

Weili Fang, Yixiao Shao, Peter E.D. Love, Timo Hartmann, Wenli Liu

## **Detecting anomalies and de-noising monitoring data from sensors: A smart data approach**

**Open Access via institutional repository of Technische Universität Berlin**

### **Document type**

Journal article | Accepted version

(i. e. final author-created version that incorporates referee comments and is the version accepted for publication; also known as: Author's Accepted Manuscript (AAM), Final Draft, Postprint)

### **This version is available at**

<https://doi.org/10.14279/depositonce-20005>

### **Citation details**

Fang, W., Shao, Y., Love, P. E. D., Hartmann, T., & Liu, W. (2023). Detecting anomalies and de-noising monitoring data from sensors: A smart data approach. In *Advanced Engineering Informatics* (Vol. 55, p. 101870). Elsevier BV. <https://doi.org/10.1016/j.aei.2022.101870>.

### **Terms of use**

This work is protected by copyright and/or related rights. You are free to use this work in any way permitted by the copyright and related rights legislation that applies to your usage. For other uses, you must obtain permission from the rights-holder(s).

1 Detecting Anomalies and De-Noising Monitoring Data from Sensors:  
2 A Smart Data Approach

3  
4 Weili Fang<sup>1</sup>, Yixiao Shao<sup>2</sup>, Peter E.D. Love<sup>3</sup>, Timo Hartmann<sup>1</sup>, Wenli Liu<sup>2</sup>

5  
6 <sup>1</sup>Department of Civil and Building Systems, Technische Universität Berlin, Gustav-  
7 Meyer-Allee 25, 13156 Berlin, Germany

8  
9 <sup>2</sup>School of Civil and Hydraulic Engineering, Huazhong University of Science and  
10 Technology, Wuhan, 430074, China

11  
12 <sup>3</sup>School of Civil and Mechanical Engineering, Curtin University, GPO Box U1987, Perth,  
13 Western Australia 6845, Australia

14  
15 Corresponding Author: Wenli Liu, Email: liu\_wenli@hust.edu.cn

27    **Detecting Anomalies and De-Noising Monitoring Data from Sensors:**  
28                                    **A Smart Data Approach**

29

30    **Abstract:** When monitoring safety levels in deep pit foundations using sensors, anomalies (e.g.,  
31    highly correlated variables) and noise (e.g., high dimensionality) exist in the extracted time  
32    series data, impacting the ability to assess geotechnical and structural safety risks. Our research  
33    aims to address the following question: *How can we detect anomalies and de-noise monitoring*  
34    *data from sensors in real time to improve its quality and use it to assess geotechnical safety*  
35    *risks?* In addressing this research question, we develop a hybrid smart data approach that  
36    integrates Extended Isolation Forest and Variational Mode Decomposition models to detect  
37    anomalies and de-noise data effectively. We use real-life data obtained from sensors to validate  
38    our smart data approach while constructing a deep pit foundation. Our smart data approach can  
39    detect anomalies with a root mean square error and signal-to-noise ratio of 0.0389 and 24.09,  
40    respectively. To this end, our smart data approach can effectively pre-process data enabling  
41    improved decision-making and the management of safety risks.

42

43    **Keywords:** Anomaly, deep pit foundations, de-noise, detection, smart data, safety risks

44

## 45 **1.0 Introduction**

46 Developments in sensing and data-processing technologies have enabled the effective  
47 monitoring of engineering data during the construction of deep pit foundations enabling  
48 geotechnical safety risks to be examined in greater detail (Zhou *et al.*, 2019a;b; Asadzadeh *et*  
49 *al.*, 2020). The data extracted and transmitted from sensors is imperfect, containing anomalies  
50 and noise, which severely jeopardizes its accuracy and completeness, which can also be  
51 exacerbated by random disturbances (Bao *et al.*, 2019; Nessa *et al.*, 2020; Li *et al.*, 2021; Liu  
52 *et al.*, 2022; Seites-Rundlett *et al.*, 2022).

53

54 Anomalies do not comply with the expected patterns (e.g., data missing) and possess various  
55 characteristics, which are common when using sensors. They generally occur due to faults,  
56 transmission errors, or structural damage (Kromanis and Kripakaran, 2013; Yi *et al.*, 2013).  
57 Consequently, anomalous data may provide false information for decision-making and  
58 determination of safety risks (e. g., assessing geotechnical conditions). Thus, monitoring data  
59 needs to be automated and accurately detect anomalies to ensure it is robust and relevant for  
60 risk assessment (Nguyen and Goulet, 2019).

61

62 However, detecting anomalies is challenging due to the high dimensionality and correlations  
63 between the extracted engineering data (Thudumu *et al.*, 2020). As sensors continuously  
64 monitor data, it is unfeasible to inspect and detect it in real time manually. Thus, several  
65 machine-learning approaches, such as neural network classifiers and decision trees, have been  
66 proposed to detect anomalies in real-time (Zuo and Xiong, 2019; Huang *et al.*, 2020). While  
67 such approaches can detect anomalies quickly, they depend on several labeled  
68 (normal/abnormal) databases for their identification (Ahmed *et al.*, 2016). Furthermore, labeled  
69 samples often contain noise generated by external disturbances. In the case of data monitoring

70 in deep pit foundations, such disturbances are attributable to movements and vibrations  
71 generated by plants and equipment.

72

73 To detect anomalies, supervised and unsupervised approaches have been adopted widely.  
74 Commonly used supervised that can detect anomalies with high levels of performance (i.e.,  
75 low false alarm rate) are the Support Vector Machine (Bhavsar and Waghmare, 2013) and  
76 Random Forest (Hasan *et al.*, 2014). The datasets that employ supervised approaches in  
77 complex engineering environments require high-quality labeling. But datasets are often  
78 incomplete, requiring labeling to be undertaken manually, which is a time-consuming process.  
79 Contrastingly, the unsupervised approach does not require labeling. Their use has been  
80 advocated for detecting anomalies in sensor data (Chen *et al.*, 2017). However, the detection  
81 rate is always low, and false-positive rates are high (Chen *et al.*, 2017). Typical techniques used  
82 to process signals are the Wavelet Transform, Fourier Transform, and Empirical Mode  
83 Decomposition (EMD), though each has limitations (Abbate *et al.*, 1997; Urciuolo and Marta,  
84 2008; Hou and Guo, 2020). For example, the Fourier transform is unsuitable for handling non-  
85 stationary, non-linear signals with frequency over time (Urciuolo and Marta 2008), and the  
86 components of EMD are prone to modal aliasing (Hou and Guo 2020).

87

88 Against this contextual backdrop, our research addresses the following question: *How can we*  
89 *detect anomalies and de-noise monitoring data in real-time to improve its quality and use it to*  
90 *assess geotechnical safety risks?* In addressing this research question, a hybrid smart data  
91 approach that integrates the EIF and Variational Mode Decomposition (VMD) models is  
92 proposed to effectively detect anomalies and de-noise monitoring data to improve its quality to  
93 assess safety risks. The Extended Isolation Forest (EIF), an unsupervised anomaly detection  
94 algorithm, can detect anomalies and performs comparably to supervised algorithms (Carrera *et*

95 *al.*, 2022).

96

97 The basic idea of EIF is similar to the Isolation Forest (IF), which does not rely on building a  
98 profile for data to find non-conforming samples and remedies the shortcomings of the IF, which  
99 arise due to biases in the way the branching of the trees takes place (Hariri *et al.*, 2021). The  
100 EMD is an algorithmic method to detect and decompose a signal into principal “modes” and is  
101 widely used in various time-frequency analysis applications. The EMD method is adaptive and  
102 applicable to non-linear and non-stationary processes (Huang *et al.*, 1998). The variational  
103 mode decomposition (VMD) (i.e., a non-recursive and noise robustness multi-resolution  
104 decomposition method) has better noise robustness performance than EMD in the application  
105 of vibration signal decomposition (Li *et al.*, 2021). Furthermore, when compared with EMD,  
106 problems such as modal aliasing and end-point effects are better avoided (Cai *et al.*, 2022), so  
107 the method is introduced in geotechnical engineering monitoring.

108

109 Our research commences by reviewing existing studies on anomaly data detection and  
110 denoising (Section 2). We then present a novel smart data approach that integrates EIF and  
111 VMD models to address anomalies and de-noise monitoring data to improve the ability to  
112 assess safety risks (Section 3). Our approach is smart as we focus on extracting *only* relevant  
113 engineering data for making decisions about geotechnical safety risks (Matthews *et al.*, 2022).  
114 Next, the feasibility and effectiveness of our proposed approach are presented (Section 4). We  
115 subsequently discuss the implications of our approach and identify its limitations (Section 5)  
116 before submitting our conclusions (Section 6).

117

## 118 **2.0 Monitoring Data from Sensors**

119 The quality of engineering data obtained from sensors plays a pivotal role in monitoring the

120 safety conditions in construction, especially in hazardous areas such as deep pit foundations.  
121 The detection of anomalous behavior from sensor data has received considerable attention in  
122 the literature. However, within the context of construction operations, research has been limited,  
123 though the problem of anomalies and noise reduction remains akin to other applications (e.g.,  
124 Rabatel *et al.*, 2011; Ahmed *et al.*, 2016; Domingues *et al.*, 2018; Hu *et al.*, 2019).

125

### 126 **2.1 Detecting Anomalies and Noise**

127 Many approaches have been designed and developed to detect anomalies and are reported in  
128 the extant literature (Hill and Minsker, 2010; Cha and Wang, 2018). Existing sensor  
129 measurement approaches can be divided into three categories: (1) rule-based; (2) supervised  
130 learning-based; and (3) unsupervised learning-based (Huang *et al.*, 2017; Cha and Wang, 2018;  
131 Huang *et al.*, 2020; Gao *et al.*, 2022). For example, Mu and Yuen (2015) formulated an outlier-  
132 resistant extended Kalman filter to detect outliers caused by measurement errors. Similarly,  
133 Cha and Wang (2018) proposed an unsupervised anomaly-identification approach by  
134 modifying the original density-based fast clustering method. In this instance, Cha and Wang  
135 (2018) improved the ability to detect the location of structural damage by using a ‘Gaussian  
136 kernel function of radius’ to calculate the local density of data points. By the same token, under  
137 the assumption that measurement noise is Gaussian distributed, Huang *et al.* (2017) presented  
138 an anomaly-identification method in the noisy subspace of Principle Component Analysis.  
139 Examples of studies detecting anomalies from sensor measurement are shown in Table 1.

140

141 Notably, several challenges arise when using the above approaches to detect anomalies. For  
142 example, the rule-based approaches fail to recognize malicious events where no rules have  
143 been specified (Thottan and Ji 2003). Indeed, rule-based systems are restricted to only  
144 identifying events where rules exist. In the case of supervised learning, training data needs to

145 be labeled, and algorithms cannot be used if this is not the case (Chandola *et al.*, 2009; Ahmed  
 146 *et al.*, 2016). However, unsupervised learning approaches can train unlabeled data (Otoum *et*  
 147 *al.*, 2018). Despite unsupervised learning addressing this problem, training has challenges  
 148 (Ahmed *et al.*, 2016). Most unsupervised machine learning approaches to detect anomalies are  
 149 evaluated on relatively small datasets in other domains (Inoue *et al.*, 2017; Otoum *et al.*, 2018).  
 150 Moreover, with data being unlabeled, normal and abnormal signals can become mixed,  
 151 rendering it difficult to demarcate the boundary between them. Therefore, normal data may  
 152 contain anomalies in some specific scenarios.

153

154 Table 1. Examples of detecting anomalies in sensor measurement studies

155

<b>Research approach</b>	<b>Description</b>	<b>Author (Year)</b>
Pattern recognition neural network	Detection of multi-type data anomaly for structural health monitoring (SHM)	Gao <i>et al.</i> (2022)
Dynamic independent component analysis	Identification of two types of data anomalies in the SHM system of a cable-stayed bridge and then infer the structural damage	Huang <i>et al.</i> (2020)
Data visualization and deep learning network	Detect seven types of data anomalies in the SHM system of a long-span bridge	Bao <i>et al.</i> (2019)
Artificial neural network	A distributed similarity test and an artificial neural network were proposed to identify drift, spikes, and bias anomalies in wireless sensor networks	Fu <i>et al.</i> (2019)
Neural network	Estimate the state and detect the anomaly in a thermal power plant via a health monitoring system with multilayer perception	Banjanovic- Mehmedovic <i>et al.</i> (2017)
Autoregressive modelling and Kalman estimator	Detection of three types of data anomalies	Chang <i>et al.</i> (2017)

156

157 To this end, we aim to address the above challenges to develop an effective anomaly detection



158 approach to improve the quality of engineering data extracted from sensors within deep pit  
 159 foundations. Thus, to detect anomalies, we propose using the EIF described below.

160

### 161 2.1.1 Extended Isolation Forest

162 The EIF model was first proposed by Hariri *et al.* (2021). It extends the model-free anomaly  
 163 detection algorithm, Isolation Forest (*iForest*). The EIF extracts features from each monitoring  
 164 dataset (e.g., the shaft force of steel shotcrete and building settlement) and builds a baseline  
 165 model by creating an extended isolation forest tree collection. When new monitoring data is  
 166 collected, it is mapped into each of these IFtrees, and an anomaly score is calculated. It will be  
 167 defined as normal if its anomaly score is under a designated threshold value (Table 1).  
 168 Otherwise, the monitoring data will be specified as abnormal.

169

170 *iForest* samples  $n$  instances as a subset from the training dataset  $\{X_1, X_2, \dots, X_N\}$ , where  $X_i =$   
 171  $[X_{i,1}, X_{i,2}, \dots, X_{i,D}]^T$  denotes one  $D$ -dimensional data instance and then generates a binary tree  
 172 from the root node. It randomly chooses one dimension from all  $D$  dimensions and randomly  
 173 samples a split value from the uniform distribution  $U(\min_{i=1, \dots, n} X_{i,d}, \max_{i=1, \dots, n} X_{i,d})$ .

174

175 Then the dataset is split into two parts: (1)  $\{X_i | X_{i,d} < \textit{split value}; i = 1, \dots, n\}$  which is  
 176 passed to the left branch of the node; and (2)  $\{X_i | X_{i,d} \geq \textit{split value}; i = 1, \dots, n\}$  which is  
 177 passed to the right branch of the node. The procedure is repeated iteratively to create each node  
 178 of the tree until only one distinct instance remains in one node or reaches the height limit. The  
 179 EIF used an axis-obliqued splitting method to solve the issue that the split process of the  
 180 original IF will generate artifacts. Specifically, EIF creates a hyperplane with the form of a  
 181 point-norm equation to split the data as shown in Eq. [1]:

$$182 \quad (X - P) \cdot n^T = 0 \quad [1]$$

183 Where the point vector  $P = \{p_j | p_j \sim U(\min_i X_{i,j}, \max_i X_{i,j}); i = 1, \dots, n; j = 1, \dots, D\}$ , the  
 184 norm vector  $n = \{n_j | n_j \sim N(0,1)\}$ . Then the split rule becomes:  $\{X_i | (X - P) \cdot n^T < 0; i =$   
 185  $1, \dots, n\} \rightarrow \textit{left branch}; \{X_i | (X - P) \cdot n^T \geq 0; i = 1, \dots, n\} \rightarrow \textit{right branch}.$

186

187 The *iForest* algorithm assumes that the anomaly instances are rare. Such instances differ from  
 188 those deemed normal in a given data set, making them more susceptible to isolation in several  
 189 binary tree structures. In a random tree, instances are partitioned repeatedly until all instances  
 190 are isolated. In contrast, nominal instances require many more splits to finally reach their leaf  
 191 nodes (Li *et al.*, 2020). For a given dataset, the algorithm takes  $n$  random samples of size  $m$ . A  
 192 binary search tree is constructed for each random example, selecting a dimension and partition  
 193 point for each comparison node in the tree. The anomaly score of a new data point is calculated  
 194 by inserting it into each  $n$  random tree.

195

196 Isolation refers to the separation of an instance. Anomalous data has the nature of ‘few and  
 197 special’, and it is easy to isolate outliers from normal data. The *iForest* algorithm isolates data  
 198 by recursively and randomly partitioning. Usually, normal data is typically dense and needs to  
 199 be divided many times to be isolated. Conversely, abnormal data are outliers and only need to  
 200 be randomly divided a few times to be isolated. In the whole process of isolation, a binary tree  
 201 can represent the process of division. The earlier a point is divided, the more likely it is an  
 202 abnormal point. An example of the partitioning process is presented in Figure 1, where the ‘red’  
 203 leaf node is most likely an outlier.

204

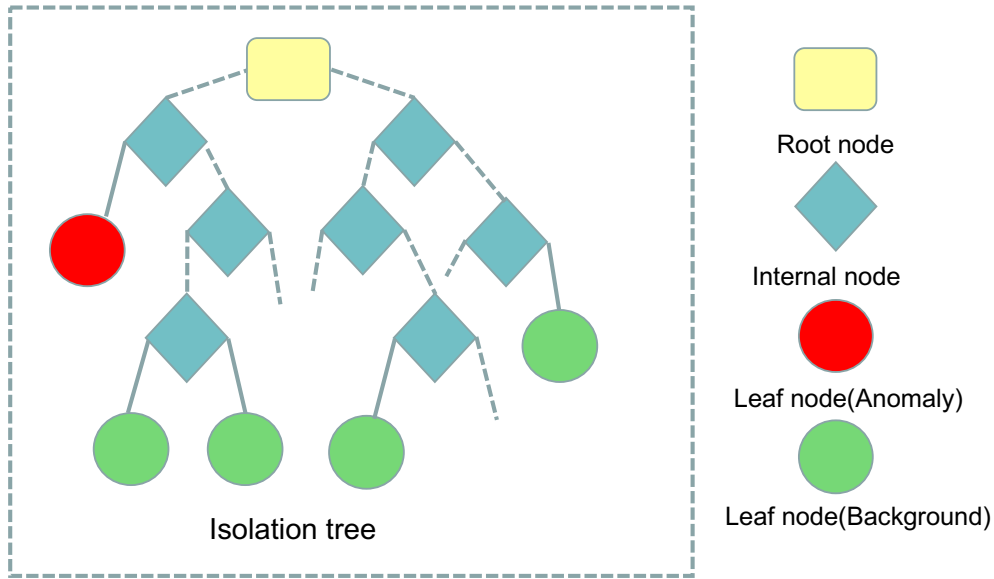


Figure 1. Example of the structure of an *iTree*

205  
206

207

208

209 The node ( $T$ ) of the isolation tree is either external with no child or internal with one test and  
 210 two daughter nodes ( $T_l, T_r$ ), where the number of external nodes is  $n$ , the number of internal  
 211 nodes is  $n-1$ , and the total number of nodes of an *iTree* is  $2n-1$  (Liu *et al.* 2012). A test consists  
 212 of an attribute  $q$  and a split value  $p$ . Given a database  $X = \{x_1, \dots, x_n\}$  of  $n$  instances from a  $d$ -  
 213 variate distribution, to build an *iTree*, we recursively divide  $X$  by randomly selecting an  
 214 attribute  $q$  and a split value  $p$ , until either: (1) the tree reaches a height limit, (2)  $|X| = 1$  or all  
 215 data in  $X$  have the same values (Liu *et al.* 2012).

216

217 The Path length ( $h(x)$ ) is determined by the number of edges  $x$  traverses an *iTree* from the root  
 218 node until the traversal is terminated at an external node. We borrow the analysis from Binary  
 219 Search Tree (BST) to estimate the average path length ( $E(h(x))$ ) of *iTree*. The anomaly score  
 220 ( $s$ ) of an instance  $x$  is defined as:

221

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad \text{Eq. [2]}$$

222

223 
$$c(n) = 2H(n - 1) - \left(\frac{2(n-1)}{n}\right) \quad \text{Eq. [3].}$$

224

225 Specific details of the assessment process can be found in Liu *et al.* (2008).

226

## 227 **2.2 De-noising Data`**

228 Due to the spatial-temporal uncertainty and complexity of working conditions in deep pit  
 229 foundations, raw monitoring data invariably contains noise. The noise will interfere with data  
 230 analysis and decision-making accuracy in this instance. Methods such as low-pass filtering (De  
 231 *et al.*, 2010), Wiener filtering (Aschero *et al.*, 2010), adaptive learning (Ortolan *et al.*, 2003),  
 232 and Kalman filtering (Singh *et al.*, 2018) are traditionally used to de-nose signals. Despite their  
 233 success, such approaches have limitations, as they filter out useful information or reduce  
 234 valuable features (Andrate *et al.*, 2006; Xiao *et al.*, 2019).

235

236 To address the above limitations, a wavelet transforms a time-frequency domain method has  
 237 been introduced to de-noise signals (Andrade *et al.*, 2006; Maier *et al.*, 2018). This method  
 238 comprises three steps: (1) signal decomposition; (2) detail coefficient thresholding; and (3)  
 239 signal reconstruction. When a wavelet transform is used, there is no requirement to incorporate  
 240 artificial components into the original signal (Andrade *et al.*, 2006; Maier *et al.*, 2018).  
 241 However, the limitation of such an approach is that the mother wavelet function must be pre-  
 242 defined. Additionally, the selection of different mother wavelet functions can affect detection  
 243 performance (Maier *et al.*, 2018; Xiao *et al.*, 2019).

244

245 Monitoring data is often mixed with multi-type noise. Thus, a novel signal analysis method is  
 246 needed to decompose a multi-component signal into several band-limited intrinsic mode

247 functions (BLIMFs). The VMD effectively determines the signal segmentation in the  
248 frequency domain and the components' separation. It has also been proven to simultaneously  
249 achieve accurate signal separation, better noise robustness, and higher computational efficiency.  
250 With this in mind, we will use the VMD in this research to de-noise the data.

251

252 Several studies have attempted to denoise the time series data using the traditional method for  
253 solving the signal denoising problem, which involves using linear time-invariant (LTI) filters  
254 (Selesnick *et al.*, 2014; Prateek *et al.*, 2021). An alternative approach uses wavelets; the main  
255 drawback of this approach is that it introduces pseudo-Gibbs artifacts at the singular points due  
256 to more local oscillations and smaller amplitude near signal discontinuities. And the sparsity-  
257 based methods, such as compressed sensing with dictionary elements from an oversampled  
258 discrete Fourier transform (DFT) matrix, cannot reconstruct the signal perfectly. The modal  
259 decomposition algorithm handles non-linear and non-smooth signals with good adaptive  
260 decomposition capability. It can decompose complex signals into intrinsic modal function  
261 forms sorted by frequency from high to low and extract the decomposed modal function to  
262 construct a filter. As a modal decomposition algorithm, VMD is selected in our research as it  
263 can separate tones of similar frequencies to represent time series characterization.

264

### 265 2.2.1 Variational Mode Decomposition

266 The VMD proposed by Dragomiretskiy and Zosso (2013) is a non-recursive decomposition  
267 method used for adaptive and quasi-orthogonal signal decomposition. It can simultaneously  
268 decompose a multi-component seismic trace into a finite number of band-limited intrinsic  
269 mode functions (IMFs). The VMD generalizes the classic Wiener filter into multiple adaptive  
270 bands. Wiener filtering is one of the most ubiquitous tools in signal processing, particularly for  
271 signal denoising and source separation. In the context of audio, it is typically applied in the

272 time-frequency domain using the short-time Fourier transform (STFT) (Samuel and James,  
 273 2008). The VMD algorithm is more robust to noise than the EMD-based adaptive  
 274 decomposition methods (Dragomiretskiy and Zosso, 2013). The concepts and theories related  
 275 to VMD are as follows.

276

277 **Definition 1:** (Intrinsic Mode Function)

278 Intrinsic Mode Functions are amplitude-modulated-frequency-modulated (AM-FM) signals,  
 279 which differs from the definition of EMD.

280

$$281 \quad \mu_k(t) = A_k(t) \cos(\phi_k(t)) \quad [4]$$

282

283 Where the phase  $A_k(t)$  is an envelope of  $\mu_k(t)$  and  $\phi_k(t)$  is a non-decreasing function. The  
 284 equation of phase  $\phi_k(t)$  and instantaneous frequency  $\omega_k(t)$  is as follow:

285

$$286 \quad \omega_k(t) = \frac{d\phi_k(t)}{dt} \geq 0 \quad [5]$$

287

288 **Definition 2:** (Total Practical IMF Bandwidth)

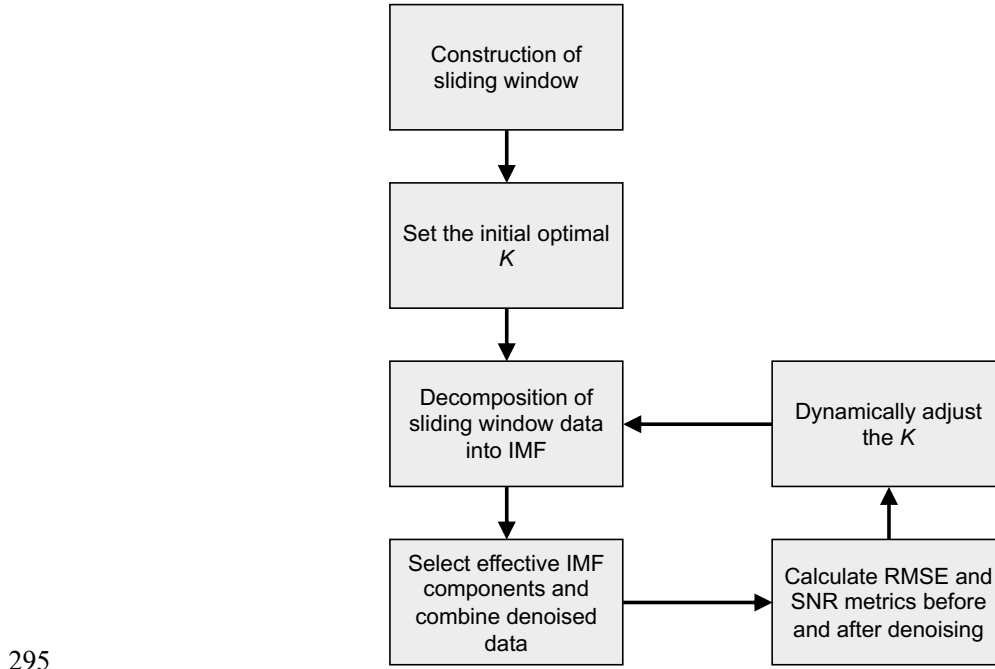
289 The total practical bandwidth of an IMF is estimated as Eq. [6]. Depending on the actual IMF,  
 290 either of these terms may be dominant.

291

$$292 \quad BW_{AM-FM} = 2(\Delta f + f_{FM} + f_{AM}) \quad [6]$$

293

294 The workflow of the VMD, which we will follow, is presented in Figure 2.



299 Figure 2. Workflow of implementation of VMD

300 The VMD comprises the following steps:

301 *Step 1: modes update.* The modes  $(\hat{u}_k^{n+1}(w))$  are updated by Eq. [7], where  $\alpha$  is the penalty  
302 factor,  $\lambda$  is the Lagrangian multiplier,  $\mu_k$  is the IMF. Wiener filtering is embedded for an  
303 update as the mode directly in the Fourier domain with a filter tuned to the current center  
304 frequency  $w_k^n$ ;

305

306 
$$\hat{u}_k^{n+1}(w) = \frac{\hat{f}(w) - \sum_{i < k} \hat{u}_i^{n+1}(w) - \sum_{i > k} \hat{u}_i^n(w) + (\hat{\lambda}^n(w)/2)}{1 + 2\alpha(\omega - \omega_k^n)^2} \quad \text{Eq. [7]}$$

307

308 *Step 2: Center frequencies update.* The center frequencies are updated as the center of gravity  
309 of the corresponding mode's power spectrum, as shown in Eq. [8]

310

$$\omega_k^{n+1} = \frac{\int_0^\omega \omega |\hat{u}_k^{n+1}(w)|^2 d\omega}{\int_0^\omega |\hat{u}_k^{n+1}(w)|^2 d\omega} \quad \text{Eq. [8]}$$

312

313 *Step 3: Dual ascent update.* For all  $\omega \geq 0$ , the Lagrangian multiplier  $\hat{\lambda}^{n+1}$  is updated by Eq.  
 314 [9] as a dual ascent to enforce exact signal reconstruction until  $\sum_k \|\hat{u}_k^{n+1}(w) - \hat{u}_k^n(w)\|_2^2 /$   
 315  $\|\hat{u}_k^n\|_2^2 < \varepsilon$ .

316

$$\hat{\lambda}^{n+1} = \hat{\lambda}^n + \tau(\hat{f} - \sum_k \hat{u}_k^{n+1}) \quad \text{Eq. [9]}$$

318

319 Additional details about the VMD can be found in the works of Dragomiretskiy and Zosso  
 320 (2013).

321

### 322 **3.0 Research Approach**

323 To recap, our research aims to develop a smart data approach to detect anomaly monitoring  
 324 data and reduce noise to improve the quality of monitoring data extracted during the  
 325 construction process of hazardous activities such as deep pit foundations. Our smart data  
 326 approach extracts data relevant for decision-making to determine safety risks and consists of  
 327 EIF and VMD. In the process of data collection, it is inevitable to produce some data that  
 328 deviates from the rest of the observations in the sample to which it belongs. The reasons mainly  
 329 include: (1) the failure of the equipment; (2) the abnormality of the collected data caused by  
 330 the dynamic working environment.

331

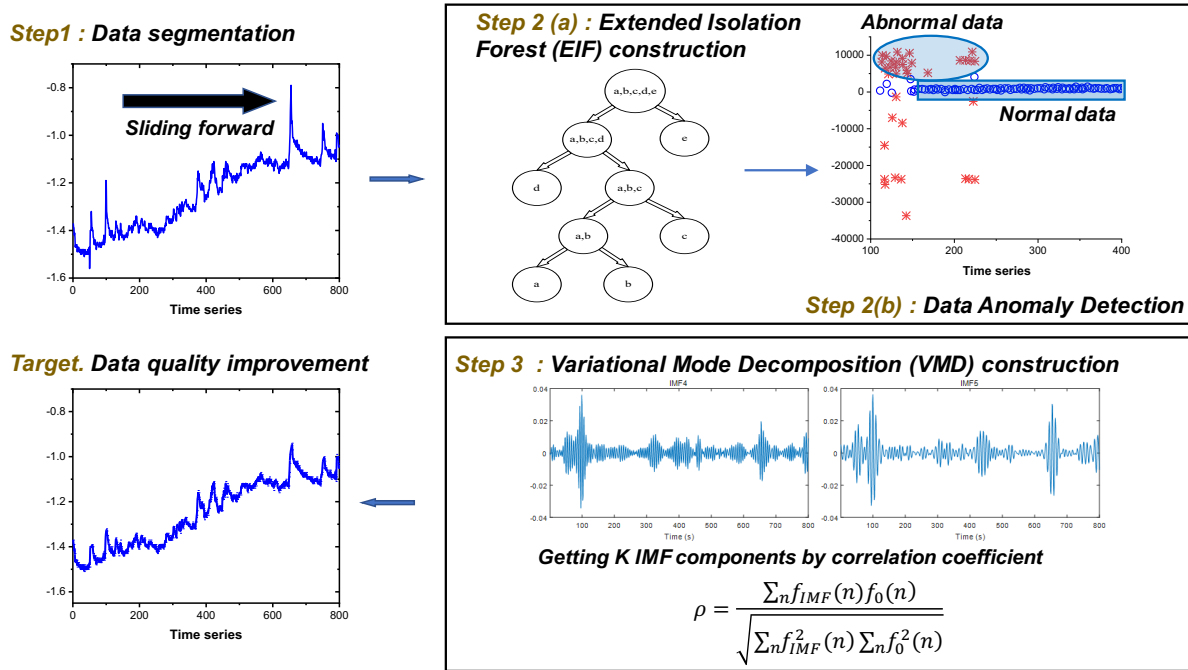
332 To obtain a high-quality time-series monitoring dataset, it is necessary to perform abnormal  
 333 processing on the data. While noise and outliers are similar in their statistical distribution and  
 334 characteristics, they originate from fundamentally different causes. The workflow of our



335 proposed method is presented in Figure 3. The research process we have adopted to develop  
336 our hybrid smart data approach consists of the following three steps (Figure 3):

337

- 338 • *Step 1 - Data segmentation:* Extracted monitoring data is divided into segments using a  
339 rectangle sliding window. Then, various numerical features, such as root mean square  
340 (RMS) and kurtosis of each data window, are determined.
- 341 • *Step 2(a) - Extended isolation forest construction:* The EIF is an outlier detector that  
342 builds an ensemble of *i*Trees for a given dataset. The EIF resolves the issues associated  
343 with assigning anomaly scores to given data points by using hyperplanes with random  
344 slopes (non-axis-parallel) to split data to create *i*Trees (Hariri *et al.*, 2021)
- 345 • *Step 2(b) - Data anomaly detection:* In an EIF, data are subsampled and processed in a  
346 tree structure based on random cuts in the values of arbitrarily selected features in each  
347 dataset. Each tree is grown until each instance is isolated into a leaf node. The samples  
348 with shorter branches indicate anomalies.
- 349 • *Step 3 - Variational model decomposition construction:* As an adaptive signal processing  
350 method, the VMD removes harmonic noise and improves data quality. The VMD  
351 algorithm concurrently decomposes the input signal into several narrow-band modes.  
352 Each mode is band-limited around its center frequency, which leads to less spectral  
353 overlapping or instantaneous frequency fluctuation is observed in the VMD results. The  
354 VMD algorithm decomposes and reconstructs the monitoring data to achieve adaptive  
355 signal decomposition and noise reduction.



356

357

358

359 We now explain in greater detail our research approach for detecting anomalies and de-noising  
 360 the engineering data extracted from sensors used to monitor deep pit foundations by focusing  
 361 on EIF and VMD.

362

#### 363 4.0 Case Study

364 We use an explanatory case study to demonstrate and validate our hybrid smart data approach  
 365 (Dubé and Paré, 2003). A deep foundations pit of a subway project in Wuhan, China, while under  
 366 construction, is selected. The project was chosen as sensors were used to monitor geotechnical  
 367 safety risks, and the researchers worked closely with contractors on several other studies.

368

#### 369 4.1 Case Description

370 The selected subway project is a T-shaped transfer between stations A and B. Subway station  
 371 A is an underground three-story double-column 13m island platform station. The total length  
 372 of the station is 239.2m, the full width of the standard section is 22.5m, the structure height is

373 22.63m-25.08m, and the roof is buried deep about 3m-4.1m, both ends of the station are shield  
 374 tunnel receiving wells. Subway station B is a 14-meter island-style station with two  
 375 underground floors and two columns. The total outsourcing of the station is 634.105m, and the  
 376 full width of the standard section is 23.1m. The landform can be classified as a denudation  
 377 accumulation ridge area (grade III terrace), and the ground elevation of the exploration area is  
 378 between 26.0 and 30.7m.

379



380

381

382 Figure 4. Example of deep pit foundation

383

#### 384 4.2 *Experimental Set-up*

385 We first install sensors while excavations are constructed to conduct this experiment, focusing  
 386 explicitly on the supporting shaft's axial forces and building settlement. The layout of these  
 387 sensors is presented as follows:

388

- 389 • *Support shaft axial force monitoring*: The axial force meter monitors the axial force of  
 390 steel support. The meter is installed at the end of the steel support. In this case, four axial  
 391 force monitoring points are established on the fourth and fifth layers of steel supports, with  
 392 two monitoring points on each layer. The monitoring equipment used is the Vibrating

393 String Axial Force Meter produced by the Shenzhen JingSheng Tech Co., LTD, model  
394 MAS-AXF-40. The maximum range is 4000KN, a MAS of 0.1%, and a precision of  $\pm 0.5$ .

395 • *Building settlement monitoring:* Based on the experience of experts, four-building  
396 settlement monitoring points were installed symmetrically on the four corners of the  
397 building, closest to the foundation pit. The monitoring equipment used is Photoelectric  
398 Static Level produced by TongWei Sensing, model ESJS-50. The measuring range is  
399 50mm, and the precision is  $\pm 0.1$ mm.

400

401 Examples of sensors installed in the case are presented in Figure 5. After the sensors are  
402 installed, the data are transmitted and stored in the web-based monitoring system, as shown in  
403 Figure 6.

404

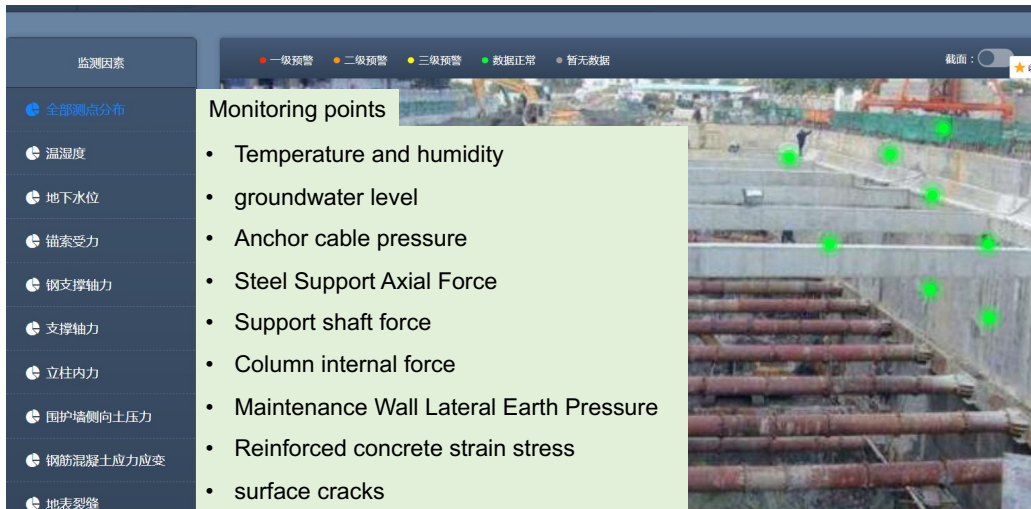


405

406

407

Figure 5. Examples of sensors installed in foundation pit



408

409

Figure 6. Web-based monitoring data system

410

### 411 4.3 Anomaly Monitoring Data Detection

412 The four installed sensors act as the monitoring points and our data source, one of which is  
 413 used as the baseline, and the other three monitoring points (CJ1, CJ2, CJ3) are for analysis. As  
 414 noted in Figure 7a, the settlement data has no apparent abnormality; The anomaly data  
 415 detection under one (i.e., CJ1), two (i.e., CJ1 and CJ2), and three (CJ1, CJ2, and CJ3)  
 416 dimensional analysis. An anomaly value is calculated for each point under different training set  
 417 sizes in anomaly data detection. The distribution of anomaly value is used to analyze the effect  
 418 of dimension selection and dataset size on anomaly detection. The results of anomaly  
 419 monitoring data detection under different training sets and dimensional analysis are presented  
 420 in Figure 7b.

421

422 Here we define the data with anomaly scores higher than 0.6 as outliers and analyze the  
 423 anomaly scores of the outliers. Figure 7b shows that one-dimensional data has a higher anomaly  
 424 score than two-dimensional and three-dimensional data. As the sizes of the training set change,  
 425 the anomaly scores of the one-dimensional data also change, but these values generally exceed

426 0.78. Conversely, the anomaly scores of two-dimensional and three-dimensional data are lower  
427 than 0.74. The maximum value of the anomaly score of the two-dimensional data fluctuates  
428 between 0.71 and 0.74. The maximum value of the anomaly score of the three-dimensional  
429 data ranges between 0.69 and 0.71. The anomaly scores for outliers in high-dimensional  
430 datasets are more concentrated with lower anomaly scores. Thus, we can conclude that the  
431 *iForest* algorithm can process high-dimensional data (i.e., settlement monitoring data). The  
432 detection of outliers is smoother than the low-dimensional data, and the abnormal value is  
433 relatively lower.

434

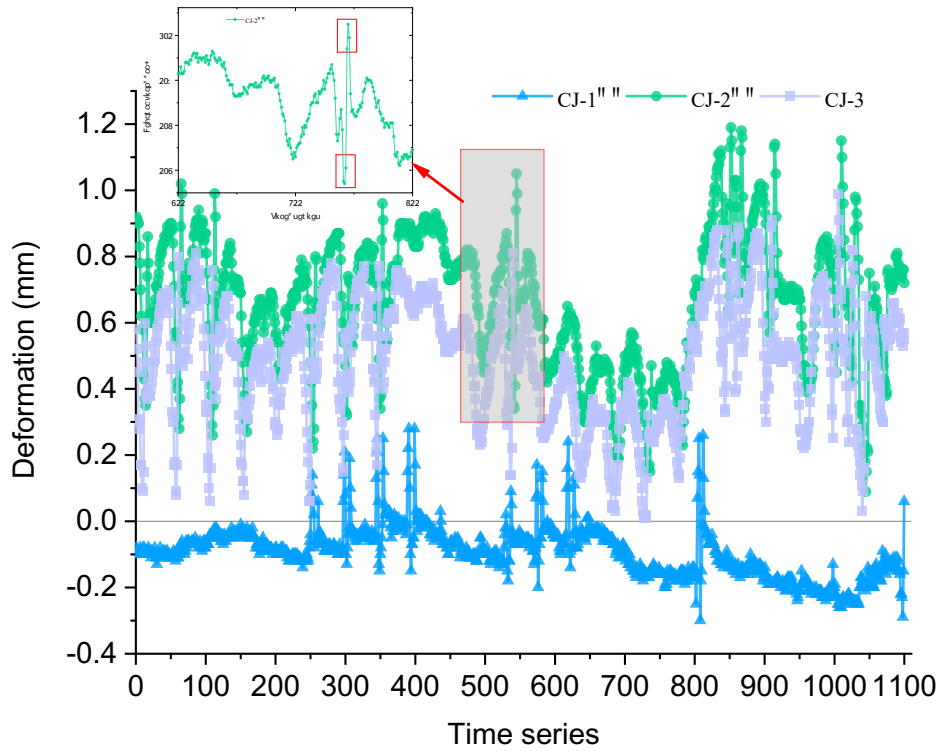
435 Again, four monitoring points are selected to analyze the steel support axial forces (ZCL-02-  
436 21, ZCL-02-22, ZCL-04-C6, ZCL-04-C7). Our results are presented in Figure 8a, and we can  
437 conclude that the monitoring data of ZCL-04-C6 is abnormal. As a result, we then analyzed the  
438 anomaly data detection under different dimensional conditions, with the results being presented  
439 in Figure 8b. As seen from Figures 8 and 9, in the detection data of the *iForest*, the higher the  
440 dimension of monitoring data, the less sensitive it is to detecting anomalies. We can find the  
441 difference between single and multi-dimensional anomalies by analyzing high-dimensional  
442 monitoring data. The higher the dimension of monitoring data, the higher the anomaly value,  
443 and the easier it is to determine the cause of the abnormality.

444

#### 445 **4.4 Evaluation Performance**

446 We compared the EIF with the *iForest* algorithm to determine which method can better identify  
447 abnormal points. Figure 9 shows the two-dimensional abnormal point detection of the steel  
448 support axial force monitoring data using EIF and a standard *iForest* algorithm. The left and  
449 right columns are the standard isolation forest and EIF algorithms.

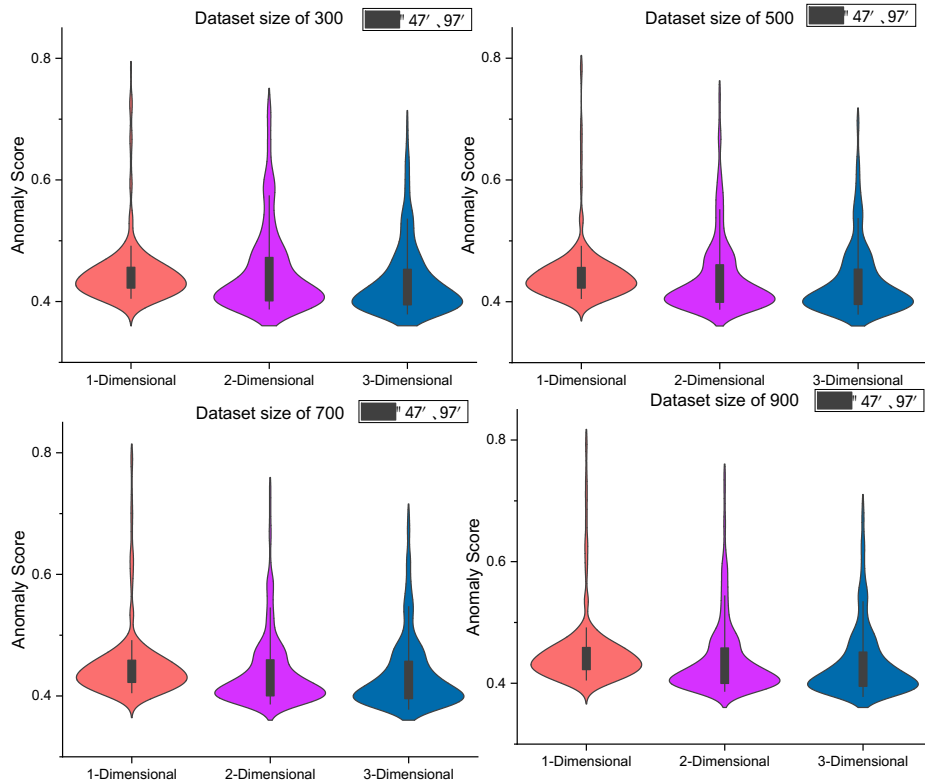
450



451

452

(a) Examples of monitoring data



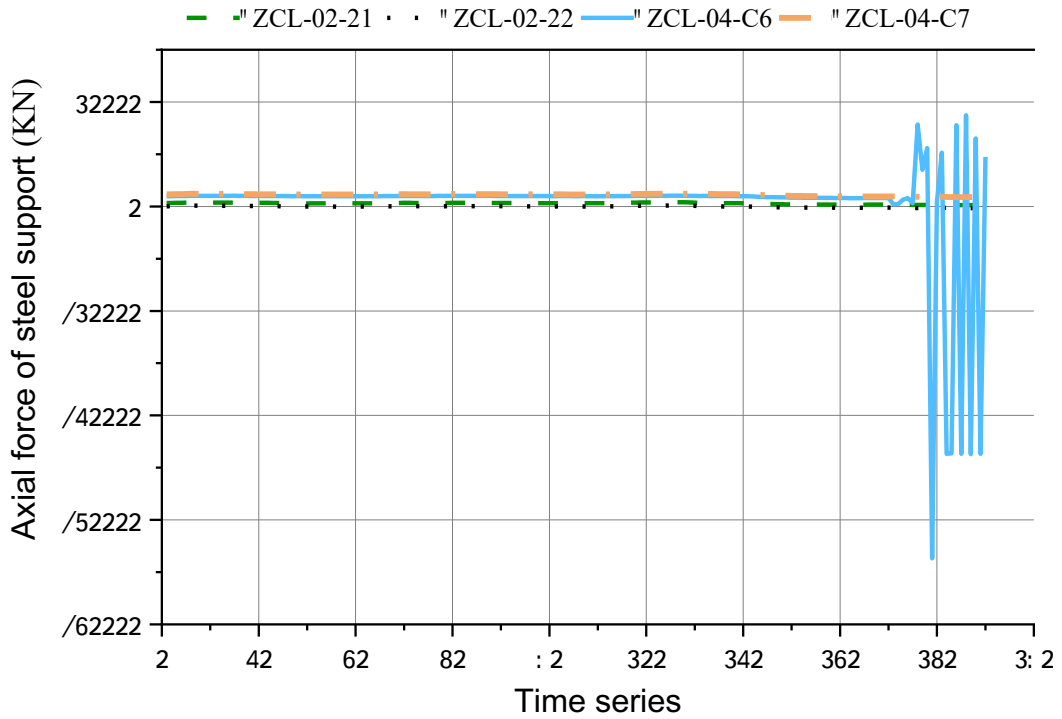
453

454

(b) Settlement data: different dimensional and training set sizes

455

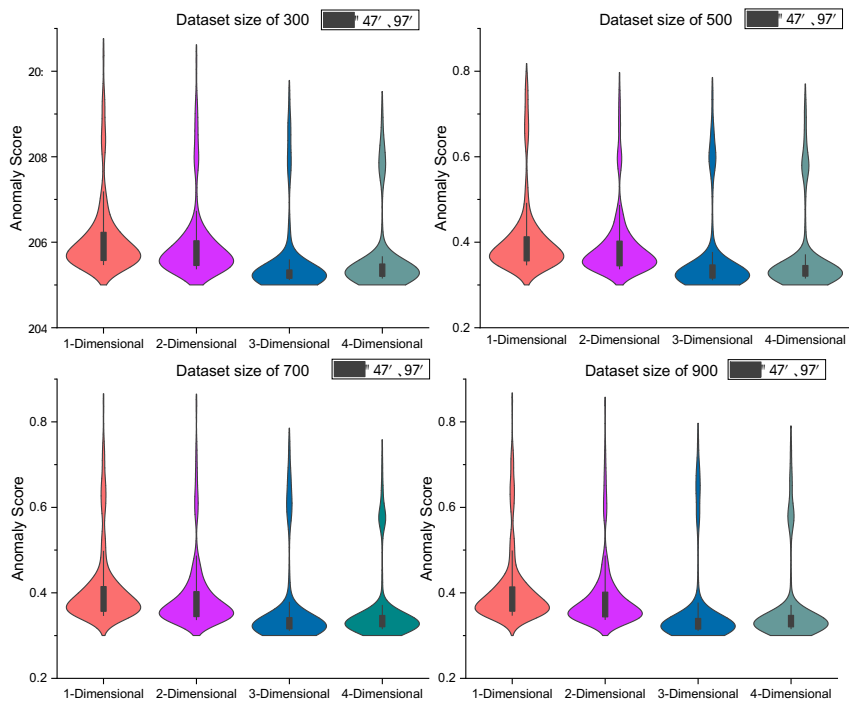
Figure 7. Examples of monitoring and settlement data



456

457

(a) Examples of monitoring data



458

459

(b) Steel support axial force: different dimensional and training set sizes

460

Figure 8. Examples of monitoring steel support axial force data



461 Figure 9 shows that when the anomaly score is 0.75, all the anomalous data cannot be identified.  
462 We suggest that the reason is that the training set in this figure includes many abnormal data,  
463 and the data selection for model training cannot directly take the sliding window of a time point.

464

#### 465 4.4.1 Quality of Training Database

466 The detection of anomalies using the EIF algorithm consists of two phases:

467

- 468 1. *Training*: An isolated tree is built based on subsamples of the training set;
- 469 2. *Testing*: An isolated tree calculates anomaly scores for each test sample. Hence, we  
470 design two group experiments to conduct this test: (i) a training database; and (ii) a  
471 training dataset without anomaly monitoring data. The size of the training database of  
472 these two group experiments is set to 800. The results are presented in Figure 10.

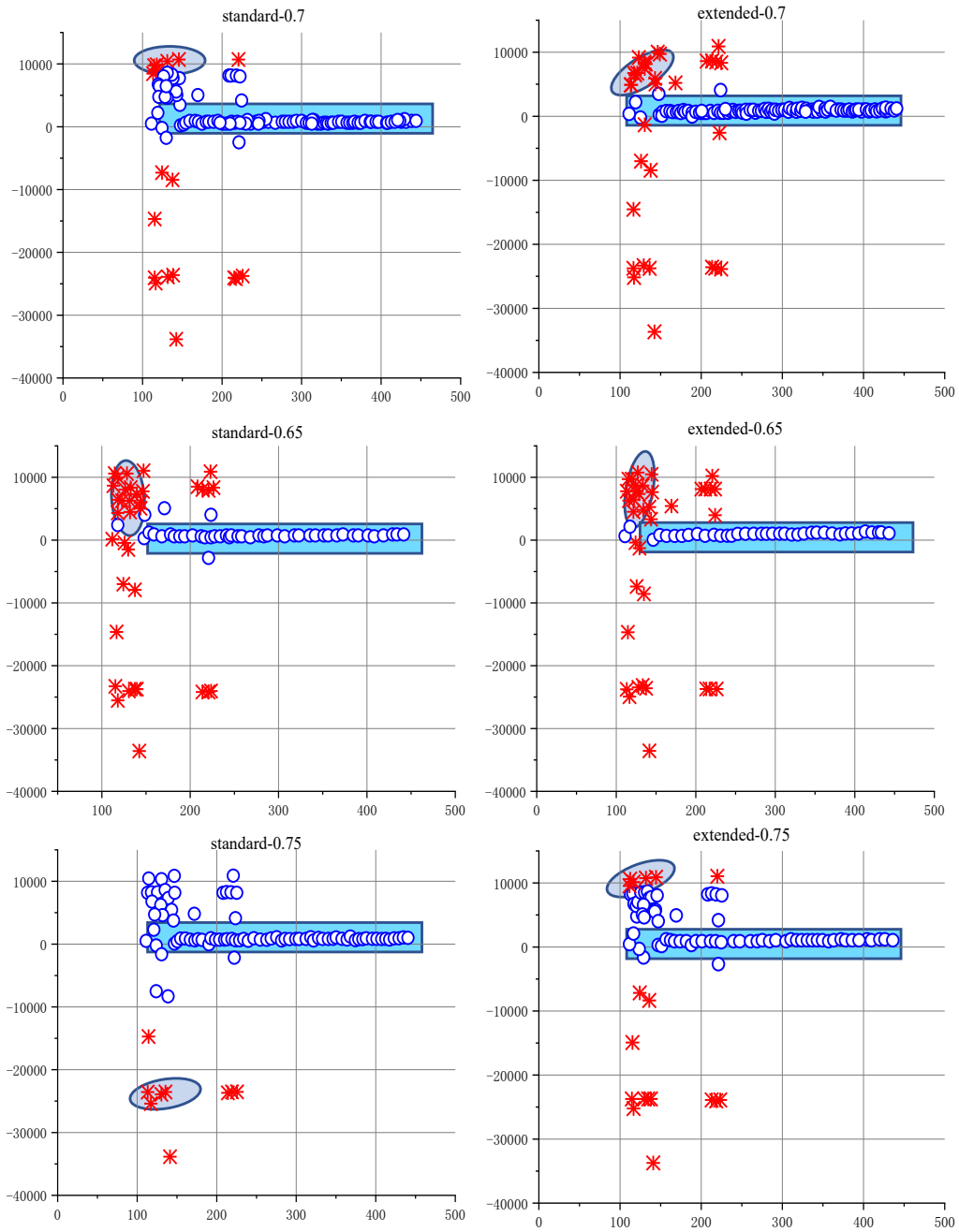
473

474 Figure 10 shows that the training set that excludes abnormal monitoring data achieves better  
475 performance on abnormal data detection. We process the data by setting a threshold for the  
476 anomaly score; that is, data with an anomaly score over 0.7 will be removed. During data  
477 analysis, we conclude that the different performance is from the abnormal points not excluded,  
478 leading to the other points being no longer 'isolated' as the other data set. Therefore, in  
479 detecting anomalies in *i*Forests, the quality of the training set should be maintained.

480

#### 481 4.4.2 Size of the Training Database

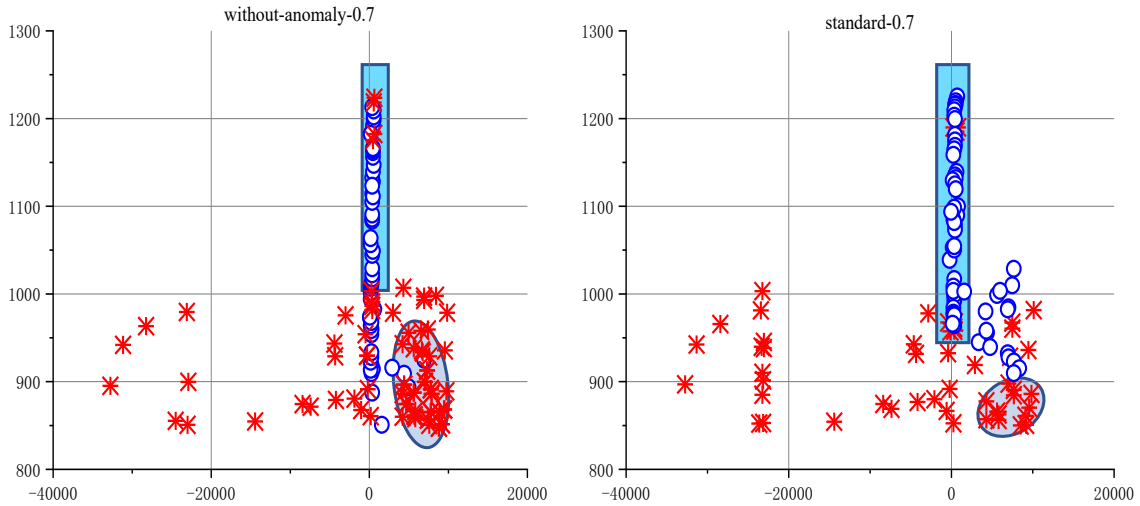
482 In this experiment, the training set size ranges from 300 to 800. It is trained with a gradient of  
483 100, with the abnormal detection results presented in Figure 11. Here we can see that with  
484 increased increments in the training set, the algorithm achieves better performance on anomaly  
485 detection until the 600 mark, where almost all anomalies are detected.



486

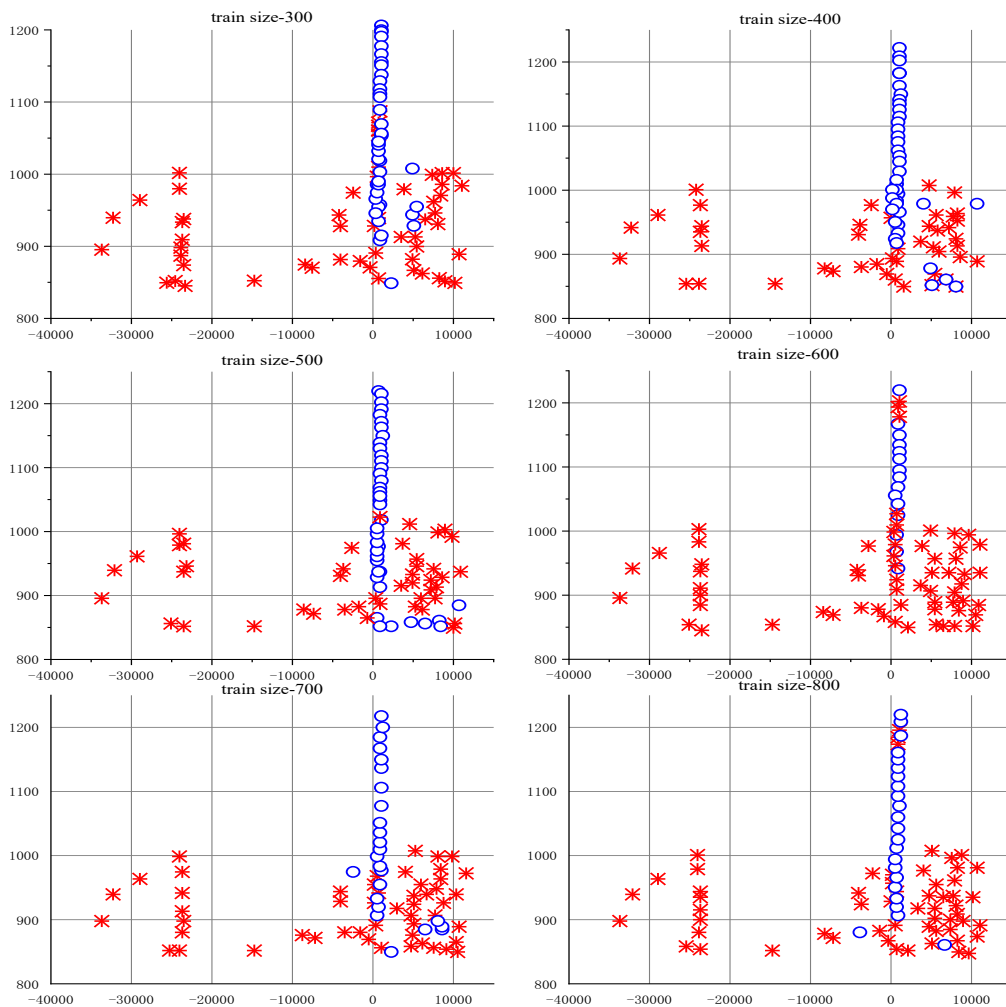
487

Figure 9. The comparison results of EIF and *i*Forest algorithms



488  
489  
490  
491

Figure 10. The comparison result of using two training databases



492  
493

Figure 11. A comparison of the results for different training databases

## 494 4.4.3 Effectiveness of EIF

495 We use the steel support axial force dataset to verify the effectiveness of the EIF algorithm  
 496 compared with the KNN algorithm and the ABOD (angle-based outlier detection) algorithm.  
 497 We calculate the results using a two-dimensional dataset and a four-dimensional dataset,  
 498 respectively, for analysis. The evaluation measures are area under the curve (AUC) and  
 499 Accuracy. The AUC is an index normally used to evaluate the efficiency of classifiers, defined  
 500 as the area under the receiver operating characteristic (ROC) curve. Accuracy is the proportion  
 501 of instances in the monitoring data detected correctly. The results are shown in Table 2.

502

503 Table 2. AUC and Accuracy of EIF, KNN, and ABOD

Algorithm	AUC	Accuracy
ABOD(2-Dimensional)	0.8629	89.5
EIF(2-Dimensional)	0.8893	94.4
KNN(2-Dimensional)	0.8827	93.7
ABOD(4-Dimensional)	0.8964	94.9
EIF(4-Dimensional)	0.8997	96.3
KNN(4-Dimensional)	0.8634	89.6

504

505 As can be seen from Table 2, the AUC and Accuracy of the four algorithms are basically  
 506 consistent, and there is no noticeable difference. The accuracy of EIF and AMOD improved  
 507 with the increase of the dataset dimension, while the accuracy of KNN decreased. The  
 508 experimental results show that the EIF algorithm can effectively improve the execution  
 509 efficiency of anomaly detection with a high-dimensional dataset. Therefore, EIF is suitable for  
 510 anomaly detection on large-scale monitoring data.

511

## 512 4.5 Monitoring Data De-noising

513 We select a dataset containing 1000 monitoring points based on the period and frequency of  
514 data collection for denoising sample data. Some studies found that the decomposition number  
515  $K$  significantly influences the decomposition results (Xia *et al.* 2021; Wang *et al.* 2019). When  
516 the number of decomposition modes ( $K$ ) is too low, under-decomposition will occur, and some  
517 ‘modes’ cannot be recognized effectively (Xia *et al.* 2021; Wang *et al.* 2019). When  $K$  is too  
518 large, a particular ‘mode’ in the signal may be ‘pulled’ into multiple IMF components, resulting  
519 in excessive decomposition (Li *et al.*, 2019). In this study, the number of decompositions  
520 includes 1 to 9 in advance to choose the best  $K$ . The waveforms of the nine IMFs decomposed  
521 by the VMD are presented in Figure 12.  $K$  is the number of modes. In the original EMD  
522 description, a mode is defined as a signal whose number of local extrema and zero-crossings  
523 differ at most by one. In later related works, the definition is slightly changed into so-called  
524 Intrinsic Mode Functions (IMF); through the decomposition of different  $K$  ensemble members,  
525 the correlation coefficient between each IMF component and the sample data is calculated,  
526 with the results being presented in Table 3.

527

528 From Table 3, we can conclude that when  $K$  is 2 and 3, the decomposed IMF components are  
529 valid according to the threshold value. It indicates that the dataset is underpinning, and the  
530 high-frequency noise is not isolated. So, these two modes are not analyzed later. We set the  
531 threshold as 0.1 of the value, which is the largest of all correlation coefficients. When the value  
532 of IMF is less than the threshold, we define the value as a failure. (Yu 2008) When  $K$  is 4~9,  
533 the effective IMFs are all 3. The first three IMF components are reconstructed. The  
534 reconstructed signal and the original data signal are calculated by the root mean square error  
535 (RMSE) and the signal-to-noise ratio (SNR). The RMSE and SNR are defined by Eq. [6] and  
536 Eq. [7]:

537 
$$RMSE = \sqrt{\frac{1}{n} \sum_n (f_0(n) - f_1(n))^2}$$
 Eq. [6]

538

539 
$$SNR = 10 \times \log_{10} \left( \frac{\frac{1}{n} \sum_n f_0^2(n)}{\frac{1}{n} \sum_n (f_0(n) - f_1(n))^2} \right)$$
 Eq. [7]

540

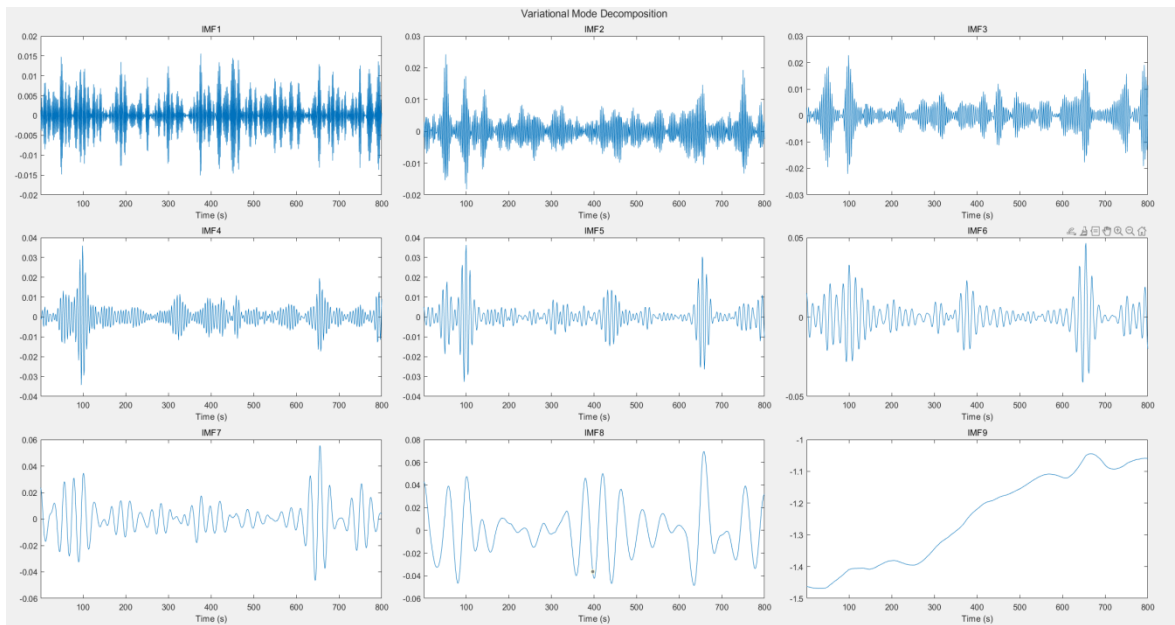
541 Where,  $f_0$  is the original signal data,  $f_1$  is the reconstructed signal data.

542 We can conclude that when  $K$  is 4, the RMSE is the smallest, and the signal-to-noise ratio SNR

543 is the largest, so the denoising effect is the best, and the optimal  $K$  value should be selected as

544 4.

545



546

547 Figure 12. Waveforms of the nine IMFs decomposed by the VMD

548

549

550

551

552

553

554 Table 3. Correlation coefficients between IMF and original signal data under different  $K$   
 555 ensemble members

556

$K$	IMF1	IMF2	IMF3	IMF4	IMF5	IMF6	IMF7	IMF8	IMF9
2	0.987	0.156							
3	0.982	0.179	0.119						
4	0.981	0.181	0.115	0.090					
5	0.980	0.181	0.111	0.085	0.075				
6	0.980	0.183	0.108	0.083	0.076	0.065			
7	0.980	0.184	0.107	0.081	0.074	0.065	0.047		
8	0.980	0.185	0.106	0.077	0.067	0.067	0.058	0.044	
9	0.980	0.185	0.106	0.077	0.065	0.063	0.059	0.046	0.038

557

558 Table 4. RMSE and signal-to-noise ratio under different  $K$  values

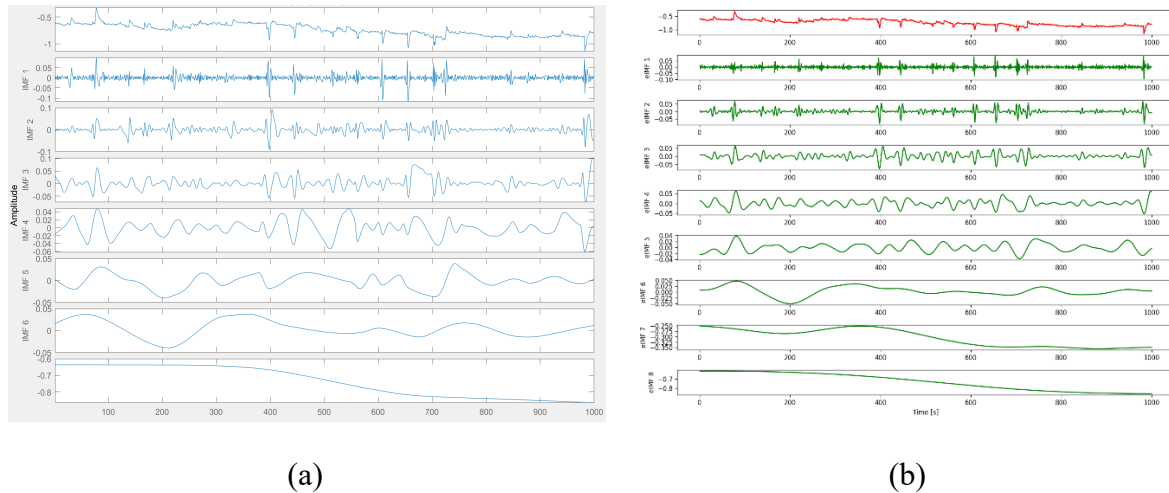
559

$K$	4	5	6	7	8	9
RMSE	0.0389	0.0413	0.0438	0.0442	0.0453	0.0455
SNR	24.09	23.58	23.07	22.98	22.77	22.74

560

561 We use the same dataset to verify the effectiveness of the VMD algorithm compared with the  
 562 EMD algorithm and the EEMD (ensemble EMD) algorithm. The evaluation measures are  
 563 RMSE and SNR. The decomposition result of the EMD and EEMD for the signal is presented  
 564 in Figure 13, and the results are shown in Table 5. As can be seen from the results, EEMD  
 565 denoising is superior to EMD, and VMD denoising is better than EEMD, which has high SNR  
 566 and low RMSE.

567



568 Figure 13. Decomposition result by EMD(a) and EEMD(b)

569  
570 Table 5. RMSE and signal-to-noise ratio of different algorithm

	VMD	EMD	EEMD
SNR	24.09	14.26	15.91
RMSE	0.0389	0.1573	0.1479

571  
572 **5.0 Discussion**

573 It has been suggested that big data analytics provides the basis to identify patterns and derive  
574 insights about safety issues in construction (Guo *et al.*, 2016; Fang *et al.*, 2020; Fang *et al.*,  
575 2021; Liu *et al.*, 2023). However, despite the espoused benefits of big data and there has been  
576 an increasing drive for construction organizations to embrace and apply its dimensions in their  
577 respective projects (Ngo *et al.*, 2020), its adoption should be treated with a degree of skepticism  
578 “as big data is not always better data” (Ghasemaghahi and Calic, 108: p.147). Many  
579 construction organizations remain unprepared to effectively utilize big data derived from  
580 sensors for assessing geotechnical safety risks (Matthews *et al.*, 2022).

581  
582 We suggest that the outcomes of our research can support decision-making in identifying



583 unsafe conditions based on big data. Our employed EIF, an anomaly detection method, is used  
584 to effectively identify anomalous data and retain the normal fluctuation characteristics within  
585 its time series. It can be helpful for subsequent data processing and provide high-quality data  
586 sources for subsequent data analysis. In addition, the denoise processing of monitoring data  
587 significantly reduces data errors and improves the accuracy of identifying unsafe conditions.  
588 Hence, the motivation to develop our hybrid smart data approach is to use monitoring data  
589 extracted from sensors to help construction organizations assess geotechnical safety risks.

590

591 One of the challenges is that a significant amount of data collected from sensors used to detect  
592 geotechnical conditions contains noise, rendering it challenging to determine the correct  
593 information needed to train algorithms and undertake risk analysis. In addressing this void, we  
594 have developed a hybrid smart data approach that can detect noise and de-noise data extracted  
595 from sensors monitoring a building's geotechnical conditions, impacting its structural safety.  
596 The contributions of our research are twofold.

597

598 Firstly, we have developed an EIF approach to detect noise in monitoring data. Existing  
599 anomaly detection algorithms detect anomalies by understanding the distribution of their  
600 properties and isolating them from a normal data sample. Our employed EIF uses a model-free  
601 algorithm that does not rely on building a profile for data to find non-conforming samples.  
602 Instead, it utilizes anomalous data with various characteristics compared with normal data  
603 samples. In this instance, our employed EIF has computationally efficient and high accuracy  
604 without a profile of normal instances demonstrated in this case.

605

606 Secondly, a VMD approach and dynamic threshold processing are used to de-noise the  
607 monitoring data to improve its validity. The value of the  $K$  has an important influence on

608 decomposing the data, as shown in Tables 2 and 3, which can prevent under-decomposition  
609 and over-decomposition problems of VMD (Dragomiretskiy and Zosso, 2013). Reconstructing  
610 the IMF components can effectively decompose the original data. By calculating the RMSE  
611 and SNR of the original data and the reconstructed signal, we get the optimal mode number of  
612 VMD. This means that when we reuse VMD for new applications, we need to pay attention to  
613 the setting of the  $K$  value to improve the effectiveness of data noise reduction.

614

### 615 **5.1 Limitations**

616 Despite the novelty of our research, it needs to be acknowledged that several limitations exist.  
617 The study was limited to a single project in the Wuhan subway and two types of monitoring  
618 geotechnical data (i.e., building settlement and steel support axial forces). Future research,  
619 therefore, is required to examine the generalisability of our approach in different projects and  
620 a broader range of activities that use sensors to monitor the geotechnical conditions that  
621 influence structural components. In addition, the experimental results demonstrate that the  
622 proposed method performs satisfactorily (i.e., RMSE and SNR are 0.0389 and 24.09,  
623 respectively). However, we did not conduct comparative experiments to evaluate our hybrid  
624 smart data approach's performance (i.e., accuracy and computational efficiency) with other  
625 state-of-the-art measurement methods (e.g., deep learning-based). We suggest this limitation  
626 can be addressed by conducting additional experiments in our future work.

627

### 628 **6.0 Conclusion**

629 Anomaly identification and denoising are necessary tasks to improve the quality of monitoring  
630 data extracted from sensors in construction. Our research aims to develop a novel smart data  
631 approach that can effectively detect anomalies and de-noise monitoring data to improve its  
632 quality to assess geotechnical safety risks. Our smart data approach consists of an:

633 • Extended Isolation Forest algorithm, which extracts features from each monitoring  
634 dataset and is used to identify abnormal points; and

635 • Variational Mode Decomposition to remove harmonic noise, thus improving data quality.

636

637 A case of the Wuhan subway project is used to validate the effectiveness and feasibility of our  
638 proposed approach. The results demonstrated that by applying EIF and VMD, a high degree of  
639 accuracy could be achieved in detecting anomalies and denoising data. Our results show that  
640 our new method can detect anomalies with an RMSE and SNR are 0.0389 and 24.09,  
641 respectively. It was revealed that the EIF and VMD could accurately detect anomalies and de-  
642 noise monitoring data.

643

644 Even though our approach could not recognize all anomalies, our hybrid smart data approach  
645 can provide site management to improve their ability to assess geotechnical safety risks.  
646 Furthermore, we suggest that our approach can improve the quality of data extracted from  
647 sensors in deep foundation pits with minimal error. Thus, our proposed novel smart data  
648 approach effectively reduces noise from monitoring data extracted from sensors.

649

## 650 **Acknowledgment**

651 The authors would like to acknowledge the financial support of the Alexander von Humboldt  
652 Foundation and the National Natural Science Foundation of China (U21A20151, 51978302).

653

## 654 **References**

655 Abbate, A. Koay, J, Frankel, J., Schroeder, S.C., and Das, P. (1997) Signal detection and noise  
656 suppression using a wavelet transform signal processor: Application to ultrasonic flaw  
657 detection. *IEEE Transactions on Ultrasonics Ferroelectrics and Frequency Control*. **44**(1),

658 pp.14-26.

659 Ahmed, M., Mahmood, A. N., and Hu, J. (2016). A survey of network anomaly detection  
660 techniques. *Journal of Network and Computer Applications*, **60**, pp.19-31.

661 Andrade, A. O., Nasuto, S., Kyberd, P., Sweeney-Reed, C. M., and Van Kanijn, F. R. (2006).  
662 EMG signal filtering is based on empirical mode decomposition. *Biomedical Signal*  
663 *Processing and Control*, **1**(1), pp.44-55.

664 Aschero, G., and Gizdulich, P. (2010). Denoising of surface EMG with a modified Wiener  
665 filtering approach. *Journal of Electromyography and Kinesiology*, **20**(2), pp.366-373.

666 Asadzadeh, A., Arashpour, M., Li, H., Ngo, T., Bab-Hadiashar, A., and Rashidi, A. (2020).  
667 Sensor-based safety management. *Automation in Construction*, **113**, 103128.

668 Bao, Y., Tang, Z., Li, H., and Zhang, Y. (2019). Computer vision and deep learning-based data  
669 anomaly detection method for structural health monitoring. *Structural Health Monitoring*,  
670 **18**(2), pp.401-421.

671 Banjanovic-Mehmedovic, L., Hajdarevic, A., Kantardzic, M., Mehmedovic, F., and  
672 Džananovic, I. (2017). Neural network-based data-driven modelling of anomaly detection  
673 in thermal power plant. *Automatika: časopis za automatiku, mjerenje, elektroniku,*  
674 *računarstvo i komunikacije*, **58**(1), pp.69-79.

675 Bhavsar, Y. B., and Waghmare, K. C. (2013). Intrusion detection system using data mining  
676 technique: Support vector machine. *International Journal of Emerging Technology and*  
677 *Advanced Engineering*, **3**(3), pp.581-586

678 Braun, S. (2011). The synchronous (time-domain) average revisited. *Mechanical Systems and*  
679 *Signal Processing*, **25**(4), pp.1087-1102.

680 Cai, L., Hu, D., Zhang, C. Yu S, Xie J. (2022). Tool vibration feature extraction method based  
681 on SSA-VMD and SVM. *Arabian Journal for Science and Engineering*, pp.1-11.

682 Carrera, F., Dentamaro, V., Galantucci, S., Iannaccone, A., Impedovo, D. and Pirlo, G. (2022).  
683 Combining unsupervised approaches for near real-time network traffic anomaly  
684 Detection. *Applied. Science*. **12**(3), pp. 1759.

685 Cha, Y. J., and Wang, Z. (2018). Unsupervised novelty detection–based structural damage  
686 localization using a density peaks-based fast clustering algorithm. *Structural Health*  
687 *Monitoring*, **17**(2), pp.313-324.

- 688 Chandola, V., Banerjee, A., and Kumar, V. (2009). Anomaly detection: A survey. *ACM*  
689 *Computing Surveys (CSUR)*, **41**(3), pp.1-58.
- 690 Chang, C. M., Chou, J. Y., Tan, P., and Wang, L. (2017). A sensor fault detection strategy for  
691 structural health monitoring systems. *Smart Structures and Systems*, **20**(1), pp.43-52.
- 692 Chen, W., Kong, F., Mei, F., Yuan, G., and Li, B. (2017). *A novel unsupervised anomaly*  
693 *detection approach for intrusion detection system*. In 2017 IEEE 3rd International  
694 Conference on Big Data Security on the Cloud, IEEE International Conference on High  
695 Performance and Smart Computing, and IEEE International Conference on Intelligent  
696 Data and Security, 26<sup>th</sup>-28<sup>th</sup> May, Beijing, China, pp. 69-73.
- 697 De Luca, C. J., Gilmore, L. D., Kuznetsov, M., and Roy, S. H. (2010). Filtering the surface  
698 EMG signal: Movement artifact and baseline noise contamination. *Journal of*  
699 *Biomechanics*, **43**(8), pp.1573-1579.
- 700 Domingues, R., Filippone, M., Michiardi, P., and Zouaoui, J. (2018). A comparative evaluation  
701 of outlier detection algorithms: Experiments and analyses. *Pattern Recognition*, **74**,  
702 pp.406-421.
- 703 Dragomiretskiy, K., and Zosso, D. (2013). Variational mode decomposition. *IEEE Transactions*  
704 *on Signal Processing*, **62**(3), pp.531-544.
- 705 Dubé, L., and Paré, G. (2003). Rigor in information systems positivist case research: Current  
706 practices, trends, and recommendations. *MIS Quarterly*, **27**(4), pp.597–636.
- 707 Fang, W., Love, P. E., Luo, H., & Ding, L. (2020). Computer vision for behaviour-based safety  
708 in construction: A review and future directions. *Advanced Engineering Informatics*, **43**,  
709 100980.
- 710 Fang, W., Love, P. E., Ding, L., Xu, S., Kong, T., & Li, H. (2021). Computer Vision and Deep  
711 Learning to Manage Safety in Construction: Matching Images of Unsafe Behavior and  
712 Semantic Rules. *IEEE Transactions on Engineering Management (In Press)*.
- 713 Fu, Y., Peng, C., Gomez, F., Narazaki, Y., and Spencer Jr, B. F. (2019). Sensor fault  
714 management techniques for wireless smart sensor networks in structural health monitoring.  
715 *Structural Control and Health Monitoring*, **26**(7), e2362.
- 716 G. Han, J. Tu, L. Liu, M. Martínez-García and C. Choi. (2022). An Intelligent Signal Processing  
717 Data Denoising Method for Control Systems Protection in the Industrial Internet of Things.

- 718 *IEEE Transactions on Industrial Informatics*, 18(4), pp. 2684-2692,
- 719 Gao, K., Chen, Z. D., Weng, S., Zhu, H. P., and Wu, L. Y. (2022). Detection of multi-type data  
720 anomaly for structural health monitoring using pattern recognition neural network. *Smart  
721 Structures and Systems*, **29**(1), pp.129-140.
- 722 Ghasemaghaei, M., and Celic, G. (2020). Assessing the impact of big data on firm innovation  
723 performance: Big data is not always better data. *Journal of Business Research*, **108**,  
724 pp.147–162
- 725 Guo, S.Y., Ding, L.Y., Luo, H., and Jiang, X.Y. (2016). A Big-Data-based platform of workers’  
726 behavior: Observations from the field. *Accident Analysis and Prevention*, **93**, pp.299-  
727 309.
- 728 Han, L., Zhang, R., Wang, X., Bao, A., and Jing, H. (2019). Multi-step wind power forecast  
729 based on VMD-LSTM. *IET Renewable Power Generation*, **13**(10), 1690-1700.
- 730 Hariri, S., Kind, M. C., and Brunner, R. J. (2021). Extended isolation forest. *IEEE Transactions  
731 on Knowledge and Data Engineering*. **33**(4), pp.1479-1489.
- 732 Hasan, M. A. M., Nasser, M., Pal, B., and Ahmad, S. (2014). Support vector machine and  
733 random forest modeling for intrusion detection system (IDS). *Journal of Intelligent  
734 Learning Systems and Applications*, **6**(1), Article ID:42869
- 735 He, Q., Wang, X., and Zhou, Q. (2014). Vibration sensor data denoising using a time-frequency  
736 manifold for machinery fault diagnosis. *Sensors*, **14**(1), pp.382-402.
- 737 Hill, D. J., and Minsker, B. S. (2010). Anomaly detection in streaming environmental sensor  
738 data: A data-driven modeling approach. *Environmental Modelling and Software*, **25**(9),  
739 pp.1014-1022.
- 740 Hou, S., and Guo, W. (2020). Optimal denoising and feature extraction methods using modified  
741 CEEMD combined with duffing system and their applications in fault line selection of  
742 non-solid-earthed network. *Symmetry*. **12**(4), pp.536.
- 743 Hu, H., Zhang, L., Yan, H., Bai, Y., and Wang, P. (2019). De-noising and baseline drift removal  
744 method of MEMS hydrophone signal based on VMD and wavelet threshold processing.  
745 *IEEE Access*, **7**, pp.59913-59922.
- 746 Huang, H. B., Yi, T. H., and Li, H. N. (2020). Anomaly identification of structural health  
747 monitoring data using dynamic independent component analysis. *ASCE Journal of*

- 748        *Computing in Civil Engineering*, **34**(5), 04020025.
- 749 Huang, H. B., Yi, T. H., and Li, H. N. (2017). Sensor fault diagnosis for structural health  
750        monitoring based on statistical hypothesis test and missing variable approach. *Journal of*  
751        *Aerospace Engineering*, **30**(2), B4015003.
- 752 Huang N E , Shen Z , Long S R , Wu MC, SHIH HH, Zheng Q, Yen NC, Tung CC, Liu HH.  
753        (1998). The empirical mode decomposition and the Hilbert spectrum for non-linear and  
754        non-stationary time series analysis. *Proceedings Mathematical Physical & Engineering*  
755        *Sciences*, 454(1971), 903-995.
- 756 Inoue, J., Yamagata, Y., Chen, Y., Poskitt, C. M., and Sun, J. (2017). Anomaly detection for a  
757        water treatment system using unsupervised machine learning. 2017 *IEEE International*  
758        *Conference on Data Mining Workshops (ICDMW)*, 18<sup>th</sup>-21<sup>st</sup> November, New Orleans,  
759        USA, pp. 1058-1065
- 760 Liu, F. T., Ting, K. M., and Zhou, Z. H. (2012). Isolation-based anomaly detection. *ACM*  
761        *Transactions on Knowledge Discovery from Data (TKDD)*, **6**(1), pp.1-39.
- 762 Li H, Xu Y, An D, Zhang L, Li S, Shi H. (2020). Application of a flat variational modal  
763        decomposition algorithm in fault diagnosis of rolling bearings. *Journal of Low Frequency*  
764        *Noise, Vibration and Active Control*, 39(2):335-351.
- 765 Li, S., Zhang, K., Duan, P., and Kang, X. (2019). Hyperspectral anomaly detection with kernel  
766        isolation forest. *IEEE Transactions on Geoscience and Remote Sensing*, **58**(1), pp.319-  
767        329.
- 768 Li, T., Zhao, Z., Sun, C., Yan, R., and Chen, X. (2021). Hierarchical attention graph  
769        convolutional network for fuse multi-sensor signals for remaining useful life prediction.  
770        *Reliability Engineering and System Safety*, **215**, 107878
- 771 Liu, F. T., Ting, K. M., and Zhou, Z. H. (2008). Isolation forest. Proceedings of the 8<sup>th</sup> IEEE  
772        International Conference on Data Mining, 15<sup>th</sup>-19<sup>th</sup> December, Pisa, Italy pp. 413-422.
- 773 Liu, D., Wang, S., and Cui, X. (2022). An artificial neural network supported Wiener process-  
774        based reliability estimation method considering individual difference and measurement  
775        error. *Reliability Engineering and System Safety*, **218**, 108162
- 776 Liu, W., Li, A., Fang, W., Love, P. E., Hartmann, T., & Luo, H. (2023). A hybrid data-driven  
777        model for geotechnical reliability analysis. *Reliability Engineering & System Safety*, **231**,

- 778 108985.
- 779 Kromanis, R., and Kripakaran, P. (2013). Support vector regression for anomaly detection from  
780 measurement histories. *Advanced Engineering Informatics*, **27**(4), pp.486-495.
- 781 Maier, J., Naber, A., and Ortiz-Catalan, M. (2017). Improved prosthetic control based on  
782 myoelectric pattern recognition via wavelet-based denoising. *IEEE Transactions on*  
783 *Neural Systems and Rehabilitation Engineering*, **26**(2), pp.506-514.
- 784 Matthews, J., Love, P. E.D, Porter, S. R., and Fang, W. (2022). Smart data and business  
785 analytics: A theoretical framework for managing rework risks in mega-projects.  
786 *International Journal of Information Management*, **65**, 102495.
- 787 Mu, H. Q., and Yuen, K. V. (2015). Novel outlier-resistant extended Kalman filter for robust  
788 online structural identification. *ASCE Journal of Engineering Mechanics*, **141**(1),  
789 04014100.
- 790 Nessa, A., Adhikari, B., Hussain, F., and Fernando, X. N. (2020). A survey of machine learning  
791 for indoor positioning. *IEEE Access*, **8**, pp.214945-214965.
- 792 Ngo, J., Hwang, B.-G., and Zhang, C. (2020). Factor-based big data and predictive analytics  
793 capability assessment tool for the construction industry. *Automation in Construction*, **110**,  
794 103042
- 795 Nguyen, L. H., and Goulet, J. A. (2019). Real-time anomaly detection with Bayesian dynamic  
796 linear models. *Structural Control and Health Monitoring*, **26**(9), e2404.
- 797 Otoum, S., Kantarci, B., and Mouftah, H. (2018). Adaptively supervised and intrusion-aware  
798 data aggregation for wireless sensor clusters in critical infrastructures. *IEEE International*  
799 *Conference on Communications (ICC)*, 20<sup>th</sup>-24<sup>th</sup> May Kansas City, USA, pp.1-6
- 800 Ortolan, R. L., Mori, R. N., Pereira, R. R., Cabral, C. M., Pereira, J. C., and Cliquet, A. (2003).  
801 Evaluation of adaptive/nonadaptive filtering and wavelet transform techniques for noise  
802 reduction in EMG mobile acquisition equipment. *IEEE Transactions on Neural Systems*  
803 *and Rehabilitation Engineering*, **11**(1), pp.60-69.
- 804 Prateek G V., Ju Y E., and Nehorai A. (2021). Sparsity-Assisted Signal Denoising and Pattern  
805 Recognition in Time-Series Data. *Circuits Systems and Signal Processing*, **2021**(9), pp.1-  
806 50.



- 807 Rabatel, J., Bringay, S., and Poncelet, P. (2011). Anomaly detection in monitoring sensor data  
808 for preventative maintenance. *Expert Systems with Applications*, **38**(6), pp.7003-7015
- 809 Rai, H. M., and Chatterjee, K. (2019). Hybrid adaptive algorithm based on wavelet transform  
810 and independent component analysis for denoising of MRI images. *Measurement*, **144**,  
811 pp.72-82.
- 812 Samuel T. Thurman and James R. Fienup. (2008). Wiener filtering of aliased imagery.  
813 *International Society for Optics and Photonics*, **7076**, pp.70760J.
- 814 Selesnick I W., Graber H L., Pfeil D S., and Barbour R. L. (2014). Simultaneous Low-Pass  
815 Filtering and Total Variation Denoising. *IEEE Transactions on Signal Processing*, **62**(5),  
816 pp. 1109-1124.
- 817 Seites-Rundlett, W., Bashar, M., Torres-Machi, C., and Corotis, R.B. (2022). Combined  
818 evidence model to enhance pavement condition prediction from highly uncertain sensor  
819 data. *Reliability Engineering and System Safety*, **217**, 108031
- 820 Singh, P., Shahnawazuddin, S., and Pradhan, G. (2018). An efficient ECG denoising technique  
821 based on non-local means estimation and modified empirical mode decomposition.  
822 *Circuits, Systems, and Signal Processing*, **37**(10), pp.4527-4547.
- 823 Thottan, M., and Ji, C. (2003). Anomaly detection in IP networks. *IEEE Transactions on Signal*  
824 *Processing*, **51**(8), pp. 2191-2204.
- 825 Thudumu, S., Branch, P., Jin, J., and Singh, J. J. (2020). A comprehensive survey of anomaly  
826 detection techniques for high dimensional big data. *Journal of Big Data*, **7**(1), pp.1-30.
- 827 Urciuolo, M. (2008). Restriction of the Fourier transform. *Revista de la Unión Matemática*  
828 *Argentina*, **49**(2), pp. 39-44.
- 829 Wang, Z., He, G., Du, W., Zhou, J., and Kou, Y. (2019). Application of parameter optimized  
830 variational mode decomposition method in fault diagnosis of gearbox. *IEEE Access*, **99**(7),  
831 pp. 44871-44882.
- 832 Xia, S.; Yang, J.; Cai, W.; Zhang, C.; Hua, L., and Zhou, Z. (2021). Adaptive Complex  
833 Variational Mode Decomposition for Micro-Motion Signal Processing Applications.  
834 *Sensors*, **21**(5), pp. 1637.
- 835 Xiao, F., Yang, D., Guo, X., and Wang, Y. (2019). VMD-based denoising methods for surface  
836 electromyography signals. *Journal of Neural Engineering*, **16**(5), 056017.

- 837 Xu, Z., Li, C., and Yang, Y. (2020). Fault diagnosis of rolling bearing of wind turbines based  
838 on the variational mode decomposition and deep convolutional neural networks. *Applied*  
839 *Soft Computing*, **95**, 106515.
- 840 Yi, T. H., Li, H. N., and Gu, M. (2013). Recent research and applications of GPS-based  
841 monitoring technology for high-rise structures. *Structural Control and Health Monitoring*,  
842 **20**(5), pp.649-670.
- 843 Yu, L, Lin L. (2008). Improvement on empirical mode decomposition based on correlation  
844 coefficient. *Computer and Digital Engineering*, **36**(12), pp.28-29.
- 845 Yue, G. D., Cui, X. S., Zou, Y. Y., Bai, X. T., Wu, Y. H., and Shi, H. T. (2019). A Bayesian  
846 wavelet packet denoising criterion for mechanical signal with non-Gaussian characteristic.  
847 *Measurement*, **138**, pp.702-712.
- 848 Zhou, C., Kong, T., Jiang, S., Chen, S., Zhou, Y., and Ding, L. (2020). Quantifying the evolution  
849 of settlement risk for surrounding environments in underground construction via complex  
850 network analysis. *Tunnelling and Underground Space Technology*, **103**, 103490.
- 851 Zhou, Y., Li, S., Zhou, C., and Luo, H. (2019a). Intelligent approach based on random forest  
852 for safety risk prediction of deep foundation pit in subway stations. *ASCE Journal of*  
853 *Computing in Civil Engineering*, **33**(1), 05018004.
- 854 Zhou, Y., Li, C., Ding, L., Sekula, P., Love, P. E.D, and Zhou, C. (2019b). Combining  
855 association rules mining with complex networks to monitor coupled risks. *Reliability*  
856 *Engineering and System Safety*, **186**, 194-208.
- 857 Zuo, R., and Xiong, Y. (2018). Big data analytics of identifying geochemical anomalies  
858 supported by machine learning methods. *Natural Resources Research*, **27**(1), pp.5-13.  
859