

1 Prediction of interface of geological formations using generalized additive model

2 Xiaohui Qi^a, Zhiyong Yang^{b*}, Jian Chu^c

3 ^aLecturer, Department of Mechanical and Construction Engineering, Northumbria University,
4 Newcastle upon Tyne, UK. Email: xiaohui.qi@northumbria.ac.uk.

5 ^bAssistant Professor, School of Civil Engineering, Sun Yat-Sen University, Guangzhou, China.
6 Email: yangzhy85@mail.sysu.edu.cn.

7 ^cProfessor, School of Civil and Environmental Engineering, Nanyang Technological University,
8 Singapore. Email: cjchu@ntu.edu.sg.

9 *Corresponding author

10 **Abstract:** Geological information such as geological interfaces is important for the design of
11 underground excavation and supporting measures. This in turn requires a method to predict
12 accurately the locations of geological interfaces for the gap areas between boreholes. This study
13 presents a generalized additive model (GAM) to predict the location of the geological interfaces.
14 The performance of the GAM method is evaluated using both simulated data and borehole data
15 for the determination of rockhead in two different geological formations in Singapore. The
16 results show that the GAM method can provide a reasonable confidence interval (CI) of the
17 mean trend and the prediction interval (PI) in the sense that the 95% CI covers about 95% of
18 the actual mean curve while the 95% PI covers around 95% of testing data. Furthermore, the
19 geological complexity can be well reflected as the prediction uncertainty in the geologically
20 complex area is larger than that in the geologically regular area. More importantly, the users
21 can impose prior information or personal judgment regarding the shape of the geological profile
22 on the model. This is an important feature to enable further improvement in the accuracy of the
23 prediction.

24 **Keywords:** spatial prediction, geological interface, generalized additive model, cubic spline

25

26 **1. Introduction**

27 Geological model including the interfaces of different geological formations is indispensable
28 information for underground constructions as it may affect the construction method and
29 supporting measures. It is necessary to predict the location of geological interfaces in the gap
30 areas between boreholes. This task is difficult as the site exploration data are always sparse and
31 limited. What makes this task more challenging is the large variability in the geological
32 interface, especially the rockhead in rock formations. The reason is that the weathering of the
33 rock is affected by many factors, including climate, topography, hydrological conditions,
34 biological systems, rock mass discontinuities, rock composition and permeability (Zhao et al.
35 1994). It is vital to find an effective method to accurately predict the location of geological
36 interfaces. The prediction should be able to provide a predicted value of the location of the
37 geological interface as well as its uncertainty.

38 A variety of methods have been used for interpolation problems in geotechnical or
39 geological engineering. These methods can be divided into two categories, deterministic
40 methods and statistical methods. Deterministic methods such as the inverse distance weighting
41 method, spline interpolation or the triangle-based tessellation method (e.g., Aswar and
42 Ullagaddi 2017; Burke et al. 2017) cannot automatically quantify the uncertainty of the
43 prediction or provide any confidence interval of the predicted property. This problem can be to
44 some extent addressed using cross-validations, as shown in Lark et al. (2013), but the results
45 highly depend on the employed testing data and the quantified uncertainty may not be reliable
46 when the testing data are limited. The statistical interpolation methods include the coupled
47 Markov chain method (Qi et al. 2016; Li et al. 2019; Liu et al. 2020), Markov random field

48 method (Wang et al. 2017, 2018), Bayesian compressive sampling method (Wang and Zhao
49 2016, 2017; Zhao, Hu, and Wang 2018), random field method (Gong et al. 2020; Zhao et al.
50 2021), geostatistical methods such as kriging and conditional random field method (Qi et al.
51 2019, 2021a). The coupled Markov chain method can characterize the geological uncertainty
52 using limited borehole data, but it can only be used when the transition of geological types has
53 a Markovian property. The Markov random field method can model complex geological
54 structures, but some of its parameters lack clear physical meaning (Mariethoz and Caers 2014).
55 The recently developed Bayesian compressive sampling method can quantify the interpolation
56 uncertainty using limited data and has a high interpolation accuracy (Wang, Akeju, and Zhao
57 2017). Moreover, it is quite versatile in that it can model both stationary and non-stationary
58 random fields (Wang, Zhao, and Phoon 2018; Wang et al. 2019). One potential problem of the
59 method is that its robustness degrades when the number of measurements is smaller than the
60 length of the discrete signal or when the measurement noise is relatively large (Huang et al.
61 2014). Geostatistical methods such as kriging or conditional random field have gained wide
62 popularity (e.g., Qi et al. 2019, 2021a). One problem with geostatistical methods is that they
63 are purely mathematic based and may not lead to realistic soil or geological profiles. For
64 example, the conditional random field or the kriging method normally produces a soil or
65 geological profile with extreme values only at known data points. Furthermore, some artificial
66 intelligence methods such as neural networks (Zhou and Wu 1994) and the support vector
67 machine (Smirnov, Boisvert, and Paradis 2008) were also applied to spatial prediction
68 problems of geological conditions. The drawback of these methods is that they lack
69 interpretability in the sense that they behave like black boxes and the effect of individual

70 explanatory variables is difficult to examine.

71 Recently, Qi et al. (2020, 2021b) applied a spline regression method to spatial predictions
72 of the location of geological interfaces. It has been shown that the method can provide a clear
73 spatial trend of the geological interface, which well reflects the geological complexity. One
74 problem of these studies is that the uncertainty in the mean trend is not well quantified or
75 distinguished from the uncertainty in the random error. Herein the uncertainty in the mean trend
76 represents the bias of the fitted curve in a regression. To be specific, if two different sets of data
77 for the same explanatory and response variables are employed to perform regression, the two
78 fittings generally produce different mean trends. The variability in the fitted curves is called
79 uncertainty in the mean trend. The uncertainty in the random error denotes the deviation of the
80 data points from the fitted curve. To address this issue, this study uses a generalized linear
81 model (GAM) to perform the spatial prediction of the geological interface. For the GAM, the
82 response variable is expressed as the weighted average of some basis functions (Wood 2017).
83 It can be viewed to be a non-parametric or semi-parametric method in the sense that the
84 structure of the model is not fixed. The uncertainty in the mean trend and random error can be
85 explicitly considered by the GAM. An additional advantage of the GAM is its interpretability,
86 which means the contribution of each independent parameter to the prediction is explicitly
87 modeled and can be readily examined.

88 In this study, the GAM is firstly briefly introduced. Secondly, the performance of the
89 GAM model is investigated based on cross-validations using simulated data. Finally, some
90 borehole data from Singapore, which reveal the rockhead elevation of two rock formations, i.e.,
91 Bukit Timah Formation and Jurong Formation, are used to predict the rockhead. Herein the

92 rockhead is the interface of soil and rock layers in a rock formation, which is mainly determined
93 according to the weathering degree of the geological layers (Qi et al. 2020, 2021b). Prediction
94 errors and the capability of GAM in characterizing the two types of uncertainty are evaluated
95 using a cross-validation procedure. The reasonableness of the prediction uncertainty, which
96 was generally ignored in existing studies, is investigated in this paper. The two types of
97 uncertainties, i.e., uncertainty in the mean trend and random error are well differentiated and
98 quantified in the investigation. The role of engineering judgment in spatial predictions is also
99 discussed in the example.

100 **2. Generalized additive model**

101 The section introduces briefly the generalized additive model to be used for the prediction of
102 the rockhead elevation. The GAM is originally developed by Hastie and Tibshirani (1986,
103 1990). It can be viewed to be a generalization of the linear regression model. The main
104 advantage of the GAM is that it can flexibly identify the nonlinear relation between explanatory
105 variables (also called predictors or covariates) and a response variable (Hastie and Tibshirani
106 1986; Wood 2017). To be specific, the users do not need to specify a particular parametric
107 function to represent the nonlinear pattern. Instead, non-parametric or semi-parametric smooth
108 functions are used to relate the predictors and responses. Another advantage of the GAM model
109 is the interpretability, which means that the effect of predictors can be examined separately and
110 explicitly (Hastie and Tibshirani 1986). The GAM has been widely applied to various
111 disciplines since its advent, such as environmental engineering (e.g., Gong et al. 2017; Ma et
112 al. 2020), soil science (e.g., de Brogniez 2015), ecology (e.g., Yee and Mitchell 1991; Simpson
113 2018), transportation engineering (e.g., Khoda Bakhshi and Ahmed 2021). In geotechnical

114 engineering, the model is applied for the determination of landslide susceptibility, as shown in
115 Goetz et al. (2015) and Bordoni et al. (2020). The spline regression methods investigated in Qi
116 et al. (2020, 2021b) are also GAM. This study extends the work in Qi et al. (2020, 2021b) by
117 taking the uncertainty in the fitted mean trend of the response variable into consideration. The
118 basic idea and the fitting of the GAM are introduced as follows.

119 **2.1 Representation of smooth functions**

120 This study intends to investigate a one-dimensional problem, namely the prediction of the
121 rockhead elevation along a line. The response variable is the rockhead elevation while the
122 explanatory variable is the distance to the leftmost point on the line. Besides, since the rockhead
123 elevation does not have any capped value, the Gaussian distribution is taken to be the
124 probability distribution of the response. In this case, the GAM can be simplified into

$$125 \quad y = f(x) + \varepsilon \quad (1)$$

126 where ε is a normally distributed random variable with a mean of 0 and variance of σ_ε^2 . The
127 smooth function is usually represented by the weighted sum of several basis functions, i.e.,

$$128 \quad f(x) = \sum_{i=1}^q b_i(x)\beta_i \quad (2)$$

129 where $b_i(x)$ is a basis function of which the expression is already known; β_i is an unknown
130 coefficient and q is the total number of basis functions, also called the dimension of basis. This
131 study adopts the commonly used cubic spline basis functions because the cubic spline
132 interpolant always provides a solution that is smooth and closely approximates to the true
133 function whatever the true function is (Wood 2017). Polynomial bases are not chosen because
134 a high-order polynomial function usually causes an oscillation problem, as shown in De Boor
135 (2001). One example of curve fitting using cubic spline basis functions is plotted in Fig. 1. In

136 Fig. 1, the 19 circles denote blow count value data from standard penetration tests taken from
 137 Baecher and Christian (2008). The vertical coordinate denotes the elevation while the
 138 horizontal coordinate denotes the blow count value. The 19 data points are fitted with three,
 139 six and twelve cubic spline basis functions using a least squared method in Figs. 1a, 1b and 1c,
 140 respectively. The basis functions are denoted as the blue dashed lines. Since the values of the
 141 basis function are too small (< 1), the basis functions are magnified by 10 times to make them
 142 discernible in the figure. Also, the locations of the knots are denoted as the red dotted lines.
 143 These knots are the connection points of two neighbouring pieces or sections of the fitted curve,
 144 each of which can be expressed by a cubic spline function. For illustration purposes, the knots
 145 are set to be evenly spaced in the elevation direction. As shown in Fig. 1, the resulting smooth
 146 function is a piecewise cubic polynomial function. The fitted curve becomes wigglier when
 147 more basis functions are used.

148 2.2 Degree of smoothing

149 After the structure of the smooth function is known, one natural question is how to determine
 150 the number and location of the knot given some data. The number and location of knots control
 151 the degree of smoothing of the resulting function. Too many knots result in a wiggly curve
 152 running across all the data points. This curve normally suffers from overfitting and performs
 153 poorly when it is used for prediction. In practice, the number and location of the knots can be
 154 determined by model selection methods or cross-validation, as shown in Qi et al. (2020, 2021b).
 155 An alternative method to control the smoothness is fixing the basis dimension at a relatively
 156 large size and adding a wiggleness penalty term in the least-squares objective (Wood 2017), i.e.,

$$157 \quad \|\sqrt{\mathbf{W}}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})\|^2 + \lambda \int [f''(x)]^2 dx \quad (3)$$

158 where $\|\cdot\|^2$ is the squared Euclidian length of a vector; $\mathbf{X}\boldsymbol{\beta}$ denotes the fitted values of the
159 smooth function in which \mathbf{X} is the model matrix denoting the values of the basis functions at
160 locations of observation data while $\boldsymbol{\beta}$ is the coefficient vector; \mathbf{W} is a diagonal matrix denoting
161 the weights of data points. Assigning a weight value of w_i for a data point is equivalent to put
162 w_i identical data points at the same location. Normally the weights for all the data points are
163 set to be 1. A value larger than 1 can be used when a data point reveals an important geological
164 feature (such as an abruptly low rockhead caused by faults) and controls the shape of the
165 geological profile. $f''(x)$ is the second derivative of the smooth function $f(x)$; λ is the
166 smoothing parameter that controls the tradeoff between the model fit and the model smoothness.
167 $\lambda = 0$ results in an un-penalized spline regression and would produce a very wiggly curve. An
168 infinitely large of λ would lead to a linear estimate of the true function (Wood 2017). The first
169 term in Eq. 3 represents the fitting errors while the second term denotes the penalty against
170 wiggleness of the fitted function. An optimal solution can be sought out by minimizing the
171 objective expression in Eq. 3. For this alternative method, the number and location of knots do
172 not significantly affect the fitted curve once the number of knots or basis function is large
173 enough (Wood 2017).

174 It is worth noting that the users of the GAM can impose their prior information, personal
175 knowledge or judgment regarding the geological profile on the model. This can be
176 accomplished by assigning suitable weights to specific data points which reveal a geological
177 feature and dominate the shape of the geological profile. The prior information can also be fed
178 into the model by setting a proper value of the smoothing parameter. For example, if an area is
179 subject to intensive tectonic activities in history, the geological profile is expected to be wavy.

180 In this case, the smoothing parameter can be set to be relatively small.

181 **2.3 Selection of smoothing parameter**

182 In addition to manually setting a value for the smoothness parameter, its value can also be
183 determined using other ways. A variety of methods have been proposed to determine the value
184 of the smoothing parameter. One simple way is to minimize certain index which denotes the
185 prediction errors, such as the Akaike information criterion (AIC), ordinary cross-validation
186 score, generalized cross-validation (GCV) score (e.g., Wood 2017; Simpson 2018). The second
187 way is to treat the smoothing parameter as a random variable and estimate its value using the
188 maximum likelihood or the restricted likelihood method. In this study, the GCV score is used
189 to determine the value of λ .

190 For a known value of λ , the coefficient parameters β can be estimated using the penalized
191 least square estimation method. The variance of the random error can be estimated as the
192 residual sum of squares divided by the residual degree of freedom. Details of the estimation of
193 these parameters can be found in Wood (2017).

194 **2.4 Mean trend uncertainty, confidence interval and prediction interval**

195 The fitted smooth function can hardly be the actual function of the response variable because
196 of various uncertainties. A reasonable practice is to provide a mean trend as well as a band that
197 denotes the uncertainty. It is worth noting that there are two different kinds of uncertainty,
198 namely the uncertainty in the mean trend and the uncertainty in a prediction (Ruppert, Wand,
199 and Carroll 2003). The former means that if two different sets of data for given explanatory
200 and response variables are used to perform the regression, the fitted curves are expected to be
201 different. This variability in the fitted curve is called the uncertainty in the mean trend. The

202 latter means that if the fitted smooth function is used for predicting a response, the predicted
 203 value will be different from the actual value. This uncertainty is caused by both the error in the
 204 mean trend and the random error which denotes the deviation of data points from the mean
 205 trend (see Fig. 1). The random error can be attributed to measurement errors in the data or other
 206 sources of errors. For example, for the rockhead elevation, the random error can be caused by
 207 the subjective judgments of engineers in determining the weathering degree of the geological
 208 layer. The standard deviation of the mean trend, SD_{μ} indicates the epistemic uncertainty in
 209 the spatial prediction, which can be reduced if more observations are available. SD_{μ} reflects
 210 the geological complexity as well as the data quantity of the investigated area. The value of
 211 SD_{μ} would be quite small if many data exist or the geological profile is very simple such as a
 212 flat curve. The standard deviation of the random errors, SD_{ε} , suggests the magnitude of the
 213 aleatoric uncertainty in the spatial prediction, which cannot be decreased even if sufficient data
 214 are available. SD_{ε} represents the minimum error or maximum accuracy that can be achieved
 215 in spatial predictions.

216 The first uncertainty is usually expressed by confidence interval (CI) while the latter by
 217 the prediction interval (PI). For example, the 95% CI of the mean trend is bounded by the mean
 218 trend minus and plus twice the estimated standard deviation of the mean trend, $2SD_{\mu}$. SD_{μ} is
 219 derived from the standard deviation of the coefficient parameter, which can be estimated using
 220 either the frequentist or Bayesian approach. For the Bayesian method, the posterior distribution
 221 of the standard deviations of the coefficient parameters has an analytical solution when the
 222 prior distribution, $f_{\beta}(\beta)$, is given by (Wood 2017)

$$223 \quad f_{\beta}(\beta) \propto e^{-\frac{1}{2}\beta^T \Sigma S_i \tau_i \beta} \quad (4)$$

224 where the τ_i are parameters controlling the dispersion of the prior distribution; S_i is an
225 element from the penalty matrix, which is a matrix of known coefficients and is derived from
226 $f''(x)$. The prior distribution in Eq. 4 gives equal probability density to all models of equal
227 smoothness, but larger probability densities to smooth models than wiggly models as normally
228 it is believed that smooth models are more likely than wiggly models. More details regarding
229 the estimation of the uncertainty in the coefficient parameter can be found in Wood (2017). The
230 95% PI of a prediction is bounded by the mean trend minus and plus twice the standard
231 deviation of a prediction, given by $2SD_p = 2\sqrt{SD_\mu^2 + SD_\varepsilon^2}$, where SD_ε denotes the standard
232 deviation of the random error. In this study, the fitting of a GAM model is performed using a
233 well-known R package, *mgcv*.

234 3. Spatial prediction using simulated data

235 This section evaluates the performance of the GAM model using cross-validation based on
236 simulated data. Firstly, dense data are simulated based on a smooth trend function and a random
237 error. Secondly, the simulated data are divided into two groups, the training group and the
238 testing group. Thirdly, the training data are used to estimate the unknown parameters of the
239 GAM model and generate the 95% CI and 95% PI. Finally, the coverage percentage of the 95%
240 CI is evaluated by counting the percentage of the input trend function covered by the 95% CI
241 (e.g., take 100 evenly spaced data points from the input trend and check how many of them are
242 covered by the 95% CI) while that of the 95% PI is assessed by computing the percentage of
243 the testing data covered by the 95% PI. These steps are repeated by 500 times and the average
244 values of the coverage percentages are computed.

245 An example taken from Wood (2017) is used to analyze the performance of the GAM,

246 including the CI and PI, the latter of which is not investigated in Wood (2017). The input trend
247 function is given by

$$248 \quad f(x)=x^{11}(10[1-x])^6+10(10x)^3(1-x)^{10} \quad (5)$$

249 The range of the explanatory variable, x is set to be $[0, 1]$ while the range of the input function
250 is scaled to the interval of $[0, 1]$ by dividing $f(x)$ by the maximum value of $f(x)$. The curve of
251 the input smooth function is plotted as the dotted line in Fig. 2(a). 500 samples of x are
252 randomly drawn from the uniform distribution with a range of $[0, 1]$. Simulated data of y were
253 generated by adding a random error to the input mean trend for the 500 samples of x . The
254 random errors are normally distributed with a mean of 0 and standard deviation of 0.2 and are
255 mutually independent at various locations. 50 data points are randomly drawn from the 500
256 data points and set as training data while the remainder as testing data. The training data are
257 used to fit the GAM model and create the 95% CI and 95% PI. The dimension of basis was set
258 to be 20 as a larger basis dimension produces quite similar results. Besides, 500 experiments
259 were carried out as such quantities of experiments are sufficient to yield a converged estimation
260 of the coverage percentages.

261 A typical experiment of the cross-validation is plotted in Fig. 2(a-c). Fig. 2(a) plots the
262 training and testing data, respectively. Fig. 2(b) plots the 95% PI and the 95% CI, which are
263 denoted by the region between the dashed lines and shaded region, respectively. Fig. 2(c) plots
264 the 95% CI of the mean trend as well as 20 samples of the mean trend. The mean trend samples
265 are simulated by first generating samples of the coefficient parameters based on their posterior
266 distribution and then multiplying the model matrix for the testing data by the coefficient vector
267 (i.e., $\mathbf{X}\boldsymbol{\beta}$ in Eq. 3). As shown, the 95% CI can cover most sections of the actual mean trend

268 while the 95% PI can cover most of the testing points, indicating the reasonableness of these
269 intervals. Also, the simulated mean curves generally capture the trend of the response variable.
270 The good performance of the 95% CI and 95% PI can also be seen from Table 1, which
271 summarizes the mean values of the 500 coverage percentages for 500 experiments. For
272 comparison, the performance for the prediction interval ignoring the uncertainty in the mean
273 trend is also studied. This prediction interval refers to the interval derived purely from the
274 random error, namely the interval bounded by the mean trend minus and plus twice the standard
275 deviation of the random errors. As shown in Table 1, both the 95% CI and 95% PI have a
276 reasonable coverage percentage, which is close to the confidence level. The 95% PI ignoring
277 the uncertainty in the mean has a coverage percentage slightly smaller than the confidence level.
278 Furthermore, the coverage percentages of the 95% CI and 95% PI for 30 training points were
279 also evaluated, which is summarized in the last row of Table 1. As shown, the coverage
280 percentage of the 95% CI and 95% PI are still close to the confidence level, 95%. However,
281 the 95% PI ignoring the uncertainty in the mean trend just has an average coverage percentage
282 of 86%, which is a little far from the theoretical value, 95%. The reason is that when the data
283 are limited, the uncertainty in the mean trend is relatively large. Fig. 2(d) plots one typical
284 experiment of cross-validation using 30 training data. As shown, the 95% CI in Fig. 2(d) is
285 considerably wider than that in Fig. 2(b). The reason is that when a relatively small number of
286 data points are used, the coefficient parameters β_i have large uncertainty. In other words,
287 relatively large values of the standard deviations of β_i are estimated in this case. The large
288 uncertainty in β_i is propagated to the mean trend, inducing a wider 95% CI of the mean trend.
289 These phenomena highlight the importance of considering the uncertainty in the mean trend

290 when the data are limited.

291 One issue that may arouse argument is that herein the residual of the simulated data is
292 assumed to be independently distributed rather than autocorrelated. Spatial autocorrelation is
293 well known to be a property of geotechnical or geological properties. To reflect this nature, it
294 may be more reasonable to model the residual of the response variable as autocorrelated
295 random variables. However, we argue that spatial autocorrelation can to some extent be viewed
296 as an artifact generated by the modelers. This is similar to what was stated by Baecher and
297 Christian (2005), i.e., the division of the spatial variation into a mean trend and a residual
298 around the mean is an artifact. As shown by Baecher and Christian (2005), the variance of the
299 residual and the associated autocorrelation function highly depend on the form of the trend
300 function. Hence, it is possible to obtain independent residuals when a suitable trend function is
301 selected. This statement is supported by Qi et al. (2020), which showed that the residual of the
302 rockhead around a mean trend described by a spline function is independent of each other.

303 **4. Spatial prediction using actual borehole data**

304 This section studied the performance of the GAM in dealing with real borehole data. Some
305 borehole data revealing the rockhead of two formations, Bukit Timah Formation and Jurong
306 Formation in Singapore are analyzed. The Bukit Timah Formation data were extracted from
307 100 boreholes while the Jurong Formation data were from 60 boreholes. The borehole data are
308 extracted from site investigation reports for the construction of two metro lines in Singapore.
309 The former data are located at the Bukit Timah Road while the latter nearby the Buona Vista
310 metro station. For both sets of data, the data are projected to a line that is approximately parallel
311 to the metro line. These borehole data have been elaborated in Qi et al. (2020) and are not

312 repeated herein. Similar to the last section, cross-validation is used to evaluate the accuracy of
313 the prediction and the reasonableness of the CI and PI. Firstly, spatial prediction is performed
314 using all the data. Secondly, the data are divided into two groups and cross-validation is
315 performed.

316 Note that there are two major differences between this study and the analyses performed
317 by Qi et al. (2020). Firstly, the two studies use different methods to determine the smoothness
318 of the fitted curve, i.e., Qi et al. (2020) using the knot number and location selection while this
319 study using the wiggleness penalty. Secondly, the uncertainty in the mean trend of the rockhead
320 elevation is ignored by Qi et al. (2020) but considered herein.

321 **4.1 Bukit Timah Formation**

322 **4.1.1 Analysis using all the data**

323 The rockhead elevation of the Bukit Timah Formation is first predicted using all the borehole
324 data, i.e., rockhead elevation from 100 boreholes. When performing the model fitting, one input
325 parameter is the dimension of the basis used to represent the smooth term, i.e., q in Eq. 2. The
326 dimension of the basis should be large enough to approximate the true, but unknown function
327 of the response parameter (Wood 2017). As a rule of thumb, the dimension of the basis is
328 considered to be sufficient when the smoothness selection criterion converges as the basis
329 dimension increases. Table 2(a) summarizes the GCV scores and estimated standard deviations
330 of the random error for various basis dimensions. As shown, when the basis dimension reaches
331 60 ~ 80, both the GCV score and the standard deviation of the random error converge. Hence,
332 the basis dimension is set to 80 and the associated fitted GAM model is plotted in Fig. 3(a). In
333 Fig. 3(a), the 95% CI of the mean trend is represented by the shaded area while the 95% PI of

334 the prediction is bounded by the two dashed lines. As shown, the fitted mean trend is generally
335 consistent with that reported in Qi et al. (2020), indicating the effectiveness of the used method.
336 Besides, the 95% CI of the mean trend has different widths at various locations. For instance,
337 the area within the distance range of [700 m, 2000 m] has a narrower confidence interval than
338 its left-hand and right-hand sides. The reason is that the geological conditions on the left-hand
339 and the right-hand sections are more complex than the middle section. This phenomenon shows
340 that the GAM model can provide an uncertainty that well reflects the geological complexity.

341 Based on the SD_{μ} for various points, it is evaluated that the average value of SD_{μ} is 3.9 m.
342 Hence, the average standard deviation of predictions is $SD_p = \sqrt{SD_{\mu}^2 + SD_{\epsilon}^2} = \sqrt{6.0^2 + 3.9^2} =$
343 7.2 m.

344 Fig. 3(b) plots the variation of the SD_{μ}/SD_p with distance. Since SD_{μ} reflects the
345 magnitude of the geological uncertainty and the data quantity, Fig. 3(b) provides quantitative
346 information on the area with large uncertainty in the trend, which indicates relatively large
347 construction risk and requires additional site investigation. Also, it is worth noting that SD_{μ}
348 denotes the epistemic uncertainty in the spatial prediction that can be reduced to 0 when
349 sufficient data are available while SD_p denotes the sum of the aleatoric and epistemic
350 uncertainties in the spatial prediction. The ratio of the two indexes indicates the potential for
351 improvement in the spatial prediction accuracy if additional data are available. For example, a
352 relatively large value of SD_{μ}/SD_p suggests that the prediction accuracy can be further
353 improved by using additional data. By contrast, if the value of the ratio approaches 0, there is
354 no need to carry out additional site investigations. As shown in Fig. 3(b), the area within the
355 distance range of [700 m, 2000 m] has a relatively small value of SD_{μ}/SD_p , indicating that

356 this area has less complex geological conditions than the remaining areas. Also, some areas
357 have a SD_{μ}/SD_p value larger than 0.6, suggesting that there is a high potential to improve the
358 prediction accuracy.

359 Furthermore, Fig. 3(c) plots the autocorrelation functions of the residuals of the rockhead
360 elevation, namely the measured rockhead elevation minus the mean trend. The horizontal
361 coordinate, lag, in Fig. 3(c) denotes the difference in serial numbers for a pair of data points.
362 For example, the lag for any pair of neighboring data points is 1 while the lag of the data pairs,
363 (1st data point, 3rd data point), (2nd data point, 4th data point), ..., is 2. In other words, the data
364 pairs with the same difference value of the serial number are used to evaluate an autocorrelation
365 coefficient. This is a rough but simple way to check the autocorrelation of the residual. As
366 shown, the autocorrelation coefficient for a lag of 1 is already negative, indicating the residual
367 of the rockhead elevation is independent. This means that there is no need to use a more
368 complex model, such as a mixed model with correlated residuals, to analyze the rockhead data.

369 Fig 3(c) plots 20 simulated samples of the mean trend of rockhead. Each sample is
370 generated by first simulating samples of the coefficient parameters based on their posterior
371 distribution and then multiplying the model matrix by the simulated coefficient vector, namely
372 $\mathbf{X}\boldsymbol{\beta}$ in Eq. 3. As shown, the generated mean trend samples generally reveal the spatial trend of
373 the rockhead elevation. Besides, a larger variability in the mean trend can be observed in the
374 left-hand and right-hand sides than that in the middle section. These simulated trend curves can
375 be readily used in numerical analyses of geological structures, such as stability analysis of
376 slopes or tunneling. This is one major benefit of using the GAM model. Note that the method
377 in Qi et al. (2020) can also produce samples of the mean trend, but these samples are generated

378 mainly by the bootstrap method, namely performing regression using randomly drawn subsets
379 of borehole data. These mean trend samples are not as accurate as those created by the GAM
380 method as the former uses fewer data.

381 **4.1.2 Cross-validation**

382 This section evaluates the performance of the GAM model using cross-validation. Similar to
383 Qi et al. (2020), 50 data points with even serial numbers are set as training data while the
384 remaining 50 points with odd serial numbers as testing data. As mentioned in section 2, one
385 feature of the GAM is that the users can impose a wiggleness constraint on the fitted curve or
386 assign different weights to the data points. Hence, three prediction schemes are considered
387 herein. Scheme 1 imposes no prior information on the smoothing parameter and weight.
388 Scheme 2 sets the smoothing parameter to be 1.5, which corresponds to a relatively wiggly
389 rockhead profile. Scheme 3 assigns larger weights to several data points at the geologically
390 complex area, namely weights for the 6th to 9th, 43rd, 45th, 47th training points = 2, and weights
391 for the remaining training points = 1. Assigning a weight value of m to a certain data point is
392 equivalent to placing m identical observations at the same location (Wood 2020). In all three
393 schemes, the dimension of the basis is set to be 40, which is obtained using the same procedure
394 as that in section 4.1.1.

395 The prediction results for various prediction schemes are plotted in Fig. 4. As shown in
396 Fig. 4(a), the predicted mean trend of the rockhead elevation cannot capture the wavy rockhead
397 profile at the two ends of the section when no prior information is imposed on the model. The
398 reason is that the data at the geologically complex area are so limited and the abnormalities in
399 these data are treated as a random error rather than counted into the mean trend. However, the

400 wavy trends can be well captured by schemes 2 and 3. The reason is that the constraint of the
401 smoothing parameter = 1.5 in scheme 2 makes the mean trend wigglier. Moreover, the larger
402 weights of the data points at the geologically complex area amount to manually adding some
403 data to the critical locations which control the shape of the rockhead profile. This phenomenon
404 shows the capability of the GAM to incorporate prior information of the engineers, such as
405 personal judgment or knowledge of the geological information in the investigated area. On the
406 other hand, even if the users do not have any prior information, they can perform some
407 sensitivity analyses and acquire multiple solutions by altering the values of smoothing
408 parameters or weights. These solutions can be subsequently submitted to experienced engineers
409 and an optimal solution can be decided based on their judgment. This feature is quite useful as
410 the solution is a joint product of the GAM method and engineer judgment. Also, by providing
411 the clients multiple scenarios of the possible rockhead profile, the risks in the construction can
412 be better appreciated and managed.

413 The prediction accuracies for the three considered schemes as well as the results in Qi et
414 al. (2020) are summarized in Table 2(b), including the mean and standard deviation (SD) of the
415 prediction errors, the estimated SD of the random error, SD_{ε} , the mean width of the 95% CI
416 and 95% PI. The mean and SD of the prediction error are evaluated from the 50 errors for the
417 50 testing points. The width of the 95% CI is four times the SD of the mean trend, $4SD_{\mu}$ while
418 the width of 95% PI is four times the SD of a prediction, i.e., $4SD_p = 4\sqrt{SD_{\mu}^2 + SD_{\varepsilon}^2}$, as
419 illustrated in section 2.4. Since Qi et al. (2020) did not consider the uncertainty in the mean
420 trend, the width of the 95% PI is set to be $4SD_{\varepsilon}$ in the last row. As shown in Table 2(b), all the
421 three considered schemes have slightly smaller prediction errors than that in Qi et al. (2020).

422 The two schemes implementing prior information have higher accuracy than the one without
423 prior information. This result shows the usefulness of imposing human judgment on the GAM
424 model. Such incorporation of human judgment is lacking in the method in Qi et al. (2020). It
425 also well demonstrates the idea that the data-driven method should incorporate engineering
426 judgment rather than replace engineering judgment, as discussed by Phoon, Ching, and Shuku
427 (2021). Besides, the 95% PIs produced by the GAM model are wider than that reported in Qi
428 et al. (2020) because the uncertainty in the mean trend is considered by the GAM. The former
429 is more rational than the latter, which can be shown by the coverage percentage of the 95% PI.
430 In Qi et al. (2020), only 42 out of 50 (i.e., 84%) testing data were covered by the 95% PI. But
431 for the GAM, 46, 46 and 45 out of the 50 (i.e., 92%, 92%, 90%) data points are covered by the
432 95% PIs for the three schemes, respectively. This observation justifies the consideration of the
433 uncertainty in the mean trend. The result indicates that the method of Qi et al. (2020)
434 underestimates the prediction uncertainty and induces unsafe designs of geotechnical structures
435 or an underestimation of the underground construction risk. By contrast, the GAM method in
436 this study can reasonably quantify the prediction uncertainty, which leads to a reasonable
437 design of geotechnical structures or underground construction scheme.

438 **4.2 Jurong Formation**

439 To further illustrate the performance of the GAM, this section performs the spatial prediction
440 using the rockhead data of the Jurong Formation. For this case, only 60 data points are
441 distributed along a line with a length of around 3400 m. Such limited data make spatial
442 prediction more challenging. Fig. 5(a) plots the fitted mean trend for the scheme imposing no
443 prior information (referred to as scheme 1) while Fig. 5(b) plots the result for the scheme with

444 the smoothing parameter set to be 5 (referred to as scheme 2). Both schemes use all the data
445 points and a basis dimension of 50, which is selected based on the same method as that in
446 section 4.1.1. Details for the selection of the basis dimension are summarized in Table 3(a). As
447 shown in Fig. 5(a, b), scheme 2 captures more local fluctuation in the mean trend of the
448 rockhead than scheme 1.

449 Cross-validations of the spatial prediction of Jurong Formation rockhead are performed
450 using 40 training points and 20 testing points. The cross-validation also adopts two schemes,
451 scheme 1 imposing no prior information and scheme 2 setting the smoothing parameter to be
452 5. The two schemes use the same training data, testing data and basis dimension, i.e., 30. The
453 predicted rockhead profile and associated 95% CI, 95% PI are plotted in Fig. 5(c, d). For
454 comparison, the mean trend obtained from all the data points and scheme 2 is plotted in both
455 Fig. 5(c) and (d). The associated prediction errors are summarized in Table 3(b). Similar
456 phenomena as those in section 4.1 can be observed in Fig. 5 and Table 3(b). These include (i)
457 imposing some prior information involving the wiggleness of the rockhead profile produces a
458 more accurate estimation of the mean trend when the data quantity is limited, (ii) the GAM
459 model can quantify the uncertainty in the mean trend of rockhead, making the 95% PI wider
460 than that reported in Qi et al. (2020).

461 **5. Conclusions**

462 This study uses the generalized additive model (GAM) for the spatial prediction of the
463 interfaces of geological formations. The performance of the GAM is evaluated using both
464 simulated data and actual borehole data for two geological formations in Singapore. The
465 prediction accuracy, the rationality of the 95% confidence intervals for the mean trend and the

466 95% prediction interval of the prediction are assessed. The benefits of the GAM are
467 summarized as follows.

468 (1) The GAM can produce a reasonable 95% confidence interval and 95% prediction interval
469 as the analyses using the simulated data show that on average the 95% confidence interval
470 covers 94% or 91% of the actual mean trend while the 95% prediction interval covers 94% or
471 92% of testing data. Ignoring the uncertainty in the mean produces a 95% prediction interval
472 with an unreasonably low coverage percentage. Furthermore, samples of the rockhead profile
473 such as the grey curves in Fig. 3(d), can be generated by the GAM. These samples can be
474 viewed to be possible rockhead profiles and be used in future numerical analyses of
475 geotechnical structures.

476 (2) Both the epistemic uncertainty and aleatoric uncertainty in the spatial predictions can be
477 quantified by the GAM method. The former refers to the uncertainty in the mean trend and can
478 be reduced if additional data are available. This uncertainty is affected by the geological
479 complexity and the data quantity. The latter refers to random errors caused by various factors
480 such as engineers' subjective judgments of weathering degrees of the geological layers. The
481 aleatoric uncertainty cannot be reduced and represents the minimum error that can be achieved
482 in spatial predictions. The relative magnitudes of the two uncertainties can be quantified by the
483 ratio of the standard deviation of the mean trend and the standard deviation of predictions.
484 From the ratio values, it is easy to determine the areas with complex geological conditions,
485 which need additional site investigations. It is also easy to check the potential for improvement
486 in the prediction accuracy when additional data are available because in theory the minimum
487 value of the ratio is 0.

488 (3) The users can apply expert judgment or knowledge on the geological profile to the model
489 by setting suitable values of the smoothing parameter or assigning suitable weights to the data
490 points at critical locations. The analyses using the actual data in the two cases show that the
491 use of expert knowledge improves the prediction accuracy and makes the resulting geological
492 profile more consistent with that obtained from more data.

493 (4) Due to the large variability of the rockhead locations, the prediction error of the rockhead
494 elevation is still relatively large. In the future, it is of interest to use additional data such as
495 geophysical data to further increase the prediction accuracy and reduce the prediction
496 uncertainties.

497 **Acknowledgment**

498 The second author would like to thank the financial support from the National Natural Science
499 Foundation of China (No. 52109144) and the Open Innovation Fund of Changjiang Institute of
500 Survey, Planning, Design and Research (No. CX2020K07).

501 **Reference**

502 Aswar, D. S. , and P. B. Ullagaddi. 2017. "An overview of 3-D geological modelling part II.
503 Summary of major 3-d geological modelling methodologies." *International Journal of*
504 *Latest Engineering and Management Research* 2 (11):15-27.

505 Baecher, G. B., and J. T. Christian. 2005. *Reliability and statistics in geotechnical engineering:*
506 *John Wiley & Sons.*

507 Baecher, G. B., and J. T. Christian. 2008. "Spatial variability and geotechnical reliability." In
508 *Reliability-based design in geotechnical engineering*, 88-145. CRC Press.

509 Bordoni, M., Y. Galanti, C. Bartelletti, M. G. Persichillo, M. Barsanti, R. Giannecchini, G. D.

510 Avanzi, et al. 2020. "The influence of the inventory on the determination of the rainfall-
511 induced shallow landslides susceptibility using generalized additive models." CATENA
512 193:104630. doi: <https://doi.org/10.1016/j.catena.2020.104630>.

513 Burke, H.F., J.R. Ford, L. Hughes, S. Thorpe, and J.R. Lee. 2017. "A 3D geological model of
514 the superficial deposits in the Selby area." In Groundwater programme commissioned
515 report CR/17/112. Nottingham British Geological Survey.

516 de Brogniez, D., C. Ballabio, A. Stevens, R. J. A. Jones, L. Montanarella, and B. van Wesemael.
517 2015. "A map of the topsoil organic carbon content of Europe generated by a generalized
518 additive model." European Journal of Soil Science 66 (1):121-34. doi:
519 <https://doi.org/10.1111/ejss.12193>.

520 De Boor, C. 2001. A practical guide to splines. Revised Edition. New York: Springer-Verlag

521 Goetz, J. N., A. Brenning, H. Petschko, and P. Leopold. 2015. "Evaluating machine learning
522 and statistical prediction techniques for landslide susceptibility modeling." Computers
523 & Geosciences 81:1-11. doi: <https://doi.org/10.1016/j.cageo.2015.04.007>.

524 Gong, X., A. Kaulfus, U. Nair, and D. A. Jaffe. 2017. "Quantifying O3 impacts in urban areas
525 due to wildfires using a generalized additive model." Environmental Science &
526 Technology 51 (22):13216-23. doi: 10.1021/acs.est.7b03130.

527 Gong, W. P., C. Zhao, C. H. Juang, H.M. Tang, H. Wang, and X.L. Hu. 2020. "Stratigraphic
528 uncertainty modelling with random field approach." Computers and Geotechnics
529 125:103681.

530 Hastie, T. J., and R. J. Tibshirani. 1986. "Generalized Additive Models." Statistical Science
531 1 (3):297-310, 14.

532 Hastie, T. J., and R. J. Tibshirani. 1990. Generalized additive models. CRC press.

533 Huang, Y., J. L. Beck, S. Wu, and H. Li. 2014. "Robust Bayesian compressive sensing for
534 signals in structural health monitoring." *Computer-Aided Civil and Infrastructure*
535 *Engineering* 29 (3):160-79.

536 Khoda Bakhshi, A., and M. M. Ahmed. 2021. "Real-time crash prediction for a long low-traffic
537 volume corridor using corrected-impurity importance and semi-parametric generalized
538 additive model." *Journal of Transportation Safety & Security*:1-35. doi:
539 10.1080/19439962.2021.1898069.

540 Lark, R.M., S.J. Mathers, S. Thorpe, S.L.B. Arkley, D.J. Morgan, and D.J.D. Lawrence. 2013.
541 "A statistical assessment of the uncertainty in a 3-D geological framework model."
542 *Proceedings of the Geologists' Association* 124 (6):946-58.

543 Li, J.H., Y.M. Cai, X.Y. Li, and L.M. Zhang. 2019. "Simulating realistic geological stratigraphy
544 using direction-dependent coupled Markov chain model." *Computers and Geotechnics*
545 115:103147.

546 Liu, L. L., Y. M. Cheng, Q. J. Pan, and D. Dias. 2020. "Incorporating stratigraphic boundary
547 uncertainty into reliability analysis of slopes in spatially variable soils using one-
548 dimensional conditional Markov chain model." *Computers and Geotechnics* 118:103321.

549 Ma, Y., B. Ma, H. Jiao, Y. Zhang, J. Xin, and Z. Yu. 2020. "An analysis of the effects of weather
550 and air pollution on tropospheric ozone using a generalized additive model in Western
551 China: Lanzhou, Gansu." *Atmospheric Environment* 224:117342. doi:
552 <https://doi.org/10.1016/j.atmosenv.2020.117342>.

553 Mariethoz, G., and J. Caers. 2014. Multiple-point geostatistics: stochastic modeling with

554 training images: John Wiley & Sons.

555 Phoon, K.K., J.Y. Ching, and T. Shuku. 2021. "Challenges in data-driven site characterization."
556 Georisk: Assessment and Management of Risk for Engineered Systems and
557 Geohazards:1-13. doi: 10.1080/17499518.2021.1896005.

558 Qi, X.H., D.Q. Li, K.K. Phoon, Z.J. Cao, and X.S. Tang. 2016. "Simulation of geologic
559 uncertainty using coupled Markov chain." *Engineering Geology* 207:129-40. doi:
560 10.1016/j.enggeo.2016.04.017.

561 Qi, X.H., and H.X. Liu. 2019. "An improved global zonation method for geotechnical
562 parameters." *Engineering Geology* 248:185-96. doi: 10.1016/j.enggeo.2018.11.013.

563 Qi, X.H., X.H. Pan, K. Chiam, Y. S. Lim, and S. G. Lau. 2020. "Comparative spatial predictions
564 of the locations of soil-rock interface." *Engineering Geology* 272:105651. doi:
565 <https://doi.org/10.1016/j.enggeo.2020.105651>.

566 Qi, X., H. Wang, J. Chu, and K. Chiam. 2021a. "Effect of autocorrelation function model on
567 spatial prediction of geological interfaces." *Canadian Geotechnical Journal*, accepted,
568 10.1139/cgj-2020-0644.

569 Qi, X.H., H. Wang, X.H. Pan, J. Chu, and K. Chiam. 2021b. "Prediction of interfaces of
570 geological formations using the multivariate adaptive regression spline method."
571 *Underground Space* 6(3): 252-66.

572 Ruppert, D., M. P. Wand, and R. J. Carroll. 2003. *Semiparametric regression*: Cambridge
573 university press.

574 Simpson, G. L. 2018. "Modelling palaeoecological time series using generalised additive
575 models." *Frontiers in Ecology and Evolution* 6:149.

576 Smirnoff, A., E. Boisvert, and S. J. Paradis. 2008. "Support vector machine for 3D modelling
577 from sparse geological information of various origins." *Computers & Geosciences* 34
578 (2):127-43.

579 Wang, H., J. F. Wellmann, Z. Li, X.R. Wang, and R. Y. Liang. 2017. "A segmentation approach
580 for stochastic geological modeling using hidden Markov random fields." *Mathematical*
581 *Geosciences* 49 (2):145-77. doi: 10.1007/s11004-016-9663-9.

582 Wang, X.R., H. Wang, R. Y. Liang, H.H. Zhu, and H.G. Di. 2018. "A hidden Markov random
583 field model based approach for probabilistic site characterization using multiple cone
584 penetration test data." *Structural Safety* 70:128-38. doi:
585 <https://doi.org/10.1016/j.strusafe.2017.10.011>.

586 Wang, Y, and T.Y. Zhao. 2016. "Interpretation of soil property profile from limited
587 measurement data: a compressive sampling perspective." *Canadian Geotechnical*
588 *Journal* 53 (9):1547-59. doi: 10.1139/cgj-2015-0545.

589 Wang, Y. , and T. Zhao. 2017. "Statistical interpretation of soil property profiles from sparse
590 data using Bayesian compressive sampling." *Géotechnique* 67 (6):523-36. doi:
591 10.1680/jgeot.16.P.143.

592 Wang, Y., O. V. Akeju, and T.Y. Zhao. 2017. "Interpolation of spatially varying but sparsely
593 measured geo-data: A comparative study." *Engineering Geology* 231:200-17. doi:
594 <https://doi.org/10.1016/j.enggeo.2017.10.019>.

595 Wang, Y., T.Y. Zhao, and K.K. Phoon. 2018. "Direct simulation of random field samples from
596 sparsely measured geotechnical data with consideration of uncertainty in interpretation."
597 *Canadian Geotechnical Journal* 55 (6):862-80.

598 Wang, Yu, T.Y. Zhao, Y. Hu, and K.K. Phoon. 2019. "Simulation of random fields with trend
599 from sparse measurements without detrending." *Journal of Engineering Mechanics* 145
600 (2):04018130.

601 Wood, S. N. 2017. *Generalized additive models: an introduction with R*: CRC press.

602 Wood, S. N. 2020. Package 'mgcv', Mixed GAM computation vehicle with automatic
603 smoothness estimation. R package version.

604 Yee, T. W., and N. D. Mitchell. 1991. "Generalized additive models in plant ecology." *Journal*
605 *of Vegetation Science* 2 (5):587-602. doi: <https://doi.org/10.2307/3236170>.

606 Zhao, J., B. B. Broms, Y. Zhou, and V. Choa. 1994. "A study of the weathering of the Bukit
607 Timah Granite Part a: Review, field observations and geophysical survey." *Bulletin of*
608 *the International Association of Engineering Geology* 49 (1):97-106. doi:
609 10.1007/BF02595006.

610 Zhao, T.Y., Y. Hu, and Y. Wang. 2018. "Statistical interpretation of spatially varying 2D geo-
611 data from sparse measurements using Bayesian compressive sampling." *Engineering*
612 *Geology* 246:162-75. doi: <https://doi.org/10.1016/j.enggeo.2018.09.022>.

613 Zhao, C., W.P. Gong, T.Z. Li, C. H. Juang, H.M. Tang, and H. Wang. 2021. "Probabilistic
614 characterization of subsurface stratigraphic configuration with modified random field
615 approach." *Engineering Geology* 288:106138.

616 Zhou, Y.X., and X.P. Wu. 1994. "Use of neural networks in the analysis and interpretation of
617 site investigation data." *Computers and Geotechnics* 16 (2):105-22.

618

Table 1 Mean values of the coverage percentages for the simulated data

Number of training point	Mean of coverage percentage for the 95% CI	Mean of coverage percentage for the 95% PI	Mean of coverage percentage for the 95% PI ignoring the uncertainty in the mean trend
50	0.94	0.94	0.91
30	0.91	0.92	0.86

Table 2 Spatial prediction of the rockhead elevation for the Bukit Timah Formation

(a) Selection of the dimension of basis for all the data

Basis dimension	20	40	60	80
GCV	74.9	60.9	58.8	58.5
Estimated standard deviation of the random error (m)	8.0	6.5	6.1	6.0

(b) Prediction results of the cross-validation

Fitting scheme	Mean of prediction error (m)	^a SD of prediction error (m)	Estimated ^a SD of random error (m)	Mean width of 95% CI (m)	Mean width of 95% PI (m)
^b Scheme 1	1.8	8.5	8.2	14.7	36.0
^c Scheme 2	1.6	7.5	6.8	21.8	35.0
^d Scheme 3	1.5	7.6	7.5	19.6	35.9
Qi et al. (2020)	2.0	9.0	6.1	–	24.4

Note: a: SD = standard deviation;

b: scheme 1 imposes no prior information on the smoothness parameter or weights;

d: scheme 2 sets the smoothness parameter to be 1.5;

d: scheme3 gives more weights to the data points at the geologically complex area, namely weight for the 6th to 9th, 43rd, 45th, 47th training points = 2, and weight for the remaining training data = 1.

Table 3 Spatial prediction of the rockhead elevation for the Jurong Formation

(a) Selection of the dimension of basis using all the data

Basis dimension	20	30	40	50
GCV	83.4	86.4	83.4	83.2
Estimated standard deviation of the random error (m)	8.2	7.0	8.2	8.2

(b) Prediction results of the cross-validation

Fitting scheme	Mean of prediction error (m)	^a SD of prediction error (m)	Estimated SD of random error (m)	Mean width of 95% CI (m)	Mean width of 95% PI (m)
^b Scheme 1	-2.0	10.6	11.4	22.1	50.9
^c Scheme 2	-2.3	10.5	10.5	30.6	52.2
Qi et al. (2020)	-3.1	12.7	8.1	–	32.4

Note: a: SD = standard deviation;

b: scheme 1 imposes no prior information on the smoothness parameter or weights;

c: scheme 2 sets the smoothness parameter to be 5.

Figure captions

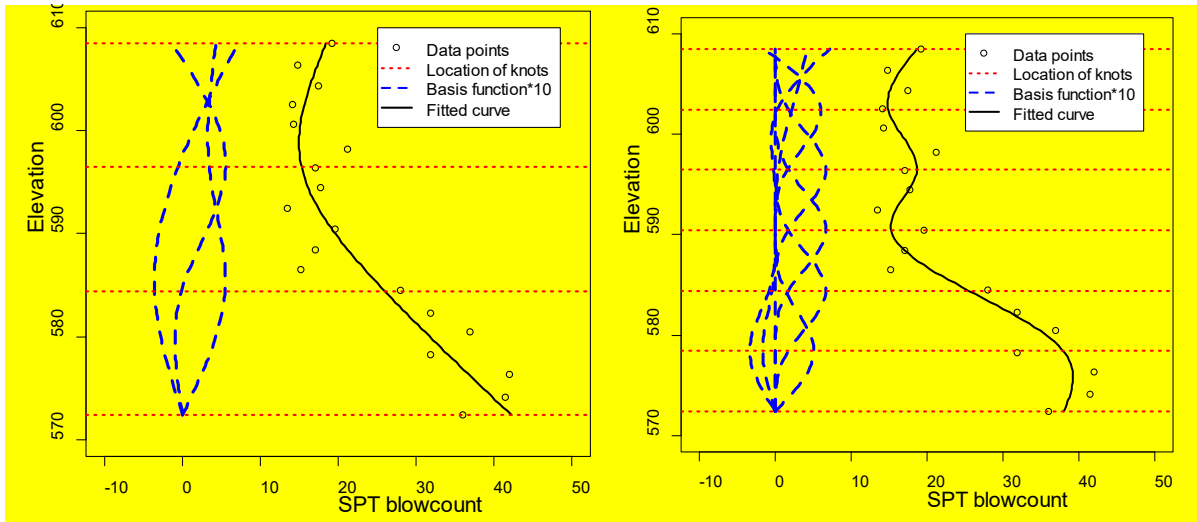
Fig. 1 Example of cubic spline basis and the fitted smooth curve

Fig. 2 Prediction using simulated data

Fig. 3 Spatial prediction of rockhead elevation for the Bukit Timah Formation using all the 100 borehole data

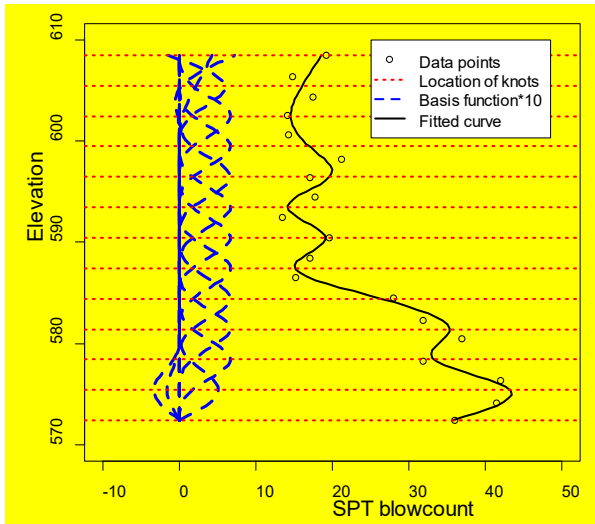
Fig. 4 Cross-validation for the spatial prediction of the rockhead elevation for the Bukit Timah Formation

Fig. 5 Spatial prediction of the rockhead elevation for the Jurong Formation



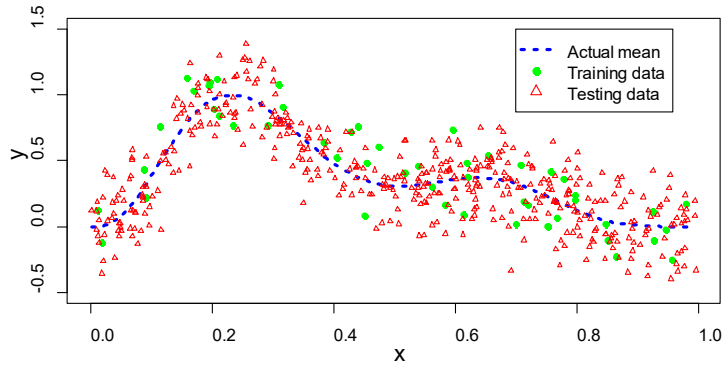
(a) Three cubic basis functions

(b) Six cubic basis functions

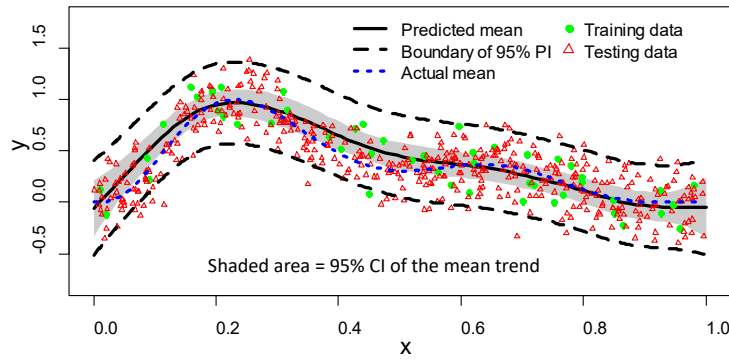


(c) Twelve basis functions

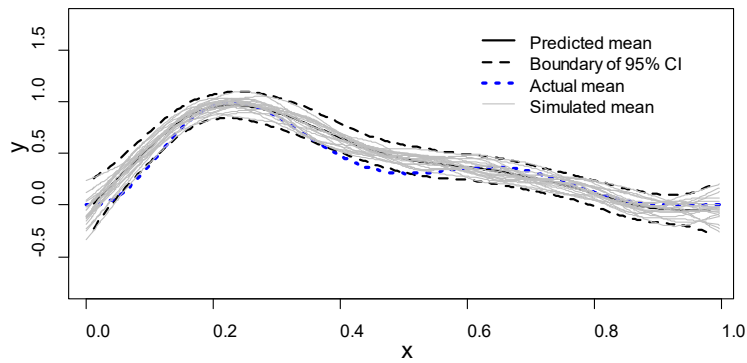
Fig. 1 Example of cubic spline basis and the fitted smooth curve



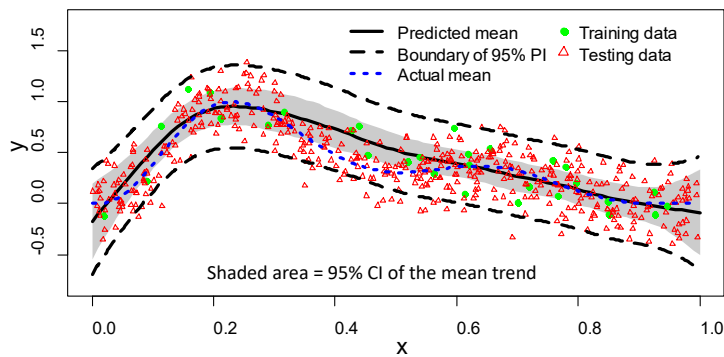
(a) Simulated data (50 training points and 450 testing points)



(b) 95% confidence interval and 95% prediction interval using 50 training points

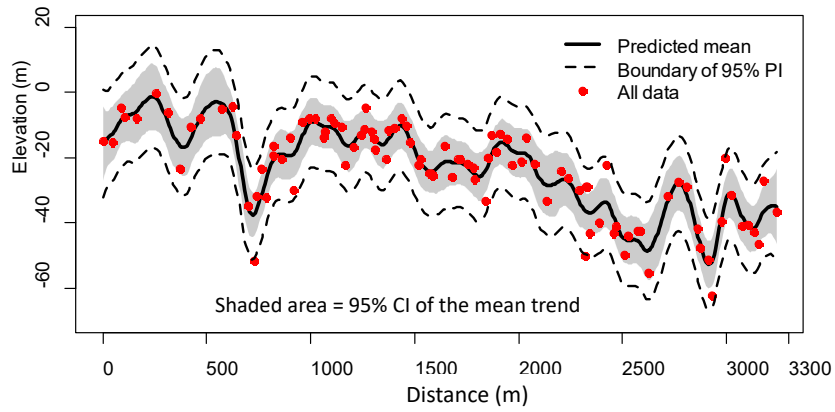


(c) Samples of the mean trend based on 50 training points

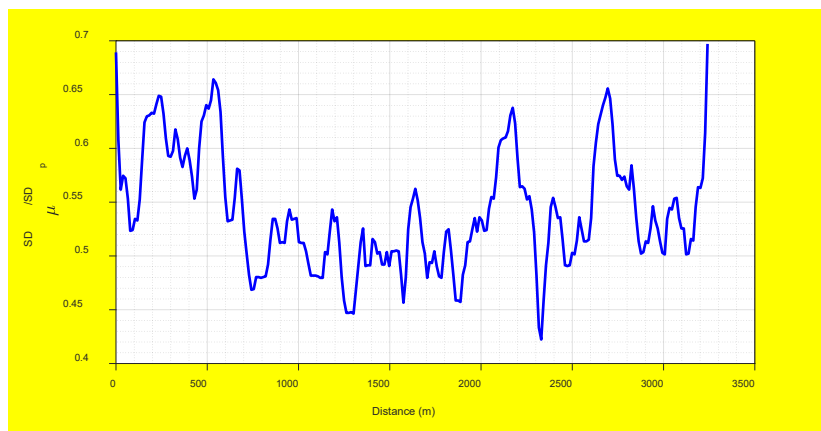


(d) 95% confidence interval and 95% prediction interval using 30 training points

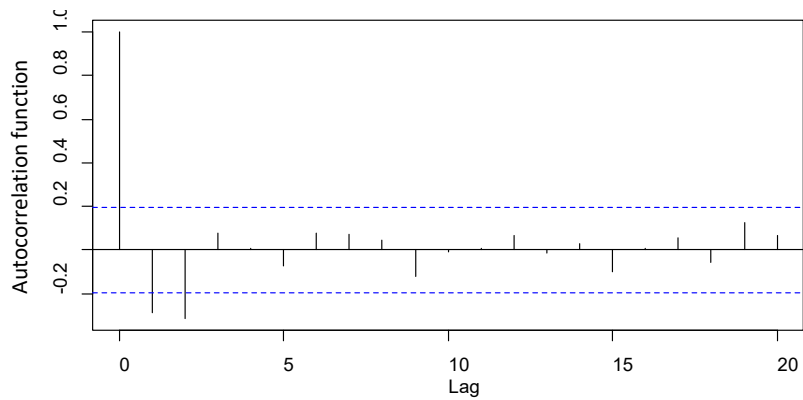
Fig. 2 Prediction using simulated data



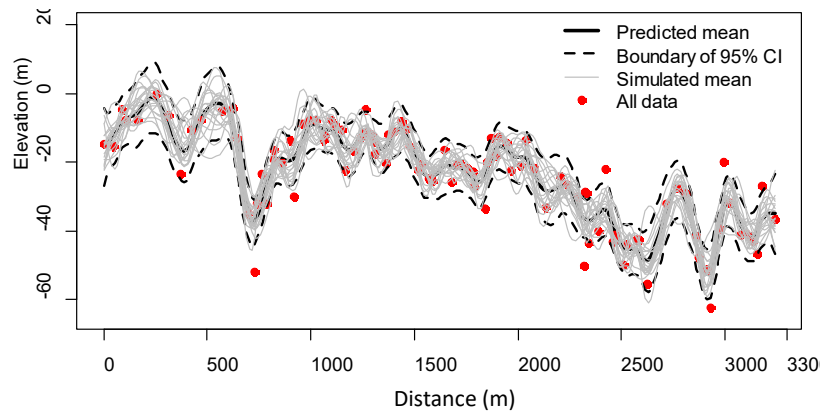
(a) Predicted mean trend, 95% confidence interval, and 95% prediction interval



(b) Ratio of SD_{μ} to SD_p

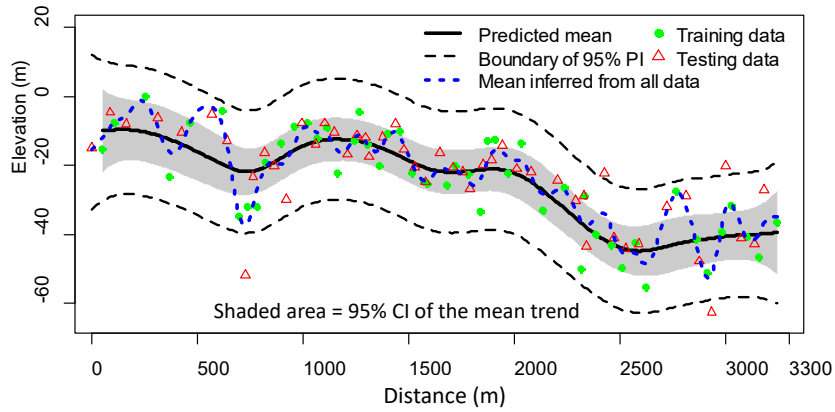


(c) Autocorrelation of the residual of rockhead elevation

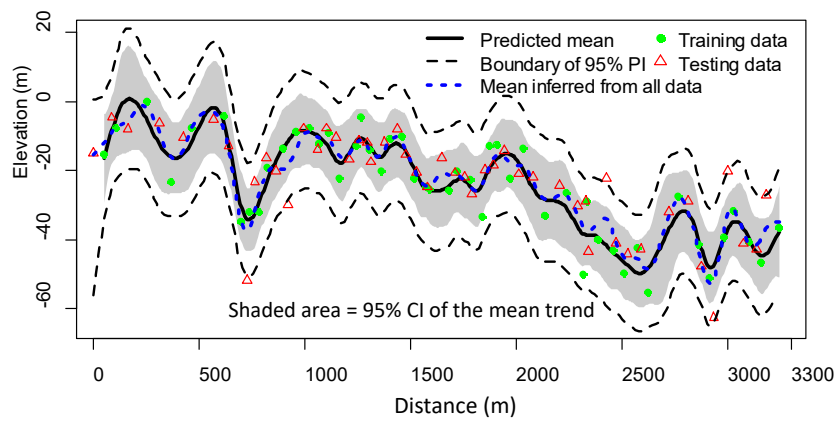


(d) Sample of the mean trend

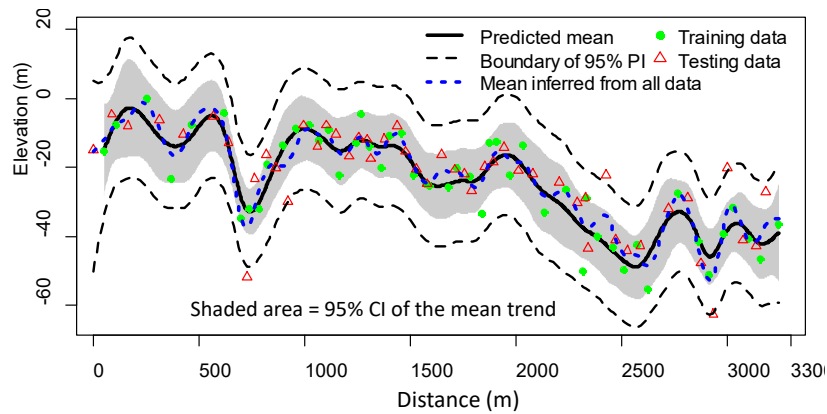
Fig. 3 Spatial prediction of rockhead elevation for the Bukit Timah Formation using all the 100 borehole data



(a) Scheme 1 (no prior information in the smoothness parameter or weights)

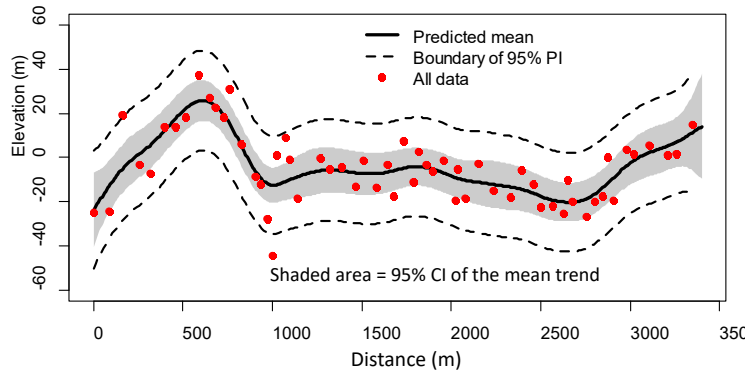


(b) Scheme 2 (the smoothness parameter fixed at 1.5)

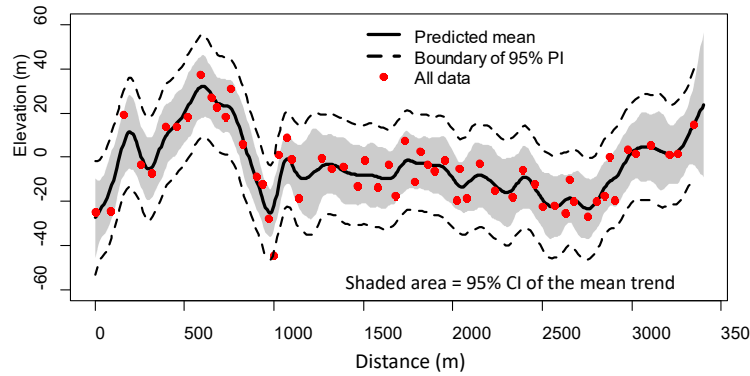


(c) Scheme 3 (Weight of 6th ~ 9th, 43rd, 45th, 47th training points = 2, and weight of remaining training points = 1)

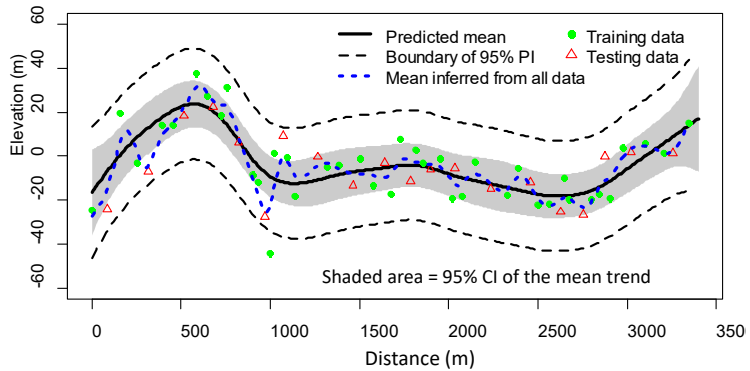
Fig. 4 Cross-validation for the spatial prediction of the rockhead elevation for the Bukit Timah Formation



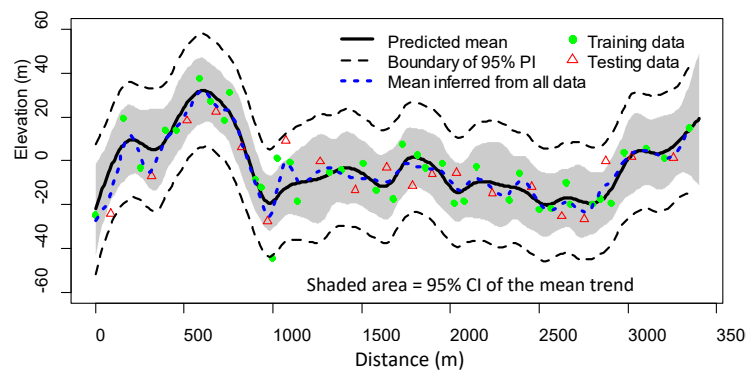
(a) Prediction using all the data (scheme 1: no prior information in the smoothness parameter)



(b) Prediction using all the data (scheme 2: smoothing parameter = 5)



(c) Cross-validation (scheme 1: no prior information in the smoothness parameter)



(d) Cross-validation (scheme 2: smoothing parameter = 5)

Fig. 5 Spatial prediction of the rockhead elevation for the Jurong Formation