{N_1} \cos\pi(s-1)\frac{\tau_2}{N_2} \tag{A.21}$$

$$(t, s \neq 1 \; concurrently)$$

801　The variance term $\sigma^2_{F(x_1,x_2)}$ can also be derived since the 2D DCT basis function $B_{t,s}(x_1, x_2)$ is

802　orthonormal and independent of each other. By definition, $\sigma^2_{F(x_1,x_2)}$ is expressed as:

$$\sigma^2_{F(x_1,x_2)} = E\left\{\left[F(x_1, x_2) - \mu_{F(x_1,x_2)}\right]^2\right\} = \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D^2} \sigma^2_{B_{t,s}(x_1,x_2)} \tag{A.22}$$

803　Note that in Equation (A.11), $\sigma^2_{B_{t,s}(x_1,x_2)}$ is derived as $\frac{1}{N_1 N_2}$. Therefore, Equation (A.22) can be

804　rewritten as

$$\sigma^2_{F(x_1,x_2)} = \frac{1}{N_1 N_2} \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D^2} \tag{A.23}$$

805　Combining Equations (A.23) and (A.21) yields the normalized DCT-based ACF of $F(x_1, x_2)$,

806　which is provided as Equation (6) in the main text:

$$\text{ACF}[F(x_1, x_2), \tau_1, \tau_2]$$

$$= \frac{1}{\sum_{t=1}^{N_1}\sum_{s=1}^{N_2} \omega_{t,s}^{2D^2}} \sum_{t=1}^{N_1} \sum_{s=1}^{N_2} \omega_{t,s}^{2D^2} \text{ACF}[B_{t,s}(x_1, x_2), \tau_1, \tau_2] \tag{A.24}$$

$$(t, s \neq 1 \; concurrently)$$

807

# References

808

Baecher, G.B. & Christian, J.T. 2003. Reliability and statistics in geotechnical engineering. John Wiley & Sons.

Blumensath, T. & Davies, M. 2006. Sparse and shift-Invariant representations of music. IEEE Transactions on Audio, Speech, and Language Processing, 14, 50-57, doi: 10.1109/TSA.2005.860346.

Brockwell, P.J., Davis, R.A. 1991. Stationary ARMA Processes. In: Time Series: Theory and Methods. Springer Series in Statistics. Springer, New York, NY. https://doi.org/10.1007/978-1-4419-0320-4_3

Candès, E.J., Romberg, J.K. & Tao, T. 2006. Stable signal recovery from incomplete and inaccurate measurements. Communications on pure and applied mathematics, 59, 1207-1223.

Candès, E.J. & Wakin, M.B. 2008. An introduction to compressive sampling. IEEE Signal Processing Magazine, 25, 21-30.

Ching, J. & Phoon, K.-K. 2017. Characterizing Uncertain Site-Specific Trend Function by Sparse Bayesian Learning. Journal of Engineering Mechanics, 143, 04017028, doi: doi:10.1061/(ASCE)EM.1943-7889.0001240.

Ching, J., Phoon, K.-K., Beck, J.L. & Huang, Y. 2017. Identifiability of Geotechnical Site-Specific Trend Functions. ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering, 3, 04017021, doi: doi:10.1061/AJRUA6.0000926.

Ching, J., Huang, W.-H. & Phoon, K.-K. 2020. 3D Probabilistic Site Characterization by Sparse Bayesian Learning. Journal of Engineering Mechanics, 146, 04020134, doi: doi:10.1061/(ASCE)EM.1943-7889.0001859.

47

831   Ching, J. & Phoon, K.-K. 2020. Measuring Similarity between Site-Specific Data and Records

832       from Other Sites. ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems,

833       Part A: Civil Engineering, 6, 04020011, doi: doi:10.1061/AJRUA6.0001046.

834   Ching, J., Uzielli, M., Phoon, K.-K. & Xu, X. 2023. Characterization of Autocovariance

835       Parameters of Detrended Cone Tip Resistance from a Global CPT Database. Journal of

836       Geotechnical   and   Geoenvironmental   Engineering,   149,   04023090,   doi:

837       doi:10.1061/JGGEFK.GTENG-11214.

838   Ching, J., Yoshida, I. & Phoon, K.-K. 2022. Comparison of trend models for geotechnical

839       spatial variability: Sparse Bayesian Learning vs. Gaussian Process Regression. Gondwana

840       Research, 123, pp.174-183.

841   Dai, H., Zhang, R. & Beer, M. 2022. A new perspective on the simulation of cross-correlated

842       random      fields.      Structural      Safety,      96,      102201,      doi:

843       https://doi.org/10.1016/j.strusafe.2022.102201.

844   Dong, Y., Sun, Z. & Jia, H. 2006. A cosine similarity-based negative selection algorithm for

845       time series novelty detection. Mechanical Systems and Signal Processing, 20, 1461-1472,

846       doi: https://doi.org/10.1016/j.ymssp.2004.12.006.

847   Donoho, D.L. 2006. Compressed sensing. IEEE Transactions on Information Theory, 52, 1289-

848       1306.

849   Einsele, G., Chough, S.K. & Shiki, T. 1996. Depositional events and their records—an

850       introduction. Sedimentary Geology, 104, 1-9, doi: https://doi.org/10.1016/0037-

851       0738(95)00117-4.

852   Fenton, G.A. 1999a. Estimation for Stochastic Soil Models. Journal of Geotechnical and

853       Geoenvironmental   Engineering,   125,   470-485,   doi:   doi:10.1061/(ASCE)1090-

854       0241(1999)125:6(470).

855    Fenton, G.A. 1999b. Random Field Modeling of CPT Data. Journal of Geotechnical and

856        Geoenvironmental Engineering, 125, 486-498, doi: doi:10.1061/(ASCE)1090-

857        0241(1999)125:6(486).

858    Guan, Z. & Wang, Y., 2020. Statistical charts for determining sample size at various levels of

859        accuracy and confidence in geotechnical site investigation. Géotechnique, 70(12), 1145-

860        1159.

861    Guan, Z. & Wang, Y. 2023. Data-driven simulation of two-dimensional cross-correlated

862        random fields from limited measurements using joint sparse representation. Reliability

863        Engineering    &    System    Safety,    238,    109408,    doi:

864        https://doi.org/10.1016/j.ress.2023.109408.

865    Guan, Z., Wang, Y. and Phoon, K.K., 2023a. Fusion of Sparse Non-co-located Measurements

866        from Multiple Sources for Geotechnical Site Investigation. Canadian Geotechnical Journal,

867        doi: https://doi.org/10.1139/cgj-2023-0289.

868    Guan, Z., Wang, Y. & Phoon, K.-K. 2023b. Dictionary learning of sparse ground investigation

869        data from a specific site and existing database from other sites. (Manuscript in preparation).

870    Han, L., Wang, L., Ding, X., Wen, H., Yuan, X. & Zhang, W. 2022. Similarity quantification

871        of soil parametric data and sites using confidence ellipses. Geoscience Frontiers, 13,

872        101280, doi: https://doi.org/10.1016/j.gsf.2021.101280.

873    Hu, Y., Zhao, T., Wang, Y., Choi, C. & Ng, C.W.W. 2019. Direct simulation of two-

874        dimensional isotropic or anisotropic random field from sparse measurement using Bayesian

875        compressive sampling. Stochastic Environmental Research and Risk Assessment, 33,

876        1477-1496, doi: 10.1007/s00477-019-01718-7.

877    Hu, Y. and Wang, Y., 2024. Evaluating statistical homogeneity of cone penetration test (CPT)

878        data profile using auto-correlation function. Computers and Geotechnics, 165, p.105852.

879 Huang, S.P., Quek, S.T. & Phoon, K.K. 2001. Convergence study of the truncated Karhunen–
880    Loeve expansion for simulation of stochastic processes. International Journal for
881    Numerical Methods in Engineering, 52, 1029-1043, doi: https://doi.org/10.1002/nme.255.

882 Itskov, M. 2007. Tensor algebra and tensor analysis for engineers. Springer.

883 Juang, C.H., Zhang, J., Shen, M. & Hu, J. 2019. Probabilistic methods for unified treatment of
884    geotechnical and geological uncertainties in a geotechnical analysis. Engineering Geology,
885    249, 148-161, doi: https://doi.org/10.1016/j.enggeo.2018.12.010.

886 Katsman, R. & Painuly, A. 2022. Influence of anisotropy in mechanical properties of muddy
887    aquatic sediment on CH4 bubble growth direction and migration pattern. Engineering
888    Geology, 299, 106565, doi: https://doi.org/10.1016/j.enggeo.2022.106565.

889 Leung, A.Y. 2023. Soil-Structure Interaction in Spatially Variable Ground. In Uncertainty,
890    Modeling, and Decision Making in Geotechnics (pp. 427-440). CRC Press.

891 Lyu, B., Hu, Y. & Wang, Y. 2023. Data-Driven Development of Three-Dimensional
892    Subsurface Models from Sparse Measurements Using Bayesian Compressive Sampling: A
893    Benchmarking Study. ASCE-ASME Journal of Risk and Uncertainty in Engineering
894    Systems, Part A: Civil Engineering, 9, 04023010.

895 Müller, S., Schüler, L., Zech, A. and Heße, F., 2022. GSTools v1. 3: a toolbox for geostatistical
896    modelling in Python. Geoscientific Model Development, 15(7), pp.3161-3182.

897 Nguyen, H.V. & Bai, L. 2011. Cosine Similarity Metric Learning for Face Verification. In:
898    Kimmel, R., Klette, R. & Sugimoto, A. (eds.) Computer Vision – ACCV 2010. Springer
899    Berlin Heidelberg, Berlin, Heidelberg, 709-720.

900 NZGD. 2023. New Zealand Geotechnical Database. In: (EQC), N.Z.E.C. (ed.),
901    https://www.nzgd.org.nz (accessed 3 October 2023)

902 Onyejekwe, S., Kang, X. & Ge, L. 2016. Evaluation of the scale of fluctuation of geotechnical

903       parameters by autocorrelation function and semivariogram function. Engineering Geology,

904       214, 43-49, doi: https://doi.org/10.1016/j.enggeo.2016.09.014.

905 Pati, Y.C., Rezaiifar, R. & Krishnaprasad, P. 1993. Orthogonal matching pursuit: Recursive

906       function approximation with applications to wavelet decomposition. *Signals, Systems and*

907       *Computers, 1993. 1993 Conference Record of The Twenty-Seventh Asilomar Conference*

908       *on*. IEEE, 40-44.

909 Phoon, K.-K. & Ching, J. 2022. Additional observations on the site recognition challenge.

910       Journal of GeoEngineering, 17, 231-247.

911 Phoon, K.-K., Ching, J. & Shuku, T. 2022. Challenges in data-driven site characterization.

912       Georisk: Assessment and Management of Risk for Engineered Systems and Geohazards,

913       16, 114-126, doi: 10.1080/17499518.2021.1896005.

914 Phoon, K.-K. & Fenton, G.A. 2004. Estimating sample autocorrelation functions using

915       bootstrap. *Ninth ASCE Specialty Conference on Probabilistic Mechanics and Structural*

916       *Reliability*, Albuquerque, New Mexico.

917 Phoon, K.-K. & Kulhawy, F.H. 1999. Characterization of geotechnical variability. Canadian

918       Geotechnical Journal, 36, 612-624.

919 Phoon, K.-K. & Zhang, W. 2023. Future of machine learning in geotechnics. Georisk:

920       Assessment and Management of Risk for Engineered Systems and Geohazards, 17, 7-22,

921       doi: 10.1080/17499518.2022.2087884.

922 Priestley, M.B. 1981. Spectral analysis and time series: probability and mathematical statistics.

923 Rafiee, J. & Tse, P.W. 2009. Use of autocorrelation of wavelet coefficients for fault diagnosis.

924       Mechanical Systems and Signal Processing, 23, 1554-1572, doi:

925       https://doi.org/10.1016/j.ymssp.2009.02.008.

926 Rao, K.R. & Yip, P. 1990. Discrete cosine transform: algorithms, advantages, applications.

927     Academic Press Professional, Inc.

928 Shalvi, O. & Weinstein, E. 1996. System identification using nonstationary signals. IEEE

929     Transactions on Signal Processing, 44, 2055-2063, doi: 10.1109/78.533725.

930 Shannon, C.E. 1948. A mathematical theory of communication. The Bell System Technical

931     Journal, 27, 379-423, doi: 10.1002/j.1538-7305.1948.tb01338.x.

932 Sharma, A., Ching, J. & Phoon, K.-K. 2022. A Hierarchical Bayesian Similarity Measure for

933     Geotechnical Site Retrieval. Journal of Engineering Mechanics, 148, 04022062.

934 Shechtman, E. & Irani, M. 2007. Matching Local Self-Similarities across Images and Videos.

935     *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 1-8.

936 Shi, C. & Wang, Y. 2021a. Development of subsurface geological cross-section from limited

937     site-specific boreholes and prior geological knowledge using iterative convolution

938     XGBoost. Journal of Geotechnical and Geoenvironmental Engineering, 147, 04021082.

939 Shi, C. & Wang, Y. 2021b. Training image selection for development of subsurface geological

940     cross-section by conditional simulations. Engineering Geology, 295, 106415.

941 Shi, C., Wang, Y. & Kamchoom, V. 2023. Data-driven multi-stage sampling strategy for a

942     three-dimensional geological domain using weighted centroidal voronoi tessellation and

943     IC-XGBoost3D.       Engineering       Geology,       325,       107301,       doi:

944     https://doi.org/10.1016/j.enggeo.2023.107301.

945 Shinozuka, M. & Deodatis, G. 1991. Simulation of Stochastic Processes by Spectral

946     Representation. Applied Mechanics Reviews, 44, 191-204, doi: 10.1115/1.3119501.

947 Shuku, T., Phoon, K.-K. & Yoshida, I. 2020. Trend estimation and layer boundary detection in

948     depth-dependent soil data using sparse Bayesian lasso. Computers and Geotechnics, 128,

949     103845, doi: https://doi.org/10.1016/j.compgeo.2020.103845.

950     Simakov, D., Caspi, Y., Shechtman, E. & Irani, M. 2008. Summarizing visual data using

951        bidirectional similarity. *2008 IEEE Conference on Computer Vision and Pattern*

952        *Recognition*, 1-8.

953     Tipping, M.E. 2001. Sparse Bayesian learning and the relevance vector machine. Journal of

954        machine learning research, 1, 211-244.

955     Vanmarcke, E. 2010. Random Fields. World Scientific.

956     Wallace, G.K. 1992. The JPEG still picture compression standard. IEEE Transactions on

957        Consumer Electronics, 38, xviii-xxxiv, doi: 10.1109/30.125072.

958     Wang, Y., Hu, Y. & Phoon, K.-K. 2022. Non-parametric modelling and simulation of

959        spatiotemporally varying geo-data. Georisk: Assessment and Management of Risk for

960        Engineered Systems and Geohazards, 16, 77-97, doi: 10.1080/17499518.2021.1971258.

961     Wang, Y. & Zhao, T. 2016. Interpretation of soil property profile from limited measurement

962        data: a compressive sampling perspective. Canadian Geotechnical Journal, 53, 1547-1559.

963     Wang, Y., Zhao, T., Hu, Y. & Phoon, K.-K. 2019. Simulation of Random Fields with Trend

964        from Sparse Measurements without Detrending. Journal of Engineering Mechanics, 145,

965        04018130, doi: doi:10.1061/(ASCE)EM.1943-7889.0001560.

966     Webster, R. & Oliver, M.A. 2007. Geostatistics for environmental scientists. John Wiley &

967        Sons.

968     Wen, Y.K. & Gu, P. 2004. Description and Simulation of Nonstationary Processes Based on

969        Hilbert Spectra. Journal of Engineering Mechanics, 130, 942-951, doi:

970        doi:10.1061/(ASCE)0733-9399(2004)130:8(942).

971     Xu, J., Wang, Y. & Zhang, L. 2021. Interpolation of extremely sparse geo-data by data fusion

972        and collaborative Bayesian compressive sampling. Computers and Geotechnics, 134,

973        104098.

974 Yang, H.Q., Zhang, L., Gao, L., Phoon, K.K. and Wei, X., 2022. On the importance of landslide

975 management: Insights from a 32-year database of landslide consequences and rainfall in

976 Hong Kong. Engineering Geology, 299, p.106578.

977 Yoshida, I., Tomizawa, Y. & Otake, Y. 2021. Estimation of trend and random components of

978 conditional random field using Gaussian process regression. Computers and Geotechnics,

979 136, 104179, doi: https://doi.org/10.1016/j.compgeo.2021.104179.

980 Zhang, Y., Li, Y.E. & Ku, T. 2021. Soil/rock interface profiling using a new passive seismic

981 survey: Autocorrelation seismic interferometry. Tunnelling and Underground Space

982 Technology, 115, 104045, doi: https://doi.org/10.1016/j.tust.2021.104045.

983 Zhao, C., Gong, W., Juang, C.H., Tang, H., Hu, X., Wang, L., 2023. Optimization of site

984 exploration program based on coupled characterization of stratigraphic and geo-properties

985 uncertainties. Engineering Geology, 317, p.107081.

986 Zhao, T., Hu, Y. & Wang, Y. 2018. Statistical interpretation of spatially varying 2D geo-data

987 from sparse measurements using Bayesian compressive sampling. Engineering Geology,

988 246, 162-175, doi: https://doi.org/10.1016/j.enggeo.2018.09.022.

989 Zhao, T. & Wang, Y. 2020. Non-parametric simulation of non-stationary non-gaussian 3D

990 random field samples directly from sparse measurements using signal decomposition and

991 Markov Chain Monte Carlo (MCMC) simulation. Reliability Engineering & System Safety,

992 203, 107087.

993