

HFN:Heterogeneous Feature Network for Multivariate Time Series Anomaly Detection

Jun Zhan^{a,b,e}, Chengkun Wu^{b,c,*}, Canqun Yang^d, Qiucheng Miao^b,
Xiandong Ma^{f,**}

^a*School of Intelligent Manufacturing, Hunan First Normal University, Changsha, 410205, China*

^b*College of Computer Science, National University of Defense Technology, Changsha, 410073, China*

^c*State Key Laboratory of High Performance Computing, Changsha, 410073, China*

^d*National SuperComputing Center in Tianjin, Tianjin, 300000, China*

^e*Key Laboratory of Industrial Equipment Intelligent Perception and Maintenance Technology in College of Hunan Province, Hunan First Normal University, Changsha, 411201, China*

^f*School of Engineering, Lancaster University, LA1 4YW, Lancaster, UK*

Abstract

As the key step of anomaly detection for multivariate time-series (MTS) data, learning the relations among different variables has been explored by many approaches. However, most of the existing approaches do not consider the heterogeneity between variables, that is, different types of variables (continuous numerical variables, discrete categorical variables or hybrid variables) may have different and distinctive edge distributions. In this paper, we propose a novel semi-supervised anomaly detection framework based on a heterogeneous feature network (HFN) for MTS. Specifically, we first combine the embedding similarity subgraph generated by sensor embedding and the feature value similarity subgraph generated by sensor values to construct a time-series heterogeneous graph, which fully utilizes the rich heterogeneous mutual information among variables. Then, a prediction model containing nodes and channel attentions is jointly optimized to obtain better time-series

*Corresponding author

**Corresponding author

Email addresses: `chengkun_wu@nudt.edu.cn` (Chengkun Wu),
`xiandong.ma@lancaster.ac.uk` (Xiandong Ma)

representations. This approach fuses the state-of-the-art technologies of heterogeneous graph structure learning (HGSL) and representation learning. **Experiments conducted on** four sensor datasets from real-world applications demonstrate that our approach detects the anomalies more accurately than those baseline approaches, thus providing a basis for the rapid positioning of anomalies.

Keywords: Heterogeneous neural network; Anomaly detection; Multi-sensor data; Multivariate time series; Deep learning

1. Introduction

As information technology develops, an increasing number of industrial systems **are** exposed to the internet, posing serious risks to their ability to operate securely [1]. Continuous monitoring the operation data of the system and precisely and effectively identifying potential attacks or the evolution of the equipment condition by **using this data** is an effective technique to handle these challenges [2]. For instance, an operation and maintenance personnel in a large power plant can quickly identify abnormal sensor behavior using the precise intrusion detection systems, which are developed by massive amounts of data collected by the supervisory control and data acquisition (SCADA) system [3], providing them a possibility to prevent potential system failures before irreversible damage. However, these monitoring data always have complicated structures, high dimensionality, and hard labeling, making manual tasks difficult to handle. Therefore, it is vitally necessary to investigate the semi-supervised or unsupervised time-series anomaly detection approach by utilizing a sizable amount of complicated unlabeled data.

Recently, deep learning technique has been applied successfully in various anomaly detection problems [4, 5, 6]. For high-dimensional MTS analysis, the temporal relations between different timestamps are considered first [7]. Because of their capability of capturing long-term dependency relations, recurrent neural network [8] and temporal convolutional network [9] were demonstrated to achieve better results on the time-series tasks involving single or multiple variables [10]. However, various sensors could be mutually coupled. The capacity of these approaches to detect abnormalities may be constrained by their modeling of solely temporal variables. Therefore, it is crucial to take into account both the temporal features of different timestamps and potential correlations among these variables [11, 12]. Combining

28 the sequential network and the convolution neural network (CNN) is an ef-
29 fective way to achieve this. Cross-correlation among high-dimensional data
30 can be extracted by using the local perception capacity of the convolutional
31 kernel [13]. However, CNN is primarily used to handle Euclid-space data,
32 such as image [14]. There exist some limitations on the MTS with different
33 attributes. In such cases, the graph neural network (GNN) has been success-
34 fully applied into the modelling of MTS due to its good structure modelling
35 capability between complex data; the most advanced results are achieved in
36 [11, 15].

37 With regards to the latent feature modeling of time-series data, the vari-
38 able attributes from the data are generally seen as homogeneous in the most
39 existing papers; that is, the data types are treated without distinction, such
40 as use of the variational autoencoders [16] and generative adversarial net-
41 works [17]. These methods model complex distribution from large-scale high-
42 dimensional datasets. After the training is finished by using the dataset from
43 normal conditions, the similar generative data are viewed as normality, while
44 the dissimilar data are viewed as anomalies. However, there are still fewer
45 works considering the heterogeneity of time-series data, although this kind
46 of data are abundant in practical situations. For instance, in a large-scale
47 water processing system [17], the information, such as flow, pressure and liq-
48 uid level collected by the sensors in the intermediate process, is collected as
49 the numeric continuous values. However, the signals, such as valve state and
50 location collected by the sensors of the actuator, are generally the categor-
51 ical discrete values. Inputting the mixed type of heterogeneous data into a
52 deep learning network may cause the useful information to be ignored and
53 therefore satisfied results cannot be obtained. The fundamental reason is
54 that there are totally different edge distributions between the variables with
55 different types [18, 19].

56 To overcome the limitation of deep learning model in such circumstances,
57 we propose a heterogeneous feature learning network for MTS, and study its
58 abnormal detection capability with the extensive real-world datasets. The
59 overall framework can be divided into three stages: 1) Heterogeneous graph
60 structure learning (HGSL) stage for MTS. We fuse the sensor embedding
61 vector similarity matrix and the feature value similarity matrix of different
62 variable categories to model the heterogeneous structural information. More-
63 over, we propose a category-based fixed-length approach to replace the widely
64 used meta-path [20] for extracting heterogeneous relation subgraphs. 2) Het-
65 erogeneous representation learning stage for MTS. We embed different kinds

66 of variables into vectors for fusion. Distinct from the previous heterogeneous
67 graph attention network [21], we further expand the channel attention on the
68 basis of node attention and semantic attention, so as to achieve a joint opti-
69 mization training of node embedding representation with different types. 3)
70 Abnormal detection and location stage. By analyzing the deviation between
71 the predicted and real values, we calculate a condition score for each sensor,
72 where the largest condition score is considered as the maximum abnormal
73 probability.

74 The major contributions of the paper are summarized as follows:

- 75 • We propose a novel HGSL approach for MTS, which learns heteroge-
76 neous graph structure information between sensor-embedding vectors
77 and category-based feature value vectors simultaneously.
- 78 • We propose a heterogeneous feature network (HFN) and apply it to
79 MTS anomaly detection. Our approach successfully learned the dy-
80 namic dependency among different variables and timestamps by uti-
81 lizing two single-level attention mechanisms, namely attention-based
82 node embedding and channel aggregation.
- 83 • The extensive experiments indicate that HFN can detect the anomalies
84 from real-world MTS datasets and is proved to outperform the most
85 existing methods. Besides, we analyze the condition scores of MTS,
86 demonstrating that the proposed method has the advantage of locating
87 the anomalies.

88 The rest of this paper is structured as follows. Section 1 describes the
89 related work of MTS anomaly detection. Section 2 presents the structure
90 and working principle of HFN-based MTS anomaly detection framework in
91 detail. Section 3 show the performance of proposed method on three real-
92 world MTS datasets. Finally, the conclusion and future improvements are
93 given in Section 4.

94 2. Related work

95 MTS anomaly detection has extensive application prospects in the fields
96 of industry, financial business, and the Internet of Things. As the key research
97 problem in this paper, we firstly review the related work for MTS anomaly
98 detection, which can generally be categorized as unsupervised, supervised,

99 and semi-supervised. We focus on studying data heterogeneity modeling of
100 MTS, especially heterogeneity representation learning from time-series data,
101 graph structure learning, and heterogeneous graph neural network.

102 *2.1. MTS anomaly detection*

103 MTS anomaly detection is typically regarded as an unsupervised learning
104 problem [22], and algorithms based on clustering [23], such as fuzzy c-means
105 [24], or spatiotemporal clustering [25], are frequently used. By grouping
106 time-series data into various clusters, these techniques can identify anoma-
107 lies by calculating the similarity or distance between the observed value [26]
108 and the cluster center [27]. However, unsupervised detection methods usu-
109 ally focus more on static data model development. In contrast, a supervised
110 abnormal detection algorithm has a higher detection accuracy. Under the
111 circumstance of high-quality labeling, the indicator accuracy can be approx-
112 imate to 100% [28]. However, the supervised detection requires that the
113 training set contains correctly both labeled positive and negative samples,
114 which is often not easy. [29]. Fortunately, in the actual cases, we have a
115 chance to obtain a large quantity of data under the normal conditions [17],
116 making the semi-supervised abnormal detection attract wide attentions [30].
117 In the latest work, Miryam et al. [31] proposes the methods to show the
118 great advantages and extensive application **prospects** of the semi-supervised
119 algorithm in MTS abnormal detection.

120 *2.2. Modeling for heterogeneous data*

121 The data heterogeneity has been widely concerned such as in the music
122 recommendation system [32], academic network [33] and social platform [34].
123 **The** heterogeneous learning method usually focuses on capturing and inte-
124 grating couplings with multiple variable types **at** the same or different levels.
125 To learn the embedding representation of heterogeneous data, the matrix de-
126 composition method is traditionally adopted [35, 36]. However, it is usually
127 very expensive and low-efficient in terms of the computation cost of decom-
128 posing a large-scale matrix [37]. Moreover, the discretization of continuous
129 features [38] or continuous data [39] are also a typical method; however this
130 transformation may ignore the correlation between variables. To solve these
131 challenges, heterogeneous graph embedding or heterogeneous graph repre-
132 sentation learning [40] has been widely studied. Its main goal is to map the
133 input data into low-dimensional space while simultaneously preserving the
134 heterogeneous structure and semantic characteristics of the data [41]. For

135 instance, for the tasks of text classification, Wang et al. [21] proposed a het-
136 erogeneous graph attention network (HAN), which aggregates the features
137 of meta-path based neighbors through a hierarchical manner to generate the
138 embedding representation of nodes. Fu et al. [42] proposed a meta-path ag-
139 gregated graph neural network (MAGNN) by designing multiple candidate
140 encoder functions to extract heterogeneous information from the meta-path.
141 Wang et al. [43] combined the heterogeneous graph neural network with com-
142 parison learning, and proposed a self-supervised heterogeneous graph neural
143 network from both heterogeneous network and meta-path for learning node
144 embedding representation. In the social or citation network, in order to cap-
145 ture the dynamic performances of heterogeneous redgraphs, Hu et al. [44]
146 proposed a heterogeneous graph transformer (HGT) by introducing a relative
147 temporal encoding technique for solving the problem where the dynamic re-
148 sult dependence is difficult to capture. Yang et al. [45] proposed a dynamic
149 heterogeneous graph (DyHAN) utilizing structural heterogeneity and time
150 revolution to learn node embedding. **In addition, contrastive self-supervised**
151 **learning has been widely employed to address the limitation of sparse la-**
152 **bel information in the potential ability of heterogeneous graph neural net-**
153 **work models for representation learning. For instance, the HGCL method**
154 **proposed by Chen et al.[46] effectively utilizes the structural information**
155 **of heterogeneous graphs to capture relationships between different types of**
156 **nodes. Zhu et al.[47] combine heterogeneous graph contrastive learning with**
157 **a structure-enhancement method, proposing the STENCIL method. This**
158 **approach introduces a novel multi-view contrastive aggregation objective to**
159 **adaptively distill information from each view. Furthermore, the method en-**
160 **riches the local structural patterns of the underlying heterogeneous graph to**
161 **better explore true and challenging negative examples in graph contrastive**
162 **learning. Although the above methods have achieved significant success in**
163 **their respective application domains, leveraging the structure of heteroge-**
164 **neous graphs to enhance data representation capabilities and demonstrating**
165 **outstanding performance through representation learning methods, their ap-**
166 **plicability may be subject to domain specificity and might not necessarily be**
167 **suitable for other areas such as multivariate time series anomaly detection.**

168 2.3. Graph structure learning

169 MTS usually **exists** in the form of tabular data [48], lacking of predefined
170 graph structure required for graph neural network [15], which constitutes the
171 challenge for the modelling [49]. Hence, it is extremely vital to learn the links

172 between edges and refine the graph from the existing time-series data [50].
173 The existing methods can mainly be divided into three categories: metric-
174 based approaches usually implemented by using kernel function [51, 52], co-
175 sine similarity [53, 54] or inner product [55] to calculate the similarity between
176 nodes as edge weights. Neural networks-based approaches have generally uti-
177 lized a complex deep neural network to model the edge weights of the given
178 node features and representations. For instance, Luo et al. [56] proposed
179 a multilayer perception-based graph structure optimization approach, where
180 the edge number of a sparse graph is punished through parameterized net-
181 work for pruning the edges that are unrelated to the tasks. Zhao et al. [11]
182 proposed a graph structure learning approach with redan attention coeffi-
183 cient, while Sun et al. [57] utilized a dot-product self-attention to model the
184 dynamic connection relations between the nodes. Direct learning approaches,
185 regarding adjacent matrix as a learnable parameter, make associative learn-
186 ing together with the follow-up tasks for optimization. For instance, Gao et
187 al. [58] proposed the graph learning neural networks (GLNNs) utilizing spec-
188 tral graph theory for graph learning. However, these approaches mostly aim
189 at learning isomorphic graph structure. To enable capture the heterogene-
190 ity between the data efficiently, Zhao et al. [41] proposed a heterogeneous
191 graph learning approach utilizing the fusion of feature similarity sub-graph,
192 feature propagation graph and semantic graph, which successfully learns an
193 appropriate graph structure for a heterogeneous graph neural network.

194 3. Proposed Frameworks

195 3.1. Problem statement

196 Generally, we define heterogeneous MTS dataset as a time-series dataset
197 with L variables, N different types of sensors, and T length, which is ex-
198 pressed as $X = \{\mathbf{x}_{1:T}^N\}$, where $N \in \{type^1, \dots, type^n\}$ denotes the set of data
199 types. Note that the variable number contained in the specified categories
200 may be larger than 1. For instance, for arbitrary data type $type^n$, all time se-
201 ries at the moment t can be denoted as $\mathbf{x}_t^{type^n} \in \{x_t^{type^n_i}, for i \in \{0, \dots, d\}\}$,
202 where d represents the number of time-series sequence in this category. In
203 this paper, we adopt the sliding window-based model training approach. At
204 the moment t , we sample a continuous subsequence with the length of ω as
205 the model input, denoted as $S^N(t) = [\mathbf{x}_{t-\omega+1}^N, \dots, \mathbf{x}_t^N]$. For the abnormal
206 detection task, our target is to predict the value of all sensors \mathbf{x}_{t+1}^N at the

207 moment $t + 1$ by utilizing the input subsequence $S^N(t)$, and obtain the pre-
 208 dicted value $\widehat{\mathbf{x}}_{t+1}^N$. The mean square error (MSE) between the predicted
 209 value and practical value is used as loss to optimize the model. According
 210 to the usual semi-supervised abnormal detection methods, in the training
 211 stage, only the data collected from normal conditions are chosen. However,
 212 in the testing stage, the deviation between the predicted value and practical
 213 value is further used for calculating the condition scores of the data, while
 214 the scores of the corresponding data over the threshold **are** judged as the
 215 anomalies, otherwise normal.

216 Specially, we divide time-series data into three data types, that is $N \in$
 217 $\{C, CD, D\}$:

- 218 • Continuous numerical variables C , where the value of data are taken
 219 from continuous real number, such as $x_t^{C_i} \in \mathbb{R}$.
- 220 • Discrete categorical variables D , where the value of data are taken from
 221 a limited set of values, such as $x_t^{D_i} \in \{0, 1, 2\}$.
- 222 • Hybrid variables S_t^{CD} which contain both numerical and categorical
 223 variables where the values of the element are taken from the above two
 224 categories.

225 We construct a heterogeneous dynamic graph to model the above MTS.
 226 Different time-series variables are viewed as the node in the graph, while their
 227 connection relation is seen as the edge. This dynamic graph can be denoted as
 228 $\mathbb{G}_{S^N(t)} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} and \mathcal{E} represent node and edge set respectively. We
 229 respectively extract categorical feature subgraph $\mathbb{G}_{S^D(t)}$, numerical feature
 230 subgraph $\mathbb{G}_{S^C(t)}$, and categorical and numerical mixed subgraph $\mathbb{G}_{S^{CD}(t)}$ for
 231 learning heterogeneous information. For the arbitrary subgraph, its adjacent
 232 matrix is $A_N \in \mathbb{R}^{|\mathcal{V}_N| \times |\mathcal{V}_N|}$, where \mathcal{V}_N represents the node set with the specific
 233 type. If there exist connection relations between two arbitrary nodes in the
 234 subgraph, the corresponding element of adjacent matrix is 1. Noted that the
 235 final node embedding integrates the node embedding representations of three
 236 different subgraphs.

237 3.2. Model Architecture

238 Our HFN-based approach aims at learning the complex correlation be-
 239 tween different types of time-series data carried by the defined dynamic graph

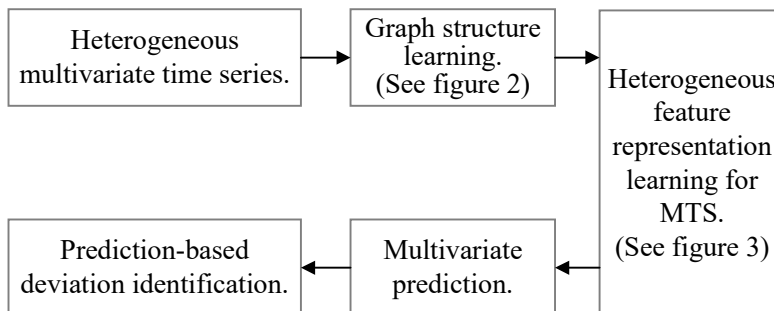


Figure 1: **Architecture of HFN-based MTS anomaly detection framework.**

240 above. For each node, the potential temporal correlation is allowed to be
 241 considered with a sliding window along the dataset.

242 Figure 1 shows the proposed HFN-based semi-supervised abnormal de-
 243 tection framework architecture. It can be seen that for a given MTS, we
 244 firstly learn a heterogeneous dynamic graph representing the structural in-
 245 formation between different variables (as shown in Figure 2), decomposing
 246 the time-series data into different graph structures. On this basis, the cat-
 247 egorical feature subgraph, the continuous numerical feature subgraph and
 248 the hybrid subgraph are extracted and then inputted into the HFN network
 249 based on graph attention function to learn the potential embedding repre-
 250 sentations of each sensor (as shown in Figure 3). Then we predict the future
 251 values of each sensor based on these embedding representations. Finally, the
 252 deviation between the predicted and practical values is used for measuring
 253 and locating the anomalies.

254 3.3. Graph structure learning pipeline

255 To learn the complex heterogeneous potential features between different
 256 types of sensors, a key process is how to map the variable correlation from
 257 MTS into the adjacent matrix of the graph. In the previous studies, all
 258 assumed that the constructed graph is the static isomorphic graph, thus re-
 259 sulting in the loss of some key information. For instance, the significance
 260 of variables exists great difference at the operating condition of full-load
 261 and partial-load of generating equipment [59]. Hence, as shown in Figure 2,
 262 we learn the potential heterogeneous graph structure of MTS from the per-
 263 spectives of global semantic correlation and local feature correlation. For
 264 the global semantic correlation, we introduce a learnable embedding vector

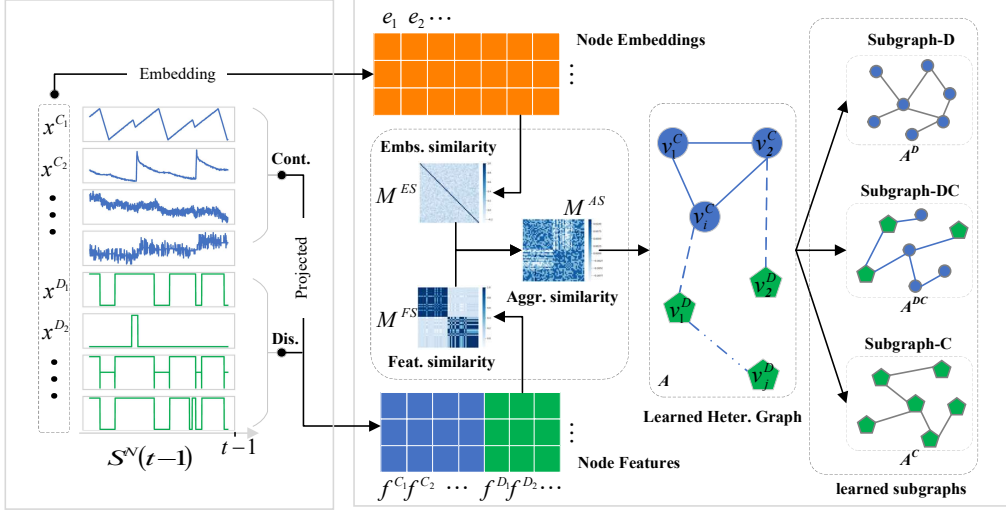


Figure 2: Structure learning of MTS heterogeneous dynamic graph.

265 for each variable, and denote it as $e_i \in \mathbb{R}^{1 \times \omega'}$. For $i \in \{0, \dots, L\}$, where
 266 ω' represents the dimension of embedding vector. This vector can be learned
 267 together with subsequent prediction network parameters. For the local fea-
 268 ture correlation, we calculate the potential structural information based on
 269 the feature values of the variables. We adopt a special mapping network
 270 to project different types of input feature vector $S^N(t)$ into a public space.
 271 Taking data type C as an example, the projected feature of the arbitrary
 272 variable x^{C_i} is denoted as $f^{C_i} \in \mathbb{R}^{1 \times \omega'}$:

$$f^{C_i} = SELU(x^{C_i} \bullet \mathbf{W}^C + \mathbf{b}^C) \quad (1)$$

273 where $x^{C_i} \in \mathbb{R}^{1 \times \omega}$ is the subset of all continuous numerical variables. ω
 274 is the time length of input feature vector. $\mathbf{W}^C \in \mathbb{R}^{\omega \times \omega'}$ is learnable weight ma-
 275 trix, and $\mathbf{b}^C \in \mathbb{R}^{1 \times \omega'}$ is biasing. Similarly, we can calculate and obtain the
 276 projected feature representation f^{D_i} of the discrete categorical variables.

277 3.4. Similarity Graphs

278 The main task of graph structure learning is to learn an adjacent matrix
 279 representing the mutual connection between nodes in the graph. There-
 280 fore, we propose a learning approach based on aggregating cosine similar-
 281 ity. According to the embedding for the variables and the mapping of
 282 variable feature vectors, we obtain the global semantic embedding matrix

283 $\mathbf{E} \in \{e_1, \dots, e_L\}$ and local feature vector representation matrix $\mathbf{F}^N \in \{f^{N_1}, \dots, f^{N_L}\}$.
 284 Clearly, these obtained **matrices** from different perspectives contain differ-
 285 ent information. Specifically, we **first** calculate cosine similarity between the
 286 elements in different matrices to obtain their connection information. After
 287 obtaining the node embedding (NE) similarity matrix $M^{E_s} \in \mathbb{R}^{L \times L}$ and node
 288 feature (NF) similarity matrix $M^{F_s} \in \mathbb{R}^{L \times L}$, we fuse them to obtain an aggre-
 289 gating similarity matrix, where the value represents the similarity between
 290 the arbitrary two nodes i and j and can be calculated as follows:

$$M^{E_s}[i, j] = \frac{e_i \bullet e_j}{e_i \times e_j} \quad (2)$$

$$M^{F_s}[i, j] = \frac{f^{N_i} \bullet f^{N_j}}{f^{N_i} \times f^{N_j}} \quad (3)$$

$$M^{A_s} = M^{E_s} \circ \mathbf{W}^{E_s} + M^{F_s} \circ \mathbf{W}^{F_s} \quad (4)$$

291 where \circ denotes Hadamard product between two matrixes. $\mathbf{W}^{E_s} \in \mathbb{R}^{L \times L}$
 292 and $\mathbf{W}^{F_s} \in \mathbb{R}^{L \times L}$ are learnable weight matrixes, which weigh the importance
 293 of different dimensions of the different similarity matrixes. In M^{A_s} , when
 294 the correlation coefficient is larger than a certain threshold, we consider that
 295 there exists a connected relation between nodes; otherwise, the connected
 296 relation does not exist. To obtain the optimal threshold, we define a learn-
 297 able parameter $\tau \in \mathbb{R}$ for automatic choice, and obtain the adjacent matrix of
 298 aggregating similarity graph through learning, which is denoted as:

$$A_{ij} = \begin{cases} 1 & \text{for } M^{A_s}[i, j] \geq \tau \\ 0 & \text{for } M^{A_s}[i, j] < \tau \end{cases} \quad (5)$$

299 In the heterogeneous dynamic graph, two objects can be connected through
 300 different semantic paths, which is called meta-path. However, the selec-
 301 tion of meta-path has a strong subjective meaning, which is difficult for
 302 complex MTS. Therefore, we propose a classifying-based fixed-length sam-
 303 pled method to replace meta-path for extracting heterogeneous relation **sub-**
 304 **graphs**. Specifically, we divide the aggregating similarity graph into the cor-
 305 responding classifying subgraphs, **including** discrete feature subgraph (DFS)
 306 $\mathbb{G}_{SD(t)}$, continuous feature subgraph (CFS) $\mathbb{G}_{SC(t)}$ and hybrid feature sub-
 307 graph (HFS) $\mathbb{G}_{SCD(t)}$ according to data types. We further make a random
 308 mask operation for the neighboring matrix of the subgraph and obtain the

309 final neighboring matrix with different relations. The transformed heteroge-
 310 neous graph structure is $A'=\{A^D, A^C, A^{CD}\}$. The random mask is conducive
 311 to exchange information between different similarity matrixes in the graph
 312 structure learning process, thus improving the accuracy of subsequent tasks
 313 and relieving the overfitting problem.

314 3.5. Graph representation learning for MTS

315 It can be seen from the learned heterogeneous graph structure that each
 316 **type of subgraph** contains different semantic properties. Hence, to aggregate
 317 the node information from different types, we introduce a graph attention-
 318 based node embedding network and an attention-based channel aggregating
 319 network to construct the HFN for MTS. The structure is shown in Figure 3.
 320 Specifically, the obtained three subgraphs A^D, A^C, A^{CD} learned by graph
 321 structure learning are inputted into three independent graph attention net-
 322 works, to learn the importance of different types of nodes for the neighbors
 323 in the subgraphs. Moreover, the important neighboring information is aggre-
 324 gated to generate a new node embedding. As shown in Figure 3, taking the
 325 continuous numerical variable channel as an example, for the arbitrary node
 326 v_i^C and its neighboring node v_j^C in subgraph A^C , we perform self-attention in
 327 the nodes. The attention coefficient representing their relation importance
 328 can be calculated as:

$$\xi_{ij} = att(\mathbf{W}f^{C_i}, \mathbf{W}f^{C_j}; A^C) \quad (6)$$

329 where $f^{C_i} \in \mathbb{R}^{1 \times \omega'}$ and $f^{C_j} \in \mathbb{R}^{1 \times \omega'}$ are mapped node feature vectors, $\mathbf{W} \in \mathbb{R}^{\omega'' \times \omega'}$ is
 330 shared weight matrix. ω' and ω'' are the calculated node feature vector di-
 331 mensions before and after the embedding. After obtaining the importance
 332 of subgraph-based node pairs, we normalize them via the SoftMax function
 333 and obtain weight coefficient α_{ij} :

$$\alpha_{ij} = softmax(\sigma_{ij}) = \frac{exp\left(\delta\left(\vec{a}^T [\mathbf{W}f^{C_i} || \mathbf{W}f^{C_j}]\right)\right)}{\sum_{\eta \in N_i^C} exp\left(\delta\left(\vec{a}^T [\mathbf{W}f^{C_i} || \mathbf{W}f^{C_\eta}]\right)\right)} \quad (7)$$

334 where δ is the activation function, and LeakyReLU function is usually
 335 adopted [55]. $\vec{a} \in \mathbb{R}^{2\omega''}$ is the learnable weight vector, which denotes the
 336 information concatenation of the two nodes. Finally, the output of each
 337 node can be obtained through aggregating its neighboring nodes. Multi-
 338 head attention mechanism is proven to be beneficial in the learning process

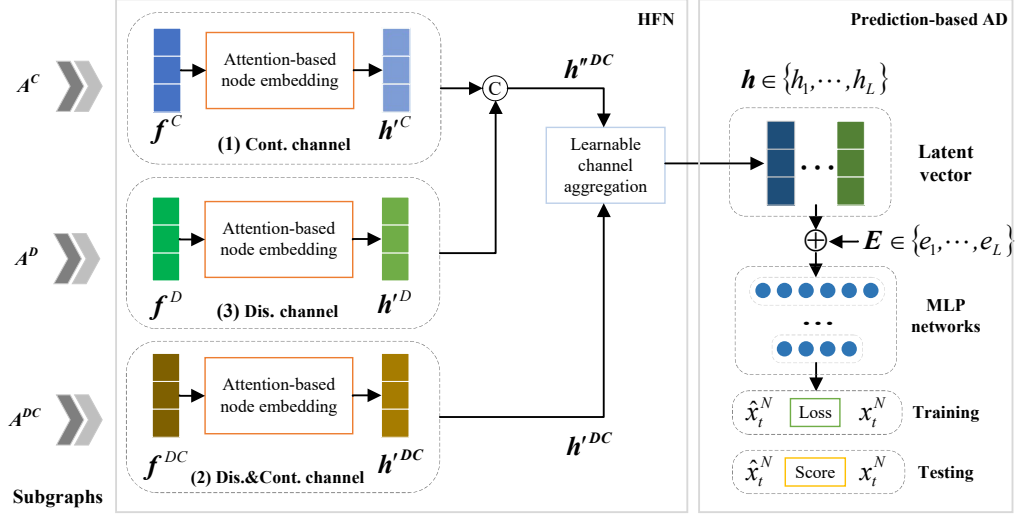


Figure 3: Heterogeneous feature network structure of MTS.

339 of stabilizing self-attention [60]. To be convenient for training, we perform an
 340 average operation to aggregate the results handled by multi-head attention.
 341 After the graph attention-based nodes embed into the network, the implicit
 342 vector can be represented as:

$$h'_i{}^C = \sigma \left(\frac{1}{H} \sum_{h=1}^H \sum_{j \in N_i^C} \alpha_{ij}^h \mathbf{W}^h f^{C_j} \right) \quad (8)$$

343 where N_i^C is the set of nodes i 's neighbors in the continuous subgraph. H denotes
 344 the number of multi-head attention mechanism head. According to the same
 345 computation method, we can obtain the node implicit vectors of discrete
 346 subgraph and mixed subgraph represented by $h'_i{}^D$ and $h'_i{}^{CD}$.

347 To address the node semantic importance of different types in the hetero-
 348 geneous graph, we put forward an attention-based multi-channel node em-
 349 bedding aggregating network. We can clearly see from Figure 3 that the node
 350 implicit vectors h'^C and h'^D singly from a continuous channel and discrete
 351 channel first are concatenated in feature dimension to obtain the global node
 352 implicit vector $h''DC$. The main purpose is to achieve the joint embedding
 353 representation learning of all nodes simultaneously. Then $h''DC$ and the node
 354 implicit vector h'^CD from the mixed channel are sent to the multi-channel
 355 node embedding aggregating network for aggregating their heterogeneous

356 information. The aggregating network automatically learns the importance
 357 degree β of the embedding vectors between different channel node implicit
 358 vectors, which can be explained as the contribution of the node correlation
 359 due to the different types of variables. The final embedding vector is com-
 360 puted as follows:

$$\mathbf{h} = \beta (\mathbf{h}'^C || \mathbf{h}'^D) + (1 - \beta) \mathbf{h}'^{CD} \quad (9)$$

361 where $\beta \in \mathbb{R}$ is a learnable parameter representing the importance degree
 362 of the embedding vectors between different channel node implicit vectors,
 363 and $||$ is a concatenation operation.

364 3.6. Prediction-based anomaly detection pipeline

365 From the above node heterogeneous feature learning network, we obtain
 366 new embedding representations of all nodes. Finally, as shown in Figure 3,
 367 we input the embedding data fused with \mathbf{h} and embedding vector \mathbf{E} into the
 368 MLP layer to have the predicted value $\hat{\mathbf{x}}_t^N$ of all sensors at the moment t :

$$\hat{\mathbf{x}}_t^N = SeLU(f(\mathbf{h} \oplus \mathbf{E})) \quad (10)$$

369 where $f(\cdot)$ is multiple layers of MLP output layer. $SeLU$ is activation
 370 function, and \oplus is addition operation.

371 At the training stage, we adopt MSE as the loss function of the model:

$$\mathcal{L}_{mse} = \frac{1}{L} \sum_i^L (\mathbf{x}_t^N - \hat{\mathbf{x}}_t^N)^2 \quad (11)$$

372 After the training is finished, we apply the network to perform real-time
 373 abnormal detection tasks. By comparing the predicted and original values of
 374 the input, we calculate the condition scores of each sample in time-series data.
 375 We define the difference between the original value and predicted value as
 376 the condition scores. To eliminate the effect of different variable dimensions,
 377 we normalize the condition scores. Finally, the condition score is computed
 378 as follows:

$$Score_i = \frac{|\mathbf{x}_t^{N_i} - \hat{\mathbf{x}}_t^{N_i}| - IQR_i}{\mu_i + 1} \quad (12)$$

379 where IQR_i denotes an interquartile range of the predicted value of the
 380 i th variable, μ_i is its median. To achieve the anomaly positioning, we take

381 the largest value of $Score_i$ as the condition score of overall record data at
382 the moment t , as denoted by $Score = \max(Score_i)$. Finally, if the $Score$ is
383 larger than the threshold, this record is judged as an anomaly. However,
384 because the threshold selection refers to complicated domain knowledge and
385 the selection methods are various depending on the applications [61], this
386 paper will not further explore the selection method for the threshold. The
387 experiment in the subsequent section will report the optimal value of each
388 evaluating metric (see details in Section 3.3).

389 3.7. Training

390 Following the application of the components introduced in the preceding
391 sections, predictions for multivariate time series can be acquired. The fun-
392 damental concept of our approach centers on maximizing the utilization of
393 diverse sensor data types within the time series, enhancing prediction accu-
394 racy, and identifying anomalies based on prediction errors. To accomplish
395 this, we collaboratively optimize a heterogeneous feature network across mul-
396 tiple channels to update the parameters of the entire network. Throughout
397 the training process, the comprehensive forward propagation procedure is
398 delineated in Algorithm 1.

Algorithm 1 HFN training procedure

Input: Heterogeneous multivariate time series training dataset $S^N(t-1) = [\mathbf{x}_{t-w}^N, \dots, \mathbf{x}_{t-1}^N]$, Batch Size \mathcal{B} , Number of Epochs \mathcal{E}

Output: Predicted values $\widehat{\mathbf{x}}_t^N$

- 1: **for** epoch=1: \mathcal{E} **do**
 - 2: Calculate the projected feature f^{C_i} and node embedding vector \mathbf{e}_i ;
 - 3: Calculate similarity matrix M^{E_s} , M^{F_s} and M^{A_s} with Eq. (2), Eq. (3) and Eq. (4);
 - 4: Calculate adjacent matrix A_N with Eq. (5);
 - 5: Extract subgraph features to obtain the node implicit vectors h_i^{C} , h_i^{D} and h_i^{CD} with Eq. (8);
 - 6: Calculate the final embedding vector \mathbf{h} with Eq. (9);
 - 7: Calculate the predicted value $\widehat{\mathbf{x}}_t^N$ with Eq. (10);
 - 8: Calculate the loss \mathcal{L}_{mse} with Eq. (11);
 - 9: Update parameters.
 - 10: **end for**
-

Table 1: Statistics of the datasets.

| Items | SWaT | WADI | WTD |
|----------------------|------------|-------------|----------|
| Time series (C/D) | 51 (25/26) | 123 (68/55) | 37(31/6) |
| Training dataset | 496800 | 784571 | 1000000 |
| Testing dataset | 449919 | 172803 | 940000 |
| Anomaly Rate (%) | 11.97% | 5.99% | 20.64% |
| Sampling Rate | 1Hz | 1Hz | 1Hz |

399 4. Experiments

400 We employ extensive experiments on two open and one private real-world
401 datasets to answer the following research questions: (1) Whether the pro-
402 posed model is more optimal than the baseline models? (2) How each com-
403 ponent of the model affects the model? (3) How the proposed approach
404 detects anomalies? (4) How the detection results locate anomalies?

405 4.1. Benchmark datasets

406 The selected three datasets contain two datasets (SWaT and WADI)
407 based on water treatment simulator testbed and a real-world dataset from a
408 large-scale wind farm (WTD). The statistical data of the datasets are given
409 in Table 1:

410 **Secure Water Treatment (SWaT) Dataset** [62]. This dataset was
411 collected from a six-stage Secure Water Treatment (SWaT) testbed. SWaT
412 represents a scaled-down version of a real-world industrial water treatment
413 plant. It took 11 days for the data collection process, which ran with nor-
414 mal operation mode during the first seven days, and constituted a training
415 dataset. During the later four days, the testbed was implemented by inter-
416 mittent network and physical attacks, which constituted the labeled testing
417 dataset. The data were collected once every second, containing 51 time-
418 series features, including 25 continuous features and 26 discrete categorical
419 features. We chose this dataset for case study, and the primary sensors or
420 actuators involved are shown in the Table 2 below.

421 **Water Distribution (WADI) Dataset** [63]. This dataset was collected
422 from a water distribution testbed (WADI). It took 16 days for the data

Table 2: Statistics of the datasets.

| No. | Name | Type | Description |
|------------|-------------|---------------------|--|
| 1 | FIT-401 | Sensor (continuous) | Flow transmitter to control the UV dechlorinator. |
| 2 | UV-401 | Actuator (discrete) | Dechlorinator to remove the chlorine from water. |
| 3 | FIT-504 | Sensor (continuous) | Flow meter, a RO re-circulation flow meter. |
| 4 | P-501 | Actuator (discrete) | Pump to pump the dechlorinated water to RO. |
| 5 | LIT-401 | Sensor (continuous) | Level transmitter to regulate the RO feed water tank level. |
| 6 | LIT-101 | Sensor (continuous) | Level transmitter to regulate the raw water tank level. |
| 7 | FIT-601 | Sensor (continuous) | Flow meter a UF backwash flow meter. |
| 8 | AIT-504 | Sensor (continuous) | RO permeate conductivity analyzer to measure the NaCl level. |
| 9 | AIT-201 | Sensor (continuous) | Conductivity analyzer to measure the NaCl level. |

423 collection process. During the last two days, the attack was launched to the
424 testbed with different intentions and time intervals, and the duration of the
425 attack lasted between 1.5 to 30 minutes to acquire the abnormal operating
426 data. The data were collected once every second, containing 123 time-series
427 features, including 68 continuous features and 55 categorical features.

428 **Wind Turbine Dataset (WTD).** This dataset was collected from a
429 large-scale wind farm [64]. It lasted 1 to 2 years for the data collection
430 process. At the training stage, there are no abnormal operating data since
431 only the time-based maintenance process was arranged for the wind turbines,
432 while at the testing stage, the abnormal operating data were detected in the
433 repairing process. All data have been labeled by the experts. The data were
434 collected once every 10 minutes, containing 37 time-series features, including
435 31 continuous features and 6 categorical features.

436 It is noteworthy that in this paper, the time scales of the time series
437 datasets are uniform, with all datasets adhering to a fixed time scale of
438 1 second. However, it is crucial to recognize that the time scale, or the
439 sampling rate of the data, can impact the identification results in time series
440 analysis. The uniformity in time scales across the datasets employed in the
441 paper ensures the effective facilitation of direct comparisons between different
442 methods.

443 4.2. Baseline models

444 We first compare the FHN model with the most advanced approaches in-
445 cluding LSTM-VAE [65], USAD [66], MAD-GAN [17], graph network based
446 MTAD-GAT [11] and GDN [54]. These approaches are extensively concerned
447 with the cross-time and cross-sequence correlation of MTS. The approaches
448 based on sequence reconstruction or prediction are used to learn the repre-
449 sentations of the whole time series. Moreover, the anomalies are judged by
450 the reconstructing or predicting errors.

451 Furthermore, we compare the proposed approach with those classic shal-
452 low anomaly detection approaches, including PCA [67], Isolation Forest (IF)
453 [68] and LightGBM [69]. These shallow detection methods are regarded as
454 the relatively direct abnormal detection methods, which usually can directly
455 locate the outlier. Moreover, to complete the anomaly detection in tempo-
456 rally related contexts has also attracted the interests of the researchers, such
457 as LSTM-NDT [1]. The idea underlying this method is to model the tem-
458 poral features of the data, predict the corresponding values, and then judge

459 whether the anomalies occur by comparing the deviation between the real
460 value and the predicted value.

461 In addition, we also conducted comparisons with the latest methods
462 based on transformer and spatiotemporal graph approaches. These include:
463 TranAD [70], an anomaly detection and diagnostic model based on deep
464 transformer networks. It employs attention-based sequence encoders for
465 rapid inference, possessing knowledge of broader temporal trends in the data;
466 FuSAGNet [71], which combines sparse autoencoder and graph neural net-
467 work. The latter predicts future time series behavior from sparse latent
468 representations learned by the former, along with graph structures learned
469 through recurrent feature embedding; MAD-SGCN [72], which effectively
470 captures the spatiotemporal correlations of input sequences using long short-
471 term memory networks (LSTMs) and spectral-based graph convolutional net-
472 works (GCNs).

473 4.3. Evaluation

474 4.3.1. Metrics

475 We select precision, recall and F1 as the evaluating metrics of the model,
476 where $Precision = \frac{TP}{TP+FP}$, $Recall = \frac{TP}{TP+FN}$, $F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$. TP ,
477 FP and FN refer to true positives, false positives, and false negatives, re-
478 spectively. These metrics are required to be obtained with a certain thresh-
479 old. Hence, due to the different threshold selection methods among different
480 tasks, there exist large differences in the metric values. Therefore, to avoid
481 introducing additional hyperparameters, we report the evaluation metrics
482 values when the optimal F1 value is obtained. The threshold value is deter-
483 mined by traversing between the maximum and the minimum scores of the
484 testing dataset.

485 We calculate the condition scores that decide the abnormal degree of
486 the overall dataset based on eq. (12). Noted that in unsupervised anomaly
487 detection for MTS (USAD) [66] and temporal hierarchical one-class network
488 (THOC) [73], the authors applied a specific evaluation method, called point
489 adjust, making F1 value higher and close to 1. It has been proved that
490 the capability of the model may be highly evaluated [74]. Hence, for the
491 comparison, we apply the open-source code of USAD, and utilize the same
492 parameters of the model in this paper to calculate the performance metrics
493 without adjustment.

494 4.3.2. Setup

495 We use Pytorch to achieve the HFN and its variants. Moreover, the model
496 is trained on a server with Intel(R) Xeon(R) Gold 5218R CPU @ 2.1GHz and
497 NVIDIA GeForce RTX 3090 graphics cards. We select Adam optimizer to
498 train the model. Meanwhile, we adopt early stopping to relieve overfitting.
499 The maximum training epoch is set to be 100. If the loss is less than 0.0001
500 after 10 epochs, the training stops automatically and the optimal model is
501 saved.

502 The proposed HFN method and the compared baseline models have
503 strived to maintain a similar level of complexity in parameter settings, en-
504 suring a fair comparison. For classical anomaly detection models, including
505 PCA and Isolation Forest, we have maintained the parameter settings at a
506 relatively standard level. The 'contamination' parameter for PCA has been
507 set to 0.05. In Isolation Forest, we opted for 100 isolation trees, each trained
508 using all features. In LightGBM, the 'num_boost_round' parameter has
509 been set to 1000 to ensure the model has a sufficient number of epochs for
510 training.

511 Regarding deep learning models, in LSTM-NDT, we employed a 4-layer
512 LSTM network, with each layer having 128 hidden nodes. Similarly, in
513 LSTM-VAE, a 4-layer LSTM network was used with 128 nodes in each hid-
514 den layer and a latent space dimension of 32. The parameter settings for the
515 DAGMM model align with those specified by the authors in the open-source
516 code, utilizing a Gaussian Mixture Model composed of four individual Gaus-
517 sian models. For the USAD model, a window length of 15 and a latent space
518 dimension of 40 were set. In the MTAD-GAT model, a convolutional kernel
519 size of 7 was chosen, and the hidden dimensions for the temporal and spatial
520 graph attention networks were set to 150. The prediction and reconstruction
521 networks comprise a 4-layer GRU network. The GDN model has a hidden
522 layer dimension of 128, an output layer with 64 hidden nodes, and a graph
523 network with 4 layers.

524 Finally, for the proposed HFN method, we selected a structure with a
525 hidden layer dimension of 64 and 4 layers in the graph network to ensure con-
526 sistency with other deep learning models. This configuration aims to provide
527 each model with similar capabilities in learning data representations, facili-
528 tating a more equitable evaluation of their performance in anomaly detection
529 tasks.

Table 3: Precision, recall and F1 values of HFN and all baseline methods on different datasets.

| Model | SWaT | | | WADI | | | WTD | | |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | Pre. | Rec. | F1 | Pre. | Rec. | F1 | Pre. | Rec. | F1 |
| PCA | 0.249 | 0.216 | 0.230 | 0.395 | 0.056 | 0.100 | 0.160 | 0.513 | 0.244 |
| IF | 0.951 | 0.588 | 0.727 | 0.299 | 0.158 | 0.207 | 0.278 | <u>0.953</u> | 0.430 |
| LightGBM | 0.783 | 0.666 | 0.719 | <u>0.989</u> | 0.153 | 0.270 | 0.237 | 0.602 | 0.340 |
| LSTM-NDT | 0.982 | 0.688 | 0.809 | 0.758 | 0.328 | 0.457 | 0.365 | 0.736 | 0.497 |
| LSTM-VAE | 0.962 | 0.599 | 0.740 | 0.878 | 0.145 | 0.250 | 0.165 | 0.550 | 0.254 |
| DAGMM | 0.470 | 0.666 | 0.551 | 0.544 | 0.267 | 0.360 | 0.164 | 0.242 | 0.195 |
| OmniAnomaly | 0.983 | 0.650 | 0.782 | 0.995 | 0.130 | 0.230 | - | - | - |
| USAD | 0.985 | 0.661 | 0.792 | 0.995 | 0.132 | 0.233 | 0.157 | 0.417 | 0.228 |
| MAD-GAN | 0.990 | 0.637 | 0.770 | 0.414 | 0.339 | 0.370 | - | - | - |
| MTAD-GAT | 0.991 | 0.633 | 0.772 | 0.988 | 0.153 | 0.265 | 0.128 | 0.397 | 0.193 |
| GDN | <u>0.994</u> | 0.681 | 0.810 | 0.975 | 0.402 | 0.570 | 0.385 | 0.937 | 0.546 |
| TranAD | 0.976 | 0.699 | 0.815 | 0.353 | 0.829 | 0.495 | 0.305 | 0.715 | 0.428 |
| FuSAGNet | 0.988 | 0.726 | 0.837 | 0.830 | 0.479 | <u>0.607</u> | - | - | - |
| MAD-SGCN | 0.986 | 0.690 | 0.823 | 0.564 | 0.399 | 0.552 | 0.416 | 0.688 | 0.518 |
| HFN | 0.973 | 0.758 | 0.852 | 0.827 | 0.413 | 0.551 | 0.505 | 0.837 | 0.630 |

530 *4.4. Experimental analysis*

531 The optimal metric values are shown in bold in Table 3. For the datasets
532 SWaT and WADI, we refer to the results in USAD [66] and graph deviation
533 network (GDN) [54]. For WTD dataset, to guarantee the objectivity of the
534 results, we only report the metrics from the obtained open code approaches.

535 *4.4.1. Performance comparison of anomaly detection*

536 To demonstrate the performance of the proposed model, we evaluated
537 the precision, recall, and F1 of all methods on the test set. We can ob-
538 serve from Table 3 that HFN shows a good abnormal detection capability
539 with remarkable performance improvements on SWaT and WTD. The im-
540 provement range of the proposed approach is 5% to 14%, as compared with
541 the optimal baseline models. The optimal baseline GDN outperforms our
542 approach in terms of F1; however, our approach has a more optimal re-
543 call rate. It is acceptable in real scenarios because we hope to detect more
544 anomalies. In short, HFN outperforms the selected baselines in terms of
545 the overall performances, because it not only concerns with the traditional
546 spatial-temporal correlation, but also obtains its heterogeneous attributes
547 from different types of data, making the model more robust. Moreover, we
548 observe that prediction-based algorithms such as HFN, GDN and LSTM-
549 NDT outperform the reconstruction-based algorithms such as LSTM-VAE
550 and USAD on these datasets, indicating that the prediction-based models
551 have an advantage in the streaming abnormal detection tasks with a single-
552 timestamp value as the target. The temporal information is also very vital
553 in the tasks for MTS abnormal detection. The results of LSTM-NDT show
554 that HFN outperforms all baselines except GDN. The PCA result is dissat-
555 isfactory, because it gives more attentions to the point anomalies without
556 spatial-temporal correlation being considered.

557 Specifically, among these abnormal detection approaches, GDN, MTAD-
558 GAT and HFN adopt the graph attention network to capture the tem-
559 poral and feature correlations. Therefore, these types of models achieve
560 good results on all datasets. GDN approach recodes multidimensional data
561 at each moment, and utilizes its strong structural learning capability of
562 graph attention network to learn coupling relations between different sensors.
563 However, it does not consider the heterogeneity of data. MTAD-GAT ap-
564 proach also captures time-dimension information through an attention mech-
565 anism. Although it considers the spatial-temporal correlation of MTS, it
566 requires a configuration of hyper-parameters for fusing the prediction-based

567 and reconstruction-based condition scores, leading to the evident differences
568 in results when this approach is applied to different datasets. Compared to
569 the recently introduced transformer-based TranAD, as well as the spatial-
570 temporal graph networks FuSAGNet and MAD-SGCN, HFN continues to
571 exhibit superior performance on the SWaT and WTD datasets. However,
572 the most recent experimental outcomes suggest that FuSAGNet achieved the
573 top results on the WADI dataset. Nonetheless, our attempts to reproduce
574 this outcome using the authors’ open-sourced code were unsuccessful.

575 Furthermore, although we processed different types of data separately,
576 our optimization efforts were predominantly concentrated on enhancing the
577 network structure without introducing a significant increase in complexity.
578 Consequently, the processing time did not exhibit a substantial increase when
579 handling the same amount of data.

580 4.4.2. Ablation experiment

581 We utilize SWaT and WADI datasets to study the necessity of five com-
582 ponents of our approach, namely, node embedding similarity matrix (NE),
583 node feature similarity matrix (NF), discrete feature subgraph (DFS), contin-
584 uous feature subgraph (CFS), and hybrid feature subgraph (HFS). As shown
585 in Figure 4, we successively exclude the corresponding component from the
586 experiments to observe its effect on the model performance. The key idea
587 of our approach is to learn the potential steady representations from het-
588 erogeneous MTS. Hence, first, we exclude NE or NF to study whether the
589 heterogeneous information is learned. Second, we discuss the anomaly detec-
590 tion performance when we only use HFS or DFS and CFS. Specifically:

- 591 • Excluding NF (expressed as "-NF") degrades the overall performance
592 of the approach and has a great influence on WADI dataset. This
593 indicates that NF is in favor of feature extraction with high-dimensional
594 dataset for the model; however, NF is not the key factor to determine
595 the model performance.
- 596 • Excluding NE (expressed as "-NE") degrades the performances clearly,
597 which implies that NE has an evident advantage in the graph structure
598 learning process.
- 599 • Excluding DFS and CFS (expressed as "-DFS" and "-CFS") degrades
600 the model performance; however, the descend range of model perfor-
601 mance is less than that of NE. This approach is actually degenerated

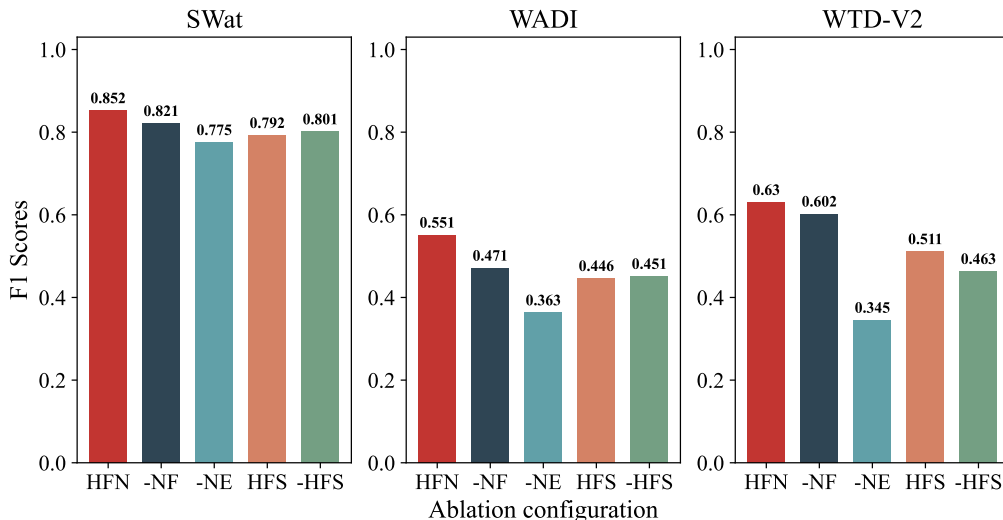


Figure 4: Effects of different HFN components on anomaly detection performance.

602 to the processing of isomorphic graphs, leading to the loss of heteroge-
 603 neous information.

- 604 • Excluding HFS (expressed as "-HFS") degrades the model performance;
 605 however, it is superior to the cases when DFS and CFS are totally ex-
 606 cluded. This indicates that the interaction between different types of
 607 sensors in the hybrid subgraph plays a complementary role in extract-
 608 ing the follow-up HFN heterogeneous information.

609 To sum up, it is necessary to extract heterogeneous structure information
 610 in the MTS datasets. The heterogeneous information can present different
 611 weights in the model according to the attention mechanism, which helps to
 612 improve the abnormal detection performance.

613 4.4.3. Case study

614 (1) Anomaly detection analysis

615 Figure 5 shows the abnormal detection results on SWaT testing dataset,
 616 where Figure 5(a) represents the actual data anomalies on this dataset, in-
 617 cluding network and physical attacks directed at the Secure Water Treatment
 618 (SWaT) testbed within the continuous four days. The data are labeled as
 619 1 if the system is attacked at a certain timestamp; otherwise, it is labeled
 620 as 0. Figure 5(b) represents the results of HFN anomaly detection, where

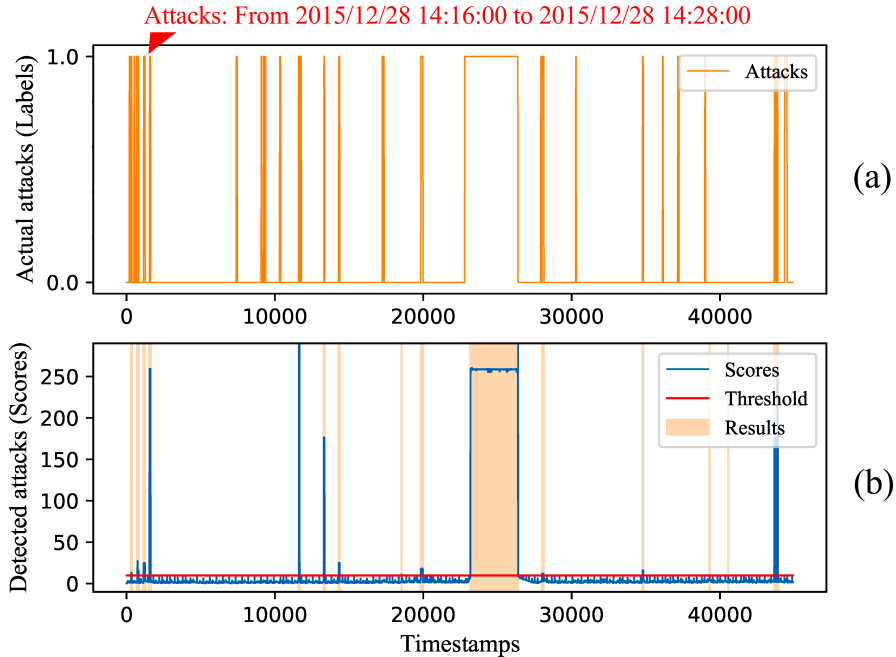


Figure 5: SWaT dataset anomaly detection results.

621 the orange shadow represents the detected anomalies, the blue curve repre-
 622 sents the condition scores calculated as described in Section 3.6, and the red
 623 straight line represents the threshold when the optimal F1 is obtained on
 624 the testing dataset. It can be seen from Figure 5(b) that, aside from a few
 625 anomalies that are very difficult to distinguish possibly due to labeling errors,
 626 our approach accurately identifies the most anomalies. According to the in-
 627 structions provided by SWAT dataset [62], we select an attack case to further
 628 interpret the abnormal detection capability of HFN. As shown in Figure 5(a),
 629 the attack starts from 14:16:00 28/12/2015 to 14:28:00 28/12/2015 against
 630 FIT401, UV401 and P501, where FIT401 is the flow transmitter for mea-
 631 suring the flow of UV de-chlorinator, UV401 is de-chlorinator for removing
 632 chlorine from water, and P501 is pump actuator for pumping the dechlori-
 633 nated water to reverse osmosis. During the attack, as shown in Figure 6, the
 634 flow value (continuous value) of FIT401 is set twice to the value deviating
 635 from the normal mode. Meanwhile, the actuators UV401 and P501 (discrete
 636 value), which should be kept to an open state, are forcefully closed.

637 Figure 6 shows the curves of actual and predicted values of attack-related
638 sensors and actuators and the HFN anomaly detection results. In order to
639 reduce the influence of data dimensions and accelerate the convergence of
640 the model, we have standardized the values of the dataset by min-max nor-
641 malization. It is worth noting that we used the same normalized parameters
642 for both the training dataset and the testing dataset, which is why the nor-
643 malized data of the testing data shown in Figure 6 has negative values. This
644 was done to reduce the impact of **testing data** information leakage on the
645 model performance. In the real water treatment process, the unit of the flow
646 sensors values are gallons per minute (GPM), while the actuators have two
647 conditions: 0 means turn on and - 1 means turn off.

648 It can be seen from Figure 6(a), (c) and (d) that before the attack, the
649 predicted values of HFN are consistent with the actual values, where the
650 prediction for both continuous variables and discrete variables achieves good
651 results. In the attack process, the flow variation arises from the prediction
652 result of FIT401 and UV401 simultaneously. This is due to the interaction
653 among these variables in the actual water treatment system. A larger devia-
654 tion between the predicted value and the actual value would provide a better
655 basis for abnormal detection. Note that although the experiment personnel
656 did not launch the attack on FIT504 sensor in the attack process, we can see
657 from Figure 6(b) and (d) that the value changes of FIT504 are still detected,
658 which is due to being abnormally closed caused by the attack on P501. We
659 can observe from the detection results in Figure 6(e) that the proposed ap-
660 proach shows a good detection capability of such complex anomalies. These
661 anomalies have been resulted from attacks to different types of sensors, in-
662 cluding continuity, discreteness and their correlation, which represent real
663 scenarios.

664 (2) Anomaly localization analysis

665 From the above analysis, we can see that our method can successfully
666 detect the occurrence of anomalies. However, we cannot assume that all
667 the variables in a real complicated system are of the same significance. In
668 other words, the variables associated with a particular system component
669 will be impacted to varying degrees of operation when that component is
670 attacked or behaves abnormally. Therefore, it is necessary to locate vari-
671 ables that have been strongly impacted by the attack, thus helping system
672 maintenance personnel to rapidly find and solve the problems. We use the
673 prediction error of each time-series sensor to represent the condition score of
674 the sequence where the sensor with the maximum score times is considered

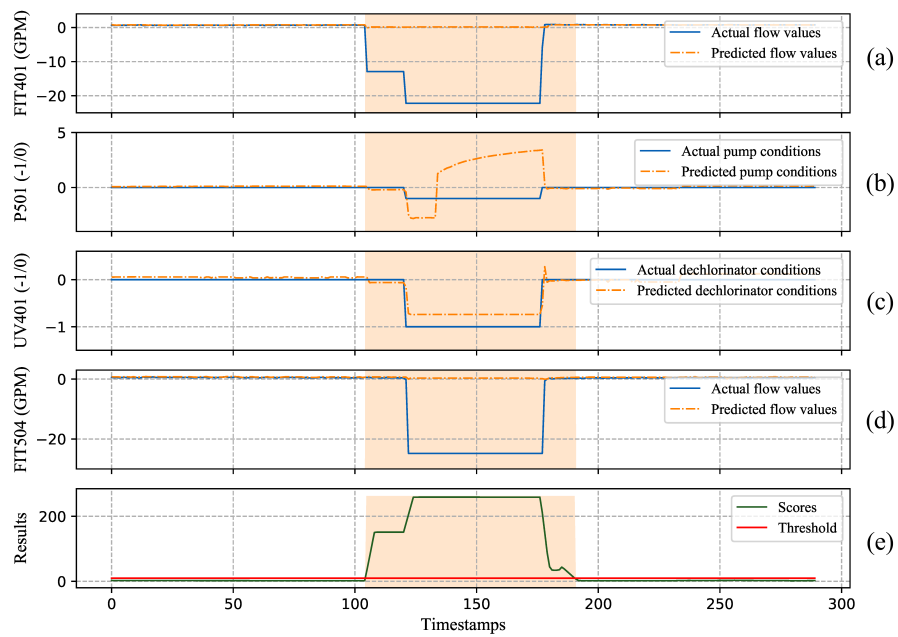


Figure 6: Abnormal detection case. The orange shadow represents the detected anomalies. The blue curve denotes the actual value of sensor or actuator. The orange dotted line represents the predicted value. The green and red curves in Fig. 6(e) represent condition score and the threshold of the optimal F1, respectively.

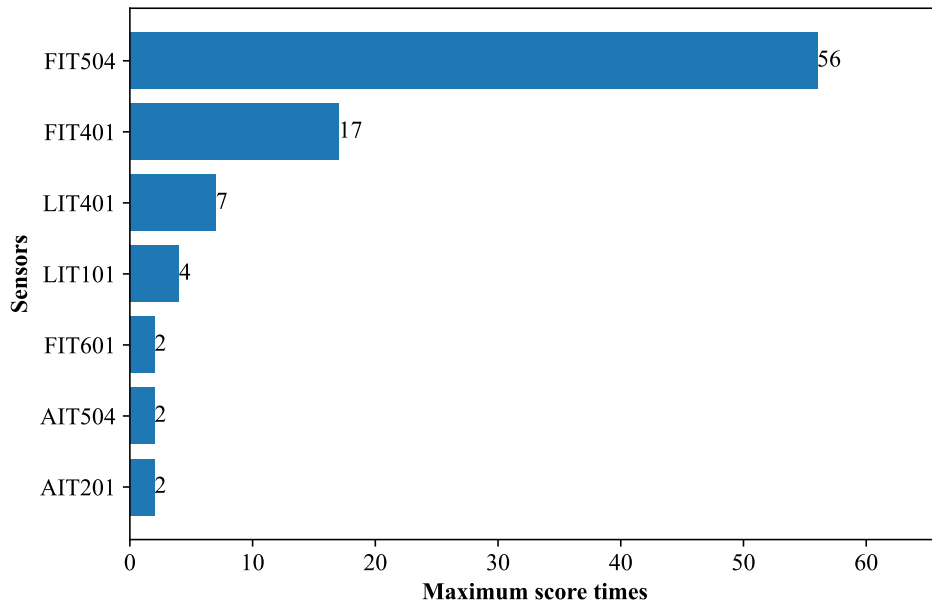


Figure 7: Maximum number of sensor scores in the abnormal dataset.

675 to have the possibility of the biggest anomalies. Figure 7 shows the number
 676 of times when the condition scores are above the threshold for different sen-
 677 sors within the attack period in the case analysis. It can be known from the
 678 figure that the sensors FIT504 and FIT401 have the maximum score times,
 679 which is consistent with the attacks where the experiment personnel made to
 680 the sensor FIT401 and pump actuator P501 during the tests. The turn-off
 681 attacks on P501 caused a sharp drop in the FIT504 flow values, as shown in
 682 Figure 6(d), since they are physically connected. On the contrary, we can
 683 also speculate which component of the system has been attacked or abnormal
 684 according to the maximum score times. In this case, during real operation
 685 and maintenance, particular **attention** should be paid to and checks should
 686 be made on the locations relating to FIT504, FIT401, and LIT401.

687 4.5. Feasibility analysis

688 To further illustrate how the heterogeneous relation in time series is
 689 learned and takes effect on the abnormal detection, we explain it through
 690 the similarity matrix before and after the anomalies due to the attacks. Fig-
 691 ure 8 and Figure 9 represent different similarity matrices before and after
 692 the attack on SWaT, respectively. Its similarity value range is $[-1,1]$, and

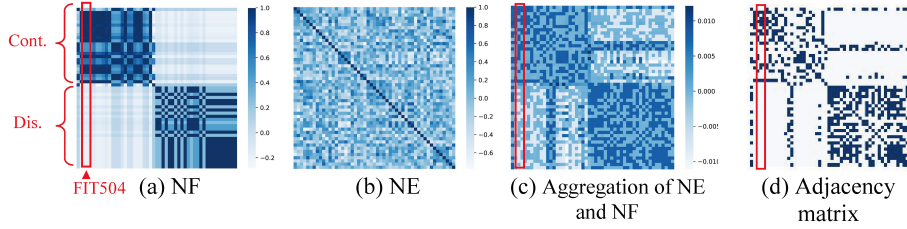


Figure 8: Example of similarity subgraphs under normal conditions.

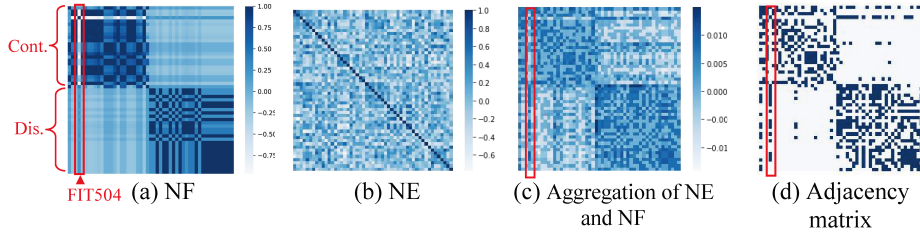


Figure 9: Example of similarity subgraphs under attack conditions.

693 the closer to 1, the stronger the similarity is. Overall, HFN aggregates the
 694 similarities of sensor signals from different perspectives to represent its het-
 695 erogeneous information. Embedding similarity matrix learns the structural
 696 information among different sensors globally from the training data. Hence,
 697 similar features are shown in Figure 8(b) and Figure 9(b) under abnormal
 698 and normal states. However, concerning the feature similarity, we can see
 699 clearly that there exist significant differences in feature similarity between
 700 Figure 8(a) and Figure 9(a) at different timestamps, because the data vary
 701 with time. Ignoring this part of information always degrades the abnormal
 702 detection performance.

703 Specifically, as shown in Figure 8(a), before the attacks on FIT401, UV401
 704 and P501, the similarity values of FIT504 flow value and other continuous
 705 variable sensor values are close to 1. However, after the attack, we can see
 706 from Figure 9(a) that the similarity value varies to -0.75. The sudden change
 707 indicates that the sensor anomalies occur, while there are slight variations in
 708 the embedding similarity. After comparing the adjacent matrix before and
 709 after the attack in Figure 8(d) and Figure 9(d), we can find that the changes
 710 in feature similarity cause the changes in the connection relation to improve
 711 the ability of the algorithm in capturing dynamic feature correlation. This
 712 further demonstrates that HFN, by aggregating the global data learning-

713 based embedding similarity matrix and the feature similarity matrix at a
714 specific timestamp, can better capture the normal and abnormal conditions
715 in MTS.

716 5. Conclusions

717 In this paper, we propose a novel heterogeneous feature network for MTS
718 anomaly detection. This approach is able to learn the complex heteroge-
719 neous structural information and temporal information between MTS data.
720 Therefore, it is suitable for abnormal detection in real scenarios where the
721 dataset comprises continuous numerical variables and discrete categorical
722 variables simultaneously. The extensive experiments indicate that our ap-
723 proach outperforms the baseline models by assessing two open datasets from
724 water treatment plants and a private dataset from a wind power plant. **Par-**
725 **ticularly noteworthy is its significant performance improvement on the SWaT**
726 **and WTD-V2 datasets, where the F1 score increased by 5% and 14%, respec-**
727 **tively, compared to the best baseline.** Furthermore, our approach demon-
728 strates a good abnormal interpretability and can help operation and main-
729 tenance personnel rapidly discover and locate the anomalies.

730 In the future, we will continue to explore various avenues to enhance
731 the proposed algorithm. We plan to extend its capabilities by incorporat-
732 ing more real and complex heterogeneous datasets, encompassing combined
733 time series data and textual information. This expansion aims to boost the
734 accuracy and practicality of the approach. While our method excels in di-
735 verse data handling, potential challenges in computational efficiency may
736 arise with larger datasets. Future efforts will be directed towards optimizing
737 the algorithm for improved scalability, especially in scenarios involving more
738 extensive network scales. Additionally, we aim to investigate the impact of
739 varying sampling intervals on our method across different datasets, thereby
740 broadening its applicability.

741 Acknowledgement

742 The work is supported by National Natural Science Foundation of China
743 (62006236), NUDT Research Project (ZK20-10), National Key Research and
744 Development Program of China (2020YFA0709803), Hunan Provincial Natu-
745 ral Science Foundation (2020JJ5673), National Science Foundation of China

746 (U1811462), National Key R&D project by Ministry of Science and Tech-
747 nology of China (2018YFB1003203), and Autonomous Project of HPCL
748 (201901-11, 202101-15).

749 **References**

- 750 [1] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, T. Soderstrom,
751 Detecting spacecraft anomalies using lstms and nonparametric dynamic
752 thresholding.
- 753 [2] H. Zhou, K. Yu, X. Zhang, G. Wu, A. Yazidi, Contrastive autoencoder
754 for anomaly detection in multivariate time series, *Information Sciences*
755 610 (2022) 266–280.
- 756 [3] J. Zhan, R. Wang, L. Yi, Y. Wang, Z. Xie, Health assessment methods
757 for wind turbines based on power prediction and mahalanobis distance,
758 *International Journal of Pattern Recognition and Artificial Intelligence*
759 33 (02) (2019) 1951001.
- 760 [4] G. Pang, C. Shen, L. Cao, A. V. D. Hengel, Deep learning for anomaly
761 detection: A review, *ACM Computing Surveys (CSUR)* 54 (2) (2021)
762 1–38.
- 763 [5] H. Mamdouh Farghaly, M. Y. Shams, T. Abd El-Hafeez, Hepatitis c
764 virus prediction based on machine learning framework: a real-world case
765 study in egypt, *Knowledge and Information Systems* 65 (6) (2023) 2595–
766 2617.
- 767 [6] P. D. Rosero-Montalvo, Z. István, P. Tözün, W. Hernandez, Hybrid
768 anomaly detection model on trusted iot devices, *IEEE Internet of Things*
769 *Journal* (2023).
- 770 [7] P. Wu, J. Liu, Learning causal temporal relation and feature discrimi-
771 nation for anomaly detection, *IEEE Transactions on Image Processing*
772 30 (2021) 3513–3527.
- 773 [8] B. Lindemann, B. Maschler, N. Sahlab, M. Weyrich, A survey on
774 anomaly detection for technical systems using lstm networks, *Computers*
775 *in Industry* 131 (2021) 103498.

- 776 [9] J. Yang, L. Zhang, C. Chen, Y. Li, R. Li, G. Wang, S. Jiang, Z. Zeng,
777 A hierarchical deep convolutional neural network and gated recurrent
778 unit framework for structural damage detection, *Information Sciences*
779 540 (2020) 117–130.
- 780 [10] A. Blázquez-García, A. Conde, U. Mori, J. A. Lozano, A review on
781 outlier/anomaly detection in time series data, *ACM Computing Surveys*
782 (CSUR) 54 (3) (2021) 1–33.
- 783 [11] H. Zhao, Y. Wang, J. Duan, C. Huang, Q. Zhang, Multivariate time-
784 series anomaly detection via graph attention network (2020) 841–850.
- 785 [12] S. Du, T. Li, Y. Yang, S.-J. Horng, Multivariate time series forecasting
786 via attention-based encoder–decoder framework, *Neurocomputing* 388
787 (2020) 269–279.
- 788 [13] T.-Y. Kim, S.-B. Cho, Predicting residential energy consumption using
789 cnn-lstm neural networks, *Energy* 182 (2019) 72–81.
- 790 [14] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image
791 recognition (2016) 770–778.
- 792 [15] Z. Wu, S. Pan, G. Long, J. Jiang, X. Chang, C. Zhang, Connecting the
793 dots: Multivariate time series forecasting with graph neural networks
794 753–763.
- 795 [16] Y. Guo, W. Liao, Q. Wang, L. Yu, T. Ji, P. Li, Multidimensional time
796 series anomaly detection: A gru-based gaussian mixture variational au-
797 toencoder approach 97–112.
- 798 [17] D. Li, D. Chen, B. Jin, L. Shi, J. Goh, S.-K. Ng, Mad-gan: Multivari-
799 ate anomaly detection for time series data with generative adversarial
800 networks (2019) 703–716.
- 801 [18] V. Borisov, T. Leemann, K. Seßler, J. Haug, M. Pawelczyk, G. Kas-
802 neci, Deep neural networks and tabular data: A survey, arXiv preprint
803 arXiv:2110.01889 (2021).
- 804 [19] A. Nazabal, P. M. Olmos, Z. Ghahramani, I. Valera, Handling incom-
805 plete heterogeneous data using vaes, *Pattern Recognition* 107 (2020)
806 107501.

- 807 [20] Y. Sun, J. Han, X. Yan, P. S. Yu, T. Wu, Pathsim: Meta path-based top-
808 k similarity search in heterogeneous information networks, Proceedings
809 of the VLDB Endowment 4 (11) (2011) 992–1003.
- 810 [21] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, P. S. Yu, Heterogeneous
811 graph attention network 2022–2032.
- 812 [22] M. Lngkvist, L. Karlsson, A. Loutfi, A review of unsupervised feature
813 learning and deep learning for time-series modeling, Pattern Recognition
814 Letters 42 (2014) 11–24.
- 815 [23] J. Li, H. Izakian, W. Pedrycz, I. Jamal, Clustering-based anomaly de-
816 tection in multivariate time series data, Applied Soft Computing 100 (4)
817 (2020) 106919.
- 818 [24] J. Zhao, K. Liu, W. Wang, Y. Liu, Adaptive fuzzy clustering based
819 anomaly data detection in energy system of steel industry, Information
820 Sciences 259 (2014) 335–345.
- 821 [25] M. Jones, D. Nikovski, M. Imamura, T. Hirata, Anomaly detection in
822 real-valued multidimensional time series.
- 823 [26] E. G. S. Nascimento, O. de Lira Tavares, A. F. De Souza, A cluster-
824 based algorithm for anomaly detection in time series using mahalanobis
825 distance 622.
- 826 [27] G. Pu, L. Wang, J. Shen, F. Dong, A hybrid unsupervised clustering-
827 based anomaly detection method, Tsinghua Science and Technology
828 26 (2) (2020) 146–153.
- 829 [28] W. Jia, R. M. Shukla, S. Sengupta, Anomaly detection using supervised
830 learning and multiple statistical methods (2019) 1291–1297.
- 831 [29] N. Görnitz, M. Kloft, K. Rieck, U. Brefeld, Toward supervised anomaly
832 detection, Journal of Artificial Intelligence Research 46 (2013) 235–262.
- 833 [30] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K.-
834 R. Müller, M. Kloft, Deep semi-supervised anomaly detection, arXiv
835 preprint arXiv:1906.02694 (2019).

- 836 [31] M. E. Villa-Pérez, M. A. Alvarez-Carmona, O. Loyola-González, M. A.
837 Medina-Pérez, J. C. Velazco-Rossell, K.-K. R. Choo, Semi-supervised
838 anomaly detection algorithms: A comparative summary and future re-
839 search directions, *Knowledge-Based Systems* 218 (2021) 106878.
- 840 [32] R. Wang, X. Ma, C. Jiang, Y. Ye, Y. Zhang, Heterogeneous informa-
841 tion network-based music recommendation system in mobile networks,
842 *Computer Communications* 150 (2020) 429–437.
- 843 [33] X. Liang, Y. Ma, G. Cheng, C. Fan, Y. Yang, Z. Liu, Meta-path-based
844 heterogeneous graph neural networks in academic network, *International*
845 *Journal of Machine Learning and Cybernetics* 13 (6) (2022) 1553–1569.
- 846 [34] X. Deng, F. Long, B. Li, D. Cao, Y. Pan, An influence model based on
847 heterogeneous online social network for influence maximization, *IEEE*
848 *Transactions on Network Science and Engineering* 7 (2) (2019) 737–749.
- 849 [35] X. Chen, J. Yin, J. Qu, L. Huang, Mdhgi: matrix decomposition and
850 heterogeneous graph inference for mirna-disease association prediction,
851 *PLoS computational biology* 14 (8) (2018) e1006418.
- 852 [36] Y. Sun, J. Gao, X. Hong, B. Mishra, B. Yin, Heterogeneous tensor de-
853 composition for clustering via manifold optimization, *IEEE transactions*
854 *on pattern analysis and machine intelligence* 38 (3) (2015) 476–489.
- 855 [37] Y. Liu, X. Luo, X. Yang, Semantics and structure based recommenda-
856 tion of similar legal cases 388–395.
- 857 [38] C. Wang, C. H. Chi, Z. Wei, R. Wong, Coupled interdependent attribute
858 analysis on mixed data (2015).
- 859 [39] S. Jian, L. Cao, G. Pang, L. Kai, G. Hang, Embedding-based represen-
860 tation of categorical data by hierarchical value coupling learning (2017).
- 861 [40] H. Cai, V. W. Zheng, K. C.-C. Chang, A comprehensive survey of graph
862 embedding: Problems, techniques, and applications, *IEEE Transactions*
863 *on Knowledge and Data Engineering* 30 (9) (2018) 1616–1637.
- 864 [41] J. Zhao, X. Wang, C. Shi, B. Hu, G. Song, Y. Ye, Heterogeneous graph
865 structure learning for graph neural networks 35 (5) (2021) 4697–4705.

- 866 [42] X. Fu, J. Zhang, Z. Meng, I. King, Magnn: Metapath aggregated graph
867 neural network for heterogeneous graph embedding 2331–2341.
- 868 [43] X. Wang, N. Liu, H. Han, C. Shi, Self-supervised heterogeneous graph
869 neural network with co-contrastive learning 1726–1736.
- 870 [44] Z. Hu, Y. Dong, K. Wang, Y. Sun, Heterogeneous graph transformer
871 2704–2710.
- 872 [45] L. Yang, Z. Xiao, W. Jiang, Y. Wei, Y. Hu, H. Wang, Dynamic hetero-
873 geneous graph embedding using hierarchical attentions (2020) 425–432.
- 874 [46] M. Chen, C. Huang, L. Xia, W. Wei, Y. Xu, R. Luo, Heterogeneous
875 graph contrastive learning for recommendation (2023) 544–552.
- 876 [47] Y. Zhu, Y. Xu, H. Cui, C. Yang, Q. Liu, S. Wu, Structure-enhanced
877 heterogeneous graph contrastive learning (2022) 82–90.
- 878 [48] P. Bloomfield, Fourier analysis of time series: an introduction, John
879 Wiley & Sons, 2004.
- 880 [49] R. Shwartz-Ziv, A. Armon, Tabular data: Deep learning is not all you
881 need, *Information Fusion* 81 (2022) 84–90.
- 882 [50] Y. Zhu, W. Xu, J. Zhang, Q. Liu, S. Wu, L. Wang, Deep graph
883 structure learning for robust representations: A survey, arXiv preprint
884 arXiv:2103.03036 (2021).
- 885 [51] R. Li, S. Wang, F. Zhu, J. Huang, Adaptive graph convolutional neural
886 networks 32.
- 887 [52] X. Wang, M. Zhu, D. Bo, P. Cui, C. Shi, J. Pei, Am-gcn: Adaptive
888 multi-channel graph convolutional networks 1243–1253.
- 889 [53] Y. Chen, L. Wu, M. Zaki, Iterative deep graph learning for graph neu-
890 ral networks: Better and robust node embeddings, *Advances in Neural
891 Information Processing Systems* 33 (2020) 19314–19326.
- 892 [54] A. Deng;, B. Hooi., Graph neural network-based anomaly detection in
893 multivariate time series, *aaai2021* (2021).

- 894 [55] D. Yu, R. Zhang, Z. Jiang, Y. Wu, Y. Yang, Graph-revised convolutional
895 network 378–393.
- 896 [56] D. Luo, W. Cheng, W. Yu, B. Zong, J. Ni, H. Chen, X. Zhang, Learning
897 to drop: Robust graph neural network via topological denoising 779–
898 787.
- 899 [57] Q. Sun, J. Li, H. Peng, J. Wu, X. Fu, C. Ji, P. S. Yu, Graph struc-
900 ture learning with variational information bottleneck, arXiv preprint
901 arXiv:2112.08903 (2021).
- 902 [58] X. Gao, W. Hu, Z. Guo, Exploring structure-adaptive graph learning
903 for robust semi-supervised classification 1–6.
- 904 [59] P. Kaewprapha, P. Prempaneerach, V. Singh, T. Tinikul, N. Intarangsi,
905 T. Kijkanjanarat, Predicting full load, partial load efficiency of a com-
906 bined cycle power plant using machine learning methods 11–16.
- 907 [60] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,
908 L. Kaiser, I. Polosukhin, Attention is all you need, arXiv (2017).
- 909 [61] H. Ren, B. Xu, Y. Wang, C. Yi, C. Huang, X. Kou, T. Xing, M. Yang,
910 J. Tong, Q. Zhang, Time-series anomaly detection service at microsoft
911 3009–3017.
- 912 [62] J. Goh, S. Adepur, K. N. Junejo, A. Mathur, A dataset to support re-
913 search in the design of secure water treatment systems (2017) 88–99.
- 914 [63] C. M. Ahmed, V. R. Palleti, A. P. Mathur, Wadi: a water distribution
915 testbed for research in the design of secure cyber physical systems.
- 916 [64] J. Zhan, S. Wang, X. Ma, C. Wu, C. Yang, D. Zeng, S. Wang, Stgat-
917 mad: Spatial-temporal graph attention network for multivariate time
918 series anomaly detection 3568–3572.
- 919 [65] D. Park, Y. Hoshi, C. C. Kemp, A multimodal anomaly detector
920 for robot-assisted feeding using an lstm-based variational autoencoder,
921 IEEE Robotics and Automation Letters PP (99) (2017).
- 922 [66] J. Audibert, P. Michiardi, F. Guyard, S. Marti, M. A. Zuluaga, Usad:
923 Unsupervised anomaly detection on multivariate time series 3395–3404.

- 924 [67] J. Camacho, A. Pérez-Villegas, P. García-Teodoro, G. Maciá-Fernández,
925 Pca-based multivariate statistical network monitoring for anomaly de-
926 tection, *Computers & Security* 59 (2016) 118–137.
- 927 [68] T. L. Fei, M. T. Kai, Z. H. Zhou, Isolation forest.
- 928 [69] M. Qi, Lightgbm: A highly efficient gradient boosting decision tree.
- 929 [70] S. Tuli, G. Casale, N. R. Jennings, Tranad: Deep transformer networks
930 for anomaly detection in multivariate time series data, arXiv preprint
931 arXiv:2201.07284 (2022).
- 932 [71] S. Han, S. S. Woo, Learning sparse latent graph representations for
933 anomaly detection in multivariate time series (2022) 2977–2986.
- 934 [72] P. Qi, D. Li, S.-K. Ng, Mad-sgcn: Multivariate anomaly detection with
935 self-learning graph convolutional networks (2022) 1232–1244.
- 936 [73] L. Shen, Z. Li, J. Kwok, Timeseries anomaly detection using temporal
937 hierarchical one-class network, *Advances in Neural Information Process-*
938 *ing Systems* 33 (2020) 13016–13026.
- 939 [74] S. Kim, K. Choi, H. S. Choi, B. Lee, S. Yoon, Towards a rigorous eval-
940 uation of time-series anomaly detection (2021).