Motif-Consistent Counterfactuals with Adversarial Refinement for Graph-Level Anomaly Detection

Chunjing Xiao Henan University China chunjingxiao@gmail.com

Yanlong Huang University of Electronic Science and Technology of China hyloong77@gmail.com Shikang Pang Henan University China pangsk0604@henu.edu.cn

Goce Trajcevski Iowa State University USA gocet25@iastate.edu Wenxin Tai
University of Electronic Science and
Technology of China
wxtai@outlook.com

Fan Zhou*
University of Electronic Science and
Technology of China
fan.zhou@uestc.edu.cn

Abstract

Graph-level anomaly detection is significant in diverse domains. To improve detection performance, counterfactual graphs have been exploited to benefit the generalization capacity by learning causal relations. Most existing studies directly introduce perturbations (e.g., flipping edges) to generate counterfactual graphs, which are prone to alter the semantics of generated examples and make them off the data manifold, resulting in sub-optimal performance. To address these issues, we propose a novel approach, Motif-consistent Counterfactuals with Adversarial Refinement (MotifCAR), for graph-level anomaly detection. The model combines the motif of one graph, the core subgraph containing the identification (category) information, and the contextual subgraph (nonmotif) of another graph to produce a raw counterfactual graph. However, the produced raw graph might be distorted and cannot satisfy the important counterfactual properties: Realism, Validity, Proximity and Sparsity. Towards that, we present a Generative Adversarial Network (GAN)-based graph optimizer to refine the raw counterfactual graphs. It adopts the discriminator to guide the generator to generate graphs close to realistic data, i.e., meet the property Realism. Further, we design the motif consistency to force the motif of the generated graphs to be consistent with the realistic graphs, meeting the property Validity. Also, we devise the contextual loss and connection loss to control the contextual subgraph and the newly added links to meet the properties *Proximity* and *Spar*sity. As a result, the model can generate high-quality counterfactual graphs. Experiments demonstrate the superiority of MotifCAR.

CCS Concepts

• Computing methodologies \rightarrow Semi-supervised learning settings; Anomaly detection.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '24, August 25–29, 2024, Barcelona, Spain.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0490-1/24/08

https://doi.org/10.1145/3637528.3672050

Keywords

Graph anomaly detection, graph neural networks, representation learning, counterfactual data augmentation

ACM Reference Format:

Chunjing Xiao, Shikang Pang, Wenxin Tai, Yanlong Huang, Goce Trajcevski, and Fan Zhou. 2024. Motif-Consistent Counterfactuals with Adversarial Refinement for Graph-Level Anomaly Detection. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24), August 25–29, 2024, Barcelona, Spain.* ACM, New York, NY, USA, 9 pages. https://doi.org/10.1145/3637528.3672050

1 Introduction

Graph-level anomaly detection aims to identify graph instances that are significantly different from the majority of graphs. As a few anomalies may cause tremendous loss, detecting anomalous data has significant implications for various domains ranging from identifying abnormal proteins in biochemistry and distinguishing brain disorders in brain networks, to uncovering fraudulent activities in online social networks [2, 30]. Numerous corresponding detection methods have been introduced by taking advantage of different deep learning techniques for anomaly detection[30, 40], such as self-supervised learning [36, 51], knowledge distillation [23, 29] and tailored Graph Neural Networks (GNNs) [36, 48].

However, these deep learning models are prone to learn dataset-dependent spurious correlations based on statistical associations [19]. This might hinder well-trained models from generalizing well to newly observed anomalies, resulting in detection errors. Counterfactual data augmentation can help the model alleviate the problem of spurious correlations by learning causal relations and enhance the generalization capacity in tabular data, image data, and text data [9, 18]. While, for graph data, research on counterfactual augmentation is insufficient due to the presence of complex structure and node information. The limited existing research principally focuses on introducing perturbations into graphs or matching counterfactual data with different treatments [5, 28, 41, 53].

Whereas, when applying to anomaly detection, these methods confronts the crucial problems: (1) Perturbations might alter the graph semantics, adversely impacting model robustness. In anomaly detection, the category of a graph (e.g., normal or anomalous label) can be determined by the presence of specific edges [30]. Perturbing these edges may change its category, as illustrated in Figure 1. However, the augmented samples will be labeled with

^{*}Corresponding Author.

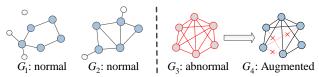


Figure 1: An example of perturbations altering semantics. G_1 and G_2 are normal graphs. While, G_3 is abnormal one as its fully connected structure deviates significantly from the normal ones. However, by pruning a few edges, abnormal graph G_3 is transformed into normal graph G_4 , leading to an altered semantic.

the same category as the original ones [11, 28, 41], resulting in incorrect labels and degraded performance. (2) These methods may generate suboptimal augmented data, leading to limited effectiveness in enhancing generalization capacity. Perturbation methods often lead to distorted examples [37], which are usually off the data manifold [18]. In such cases, deep models can be deceived because they do not generalize well to unseen test data [18]. Also, matching approaches only seek desired samples from existing data [5, 53], which cannot expand the training data.

In this paper, we propose a Motif-consistent Counterfactual data generation model with Adversarial Refinement (MotifCAR) to generate high-quality counterfactual samples for graph-level anomaly detection. In this model, to address the problem (1), instead of conducting perturbations on graphs, we introduce the discriminative motif for counterfactual data generation. According to [15], the discriminative motif is a subgraph of a graph that decides the category of this graph. Hence, it can be regarded as the core subgraph containing the identification (category) information of a given graph. Correspondingly, the remaining nodes can form a subgraph as the contextual graph. For brevity, we refer to the discriminative motif simply as motif. Based on this concept, we combine the motif of one graph and the contextual subgraph of another graph to form a raw counterfactual graph, i.e., unseen combinations of the motif and the contextual subgraph. These two source graphs can be from the same/different categories/clusters.

To address the problem (2), we design a Generative Adversarial Network (GAN)-based graph optimizer along with our tailored losses to refine the raw counterfactual graphs into the high-quality ones. According to [9, 33], good counterfactual data requires satisfying a set of properties (a.k.a. counterfactual properties), including (1) Realism: the counterfactual examples should lie close to the data manifold so that they appear realistic. (2) Validity: the model should assign the counterfactual examples to the corresponding category label in order to be valid; (3) Proximity: the distance of a counterfactual and original data should be close; (4) Sparsity: the number of perturbations on nodes/edges should be sparse. To meet the desired counterfactual properties, we design three specific losses for the graph optimizer. The refined counterfactual graphs can enlarge the distribution of training data and help the model handle the situation with varying environments.

More succinctly, our MotifCAR model is composed of two fundamental components: the raw counterfactual graph producer and the GAN-based graph optimizer. The graph producer takes as inputs two graphs, G and H, and produces a raw counterfactual graph by combining the motif of G and the contextual subgraph (non-motif) of H. This producer first merges these two graphs into one graph,

and then discards the contextual subgraph of G and the motif of H. The remaining parts (the motif of G and the contextual subgraph of H) are connected randomly to form a raw counterfactual graph. The generated graph preserves the identification (category) information of G but involves the environment characteristics of H. These unseen combinations can enlarge the distribution of training data and help the model learn transferable relations across different environments. However, the random links between both subgraphs in the raw counterfactual graphs may distort the graph structure, which does not adhere well to the desired counterfactual properties.

Hence, we further design the GAN-based graph optimizer to refine the raw graphs by adjusting the edges. The optimizer is composed of the graph generator and the discriminator. The generator aims to adjust the edges of the graphs and the discriminator tries to distinguish generated graphs from realistic ones. This adversarial training will encourage the generated graphs to conform to the patterns of realistic graphs, i.e., meet the requirement of the property (1) Realism. On the basis of the basic GAN framework, we further design three losses to improve the quality of the generated graphs. Firstly, we propose a motif consistency loss to force the motif of the generated graphs to be consistent with the realistic graphs. Since the motif contains the identification information of the graphs, this loss ensures the identification information invariant with the realistic graphs, i.e., meets the property (2) Validity. Secondly, we devise the contextual loss to keep the degree distribution of the generated contextual subgraphs to be similar to the realistic ones. Coupled with the motif consistency loss, these losses facilitate the proximity of the counterfactual graphs to the original ones, i.e., the property (3) *Proximity*. Thirdly, we present the connection loss to control the new links between the motif and the contextual subgraph which, in turn, can control the number of perturbations, i.e., meet the property (4) Sparsity. Through the adversarial training with these losses, the model can generate counterfactual graphs that conform to the counterfactual properties. Finally, these generated counterfactual graphs and realistic graphs are adopted to train a robust anomaly detection model. The contributions of this study are threefold:

- We propose a new framework for graph-level anomaly detection, MotifCAR, which introduces a novel paradigm for generating counterfactual graphs using the motif.
- We present the GAN-based graph optimizer, that develops the three tailored losses to meet the requirement of desired counterfactual properties and generate high-quality samples.
- Extensive experiments demonstrate that MotifCAR significantly improves the detection performance and counterfactual data quality compared to state-of-the-art baselines.

2 Preliminaries

We now provide a preliminary overview of the graph-level anomaly detection problem, the graphon, and the discriminative motif, which will serve as important foundational elements in our model.

2.1 Problem Formulation

In general, a graph can be represented as $G(\mathcal{V}, \mathcal{E})$, where \mathcal{V} and \mathcal{E} denote its node set and edge set, respectively. We adopt $G, H, I/\mathcal{G}, \mathcal{H}, \mathcal{I}$ to denote graphs/graph set. $\mathbf{y}_G \in \mathbb{R}^C$ denotes the label of graph G,

where C is number of classes of graphs. Accordingly, the problem of graph-level anomaly detection is defined as follows:

PROBLEM 1. Graph-level anomaly detection aims to identify individual graphs within a given set G that exhibit anomalous behavior. In this study, the problem of graph-level anomaly detection can be regarded as the task of identifying anomalous graphs G_i based on a limited set of graph labels $y_i = \{0, 1\}$, where 0 signifies that the graph is considered an anomaly.

2.2 Graphon

We will leverage grahpons to produce counterfactual graphs. Here we present an introduction to the graphon and estimation method. **Graphon**. A graphon [1] is a continuous, bounded and symmetric function $W: [0,1]^2 \to [0,1]$ which may be thought of as the weight matrix of a graph with infinite number of nodes. Then, given two points $v_i, v_j \in [0,1]$, W(i,j) represents the probability that nodes i and j are related with an edge. For a set of graphs $\mathcal G$ with a given category label, a graphon can be estimated based on these graphs. Conversely, arbitrarily sized graphs can be sampled from this graphon, and sampled graphs will preserve the same category label with this graphon [15].

Graphon Estimation. Give a graph set $\mathcal{G} = \{G_1, G_2, \dots, G_n\}$, graphon estimation aims to deduce a graphon W_G based on \mathcal{G} . It is intractable because a graphon is an unknown function without a closed-form expression for real-world graphs [15]. Hence, the step function estimation methods are generally adopted to approximate graphons [15, 42]. The typical step function estimation method is composed of two stages: aligning the nodes in a set of graphs and estimating the step function from all the aligned adjacency matrices [15]. For the node alignment, this method first aligns multiple graphs based on node measurements (e.g., degree), and then selects the top K nodes from the aligned graphs and calculates the average matrix of their adjacency matrices. For function estimation, the goal is to obtain a matrix $\mathbf{W} = [w_{kk'}] \in [0, 1]^{K \times K}$, where w_{ii} denotes the probability of an edge existing between node i and node j and K is the number of the selected nodes, with the default being the average number of nodes in all graphs[15, 42]. In particular, the step function $\mathbf{W}^P : [0,1]^2 \mapsto [0,1]$ is defined as follows:

$$\mathbf{W}^{P}(x,y) = \sum_{k,k'=1}^{K} w_{kk'} \mathbb{1}_{P_k \times P_{k'}}(x,y), \tag{1}$$

 $\mathcal{P} = (P_1, ..., P_K)$ represents the partition of [0, 1] into K adjacent intervals of length 1/K, $w_{kk'} \in [0, 1]$, and the indicator function $\mathbb{1}_{P_k \times P_{k'}}(x, y)$ equals 1 if $(x, y) \in P_k \times P_{k'}$ and 0 otherwise.

2.3 Discriminative Motif

Motif, also called interest or frequent subgraph, refers to a specific subgraph structure that frequently occurs in a graph such as edges, triangles and quadrilaterals [47]. On the basis of the motif concept, similar to the work [15], we define the discriminative motif as:

DEFINITION 1. A discriminative motif F_G of the graph G is the subgraph, which can decide the category of the graph G.

Intuitively, the discriminative motif is the core subgraph of a graph and every graph has a discriminative motif, which generally

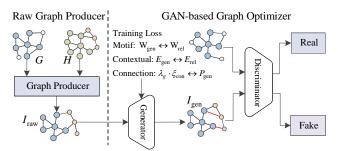


Figure 2: The Overview of the MotifCAR framework.

has a smaller number of nodes and edges. For the relationship of discriminative motifs and graphons, we have the theorem:

THEOREM 1. For a graph G which is sampled from the graphon W_G , its discriminative motif exists in W_G .

Proof Sketch. We verify this by stating that the homomorphism density of the discriminative motif in generated (sampled) graphs will be approximately equal to that in the graphon with high probability. In other words, the sampled graphs will preserve the discriminative motif of the graphon with a very high probability. Here, homomorphism density is defined to measure the relative frequency that the graph H appears in graph G [15, 26].

3 Method

As mentioned in Section 1, MotifCAR is composed of two principal components: the raw counterfactual graph producer and the GAN-based graph optimizer, which are illustrated in Figure 2. In the rest of this section, we discuss the technical details of MotifCAR.

3.1 Producing Raw Counterfactual Graph

Here we present the raw graph producer, which takes two graphs as inputs and extracts the motif subgraph from one graph and contextual subgraph from another to generate a raw counterfactual graph. The generated graph possesses the label of the first graph while involving the environment information of the second graph. To this end, this producer first merges the two graphs into one graph, and then builds a masked matrix based on graphons of the categories of these two graphs to filter the motif and the contextual subgraph. Finally, the raw counterfactual graphs are produced by combining the adjacency matrix of the merged graph and the masked matrix. We adopt this generation way because: (1) this fixed merging pattern ensures a stable outcome for each merging attempt, and (2) it enables uniform expansion of the node counts when the two graphs have different numbers of nodes.

Formally, suppose that a graph G with n_g nodes is from the graph set \mathcal{G} and graph H with n_h nodes from the graph set \mathcal{H} . Here G and H can have the same/different cluster or category label and their graphons can be computed using Equation 1, denoted as W_G and W_H , individually. To produce raw counterfactual graphs, we first merge two graphs G and H into one. Since G and H might have different numbers of nodes, we extend their node sets to a set $V_G \cup V_H$, and their adjacency matrices can be computed as:

$$A_G^{\text{ext}} = \begin{bmatrix} A_G & 0 \\ 0 & 0_H \end{bmatrix}, A_H^{\text{ext}} = \begin{bmatrix} 0_G & 0 \\ 0 & A_H \end{bmatrix}, \tag{2}$$

where 0_G and 0_H are zero matrices with shapes $n_g \times n_g$ and $n_h \times n_h$, respectively. Then, the two adjacency matrices are merged into A_B :

$$A_p = A_G^{\text{ext}} + A_H^{\text{ext}} = \begin{bmatrix} A_G & A_C \\ A_C^T & A_H \end{bmatrix}, \tag{3}$$

where A_C is a matrix indicating the cross-graph connectivity between the nodes in G and H. A_G and A_H are the aligned adjacency matrices, i.e., their first $|W_G|$ and $|W_H|$ elements are aligned with their corresponding graphon nodes. We randomly sample η edges to connect G and H to ensure the produced graph is connected.

To remove the contextual (non-motif) nodes in G and remove the motif nodes in H, we build a mask matrix which is used to perform an XNOR operation with A_p to eliminate these nodes. Since the motifs of G and H exist in the graphons W_G and W_H , individually, we build this mask matrix based on W_G and W_H . To this end, we first extend row and column numbers of W_G and W_H to be the same with A_G and A_H , respectively:

$$W_G^{\text{ext}} = \begin{bmatrix} W_G & 0 \\ 0 & 0_{\overline{G}} \end{bmatrix}, W_H^{\text{ext}} = \begin{bmatrix} W_H & 0 \\ 0 & 0_{\overline{H}} \end{bmatrix}, \tag{4}$$

where $0_{\overline{G}}$ and $0_{\overline{H}}$ are zero matrices with shapes $n_{\overline{g}} \times n_{\overline{g}}$ and $n_{\overline{h}} \times n_{\overline{h}}$, individually. Here $n_{\overline{g}} = \big| |A_G| - |W_G| \big|$ and $n_{\overline{h}} = \big| |A_H| - |W_H| \big|$. Then, we build this mask matrix as:

$$W^{m} = \begin{vmatrix} \hat{W}_{G}^{\text{ext}} & A_{C} \\ A_{C}^{T} & A_{H} - \hat{W}_{H}^{\text{ext}} \end{vmatrix}, \tag{5}$$

where A_C is the same with that in Equation 3. \hat{W}_G^{ext} and \hat{W}_H^{ext} are the matrices obtained after binary conversion of W_G^{ext} and W_H^{ext} :

$$\hat{W}_G^{\text{ext}} = Bern(W_G^{\text{ext}}), \hat{W}_H^{\text{ext}} = Bern(W_H^{\text{ext}}),$$
 (6)

where $Bern(\cdot)$ refers to the Bernoulli function which binarizes the values of the matrix. As a result, we conduct the XNOR operation between W^m and A_p to obtain the adjacency matrix of the counterfactual graph:

$$A_p^m = W^m \odot A_p. (7)$$

According to the adjacency matrix A_p^m , we can obtain the raw counterfactual graph $I_{\text{raw}}(\mathcal{V}, \mathcal{E})$, where \mathcal{E} is derived from A_p^m and \mathcal{V} is composed of the motif nodes of G and contextual nodes of H.

Since $I_{\text{raw}}(\mathcal{V}, \mathcal{E})$ is produced based on the motif of one graph and the contextual subgraph of another graph, it should possess the identification of the first graph and also involve the environment characteristics of the second graph. However, the motif nodes and the contextual nodes are connected randomly, which may result in distorted links. Hence, we further feed the produced raw counterfactual graphs into the graph optimizer to refine the graphs.

3.2 Optimizing Counterfactual Graphs

In the raw counterfactual graph $I_{\rm raw}(\mathcal{V},\mathcal{E})$, the edges between the motif nodes and the contextual nodes (i.e., A_C in Equation 3) are randomly assigned, which might not conform to the distribution of realistic graphs. To obtain high-quality graphs, we here present a GAN-based graph optimizer to refine the raw graphs to meet the four counterfactual properties (*Realism*, *Validity*, *Proximity* and *Sparsity*). This optimizer is composed of two essential modules: the graph generator and the graph discriminator. The generator aims to

generate refined graphs through adjusting edges. The discriminator is designed to distinguish between the graphs generated by the generator and realistic graphs. The graph generator and the graph discriminator are trained with our designed losses in an adversarial style to generate high-quality counterfactual graphs. These generated graphs will be used to train a robust classification model for better anomaly detection.

3.2.1 Graph Generator. For the raw graph $I_{\text{raw}}(\mathcal{V}, \mathcal{E})$, the graph generator is designed to refine the edges in order to make the generated graph more closely resemble the real graph. We assume that each edge $\mathcal{E}(v_i, v_j)$ in I_{raw} is associated with a random variable $P_{i,j} \sim Bern(\mathcal{W})$, where $\mathcal{W} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ is a learnable matrix, P is a binary matrix with size $|\mathcal{V}| \times |\mathcal{V}|$, $\mathcal{E}(v_i, v_j)$ is in I_{raw} if $P_{i,j} = 1$ and is dropped otherwise.

Inspired by the work [39], we relax the discrete $P_{i,j}$ to a continuous variable with values in the range (0, 1) to facilitate end-to-end training of the generator:

$$P = \sigma(\frac{W - X_g}{\tau_a}),\tag{8}$$

where $X_g \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ represents a random matrix with each element sampled from a uniform distribution: $X_{i,j} \sim U(0,1)$. $\sigma(x) = \frac{1}{1+e^{-x}}$ is the Sigmoid function, and $\tau_g \in (0,1]$ is a hyperparameter used to make $P_{i,j}$ approach either 0 or 1. Here, P can be regarded as an approximation of the generated adjacency matrix.

Motif Consistency Loss. The purpose of this loss is to ensure that the identification information of the generated graph is consistent with that of the given realistic graph. Since the graphon contains the motif of the graph and the motif decides the identification information, we adopt the graphon to build this consistency loss, i.e., we force the graphon of the generated graphs to be close to that of the realistic graphs. For the generated graph set $\mathcal{I}_{\text{gen}} = \{I_{\text{gen}}^1, I_{\text{gen}}^2, \cdots, I_{\text{gen}}^n\}$ and the realistic graph set $\mathcal{G}_{\text{rel}} : G_{\text{rel}}^1, G_{\text{rel}}^2, \cdots, G_{\text{rel}}^n\}$, their graphons W_{gen} and W_{rel} can be computed by Equation 1. Then, the loss is defined as:

$$\mathcal{L}_{\text{motif}} = \sum_{i,j} ||W_{\text{gen}}(i,j) - W_{\text{rel}}(i,j)||_F^2.$$
(9)

This loss forces the generated graph to have the same motif as the realistic graphs. Since the motif decides the category label, this loss encourages the generated graphs to meet the property *Validity*.

Contextual loss. This loss is proposed to force the contextual subgraph to be similar to the realistic contextual one. Since the generator mainly adjusts the edges, we consider the edge-related feature – the degree distribution – for this loss, i.e., we force the degree distributions of the generated and realistic contextual subgraphs to be similar. We introduce the degree entropy E to measure the degree distribution of the subgraph. Supposing that d(k) represents the degree number of node v_k in a generated contextual subgraph, the degree entropy is computed as:

$$E_{\text{gen}} = -\frac{1}{\ln(n_c)} \cdot \sum_{k=0}^{n_c} \frac{d(k)}{\sum_{m=0}^{n_c} d(m)} \ln\left(\frac{d(k)}{\sum_{m=0}^{n_c} d(m)}\right), \tag{10}$$

where n_c refers to the number of nodes in this contextual subgraph. This entropy indicates the average level of the node degree. Both

entropies of the generated and realistic contextual subgraphs should be similar. Hence, the contextual loss is defined as:

$$\mathcal{L}_{\text{context}} = \sum_{i=1}^{n} \left| E_{\text{gen}}^{i} - E_{\text{rel}}^{i} \right|, \tag{11}$$

where $E_{\rm rel}$ is the entropy of the realistic contextual subgraph and n denotes the number of the subgraphs. Both the contextual loss and the motif consistency loss compel the generated contextual subgraph and the motif to approach the realistic ones, which can meet the property Proximity.

Connection loss. This objective is designed to control the number of edges connecting the motif nodes and the contextual nodes in the generated graph. Inspired by the edge dropout [55], we set a ratio λ_g and train $\mathcal W$ to generate graphs with $\lambda_g \cdot |\mathcal E_{\rm con}|$ connection edges. Here $\mathcal E_{\rm con}$ is the edge set connecting the motif nodes and the contextual nodes in the realistic graph. For a batch of generated graphs $\mathcal I_{\rm gen} = \{I_{\rm gen}^1, I_{\rm gen}^2, ..., I_{\rm gen}^n\}$, we calculate their average as the connection loss. Formally, this loss is computed as:

$$\mathcal{L}_{\text{con}} = \frac{1}{n} \sum_{k=1}^{n} \left| \lambda_g \cdot |\mathcal{E}_{\text{con}}^k| - P_{\text{gen}}^k \right|, \tag{12}$$

where the P_{gen} represents the sum of P in the connection edge set of the generated graph, which is computed using Equation 8. This loss tries to limit the new connection edges, which can control the property *Sparsity* of the generated graphs.

As a result, the regularization loss is the combination of the above losses:

$$\mathcal{L}_{\text{reg}} = \lambda_1 \cdot \mathcal{L}_{\text{motif}} + \lambda_2 \cdot \mathcal{L}_{\text{context}} + \lambda_3 \cdot \mathcal{L}_{\text{con}}, \tag{13}$$

where λ_1, λ_2 and λ_3 are hyper-parameters to balance the influences of these losses.

Moreover, for the sake of efficiency, we initialize \mathcal{W} rather than training from scratch. To be specific, we introduce an initialization rate $\gamma \in [0,1]$ to constrain the number of the connection edges at the beginning. Also the similarity of nodes is considered for the initialized values. Consequently, \mathcal{W} is initialized as follows:

$$\mathcal{W}_{i,j} = \begin{cases} 1 &, & \text{if } (v_i, v_j) \in \mathcal{E}_t; \\ \frac{\gamma \cdot \lambda_g \cdot |\mathcal{E}_{\text{con}}| \cdot S_{i,j}}{|C|} &, & \text{if } (v_i, v_j) \in \mathcal{C}; \\ 0 &, & \text{otherwise,} \end{cases}$$
(14)

where \mathcal{E}_t is the edge set in the motif and the contextual subgraph, \mathcal{E}_{con} and λ_g are the same to Equation 12, $S_{i,j}$ is the node similarity which can be computed based on node embedding or node attributes, and C is a candidate set of new connection edges.

3.2.2 Graph Discriminator. The discriminator is a graph-level classifier designed to distinguish between real and generated graphs. This can help the generator produce samples that are close to realistic data, i.e., meet the property Realism. Specifically, this discriminator takes a graph as input and determines whether it is real or fake. Supposing that $\mathcal{G}_{\rm rel}$ and $\mathcal{I}_{\rm gen}$ denote the realistic graph set and the generated graph set, individually, for each $G \in \mathcal{G}_{\rm rel} \cup \mathcal{I}_{\rm gen}$, we utilize a GNN encoder f to encode the representations of each node:

$$\{\mathbf{z}_v|v\in\mathcal{V}\}=f(G). \tag{15}$$

Table 1: Statistics of Datasets.

Datasets	# graph	# class	avg.# V	avg.# E	anomaly(%)
IMDB-B	1,000	2	19.77	96.53	9.0%
IMDB-M	1,500	3	13.00	65.94	4.7%
REDDIT-B	2,000	2	429.63	497.75	9.0%
REDDIT-M	5,000	5	508.52	594.87	4.7%

Then, we calculate the graph representation by concatenating the mean pooling and the maximum pooling of the node representations:

$$\mathbf{z}_{G} = \left(\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbf{z}_{v}\right) \oplus \text{MaxPool}\left(\left\{\mathbf{z}_{v} \middle| v \in \mathcal{V}\right\}\right),\tag{16}$$

where \oplus is the concatenate operation. With the graph representation, we compute the probability of G using an L_h -layer Multilayer Perceptron (MLP) h:

$$p_G = h(\mathbf{z}_G). \tag{17}$$

To train the discriminator, we label the graphs in \mathcal{G}_{rel} with 1 and the graphs in \mathcal{I}_{gen} with 0. Suppose that the label of G is l_G . Then, the classification loss is defined as follows:

$$\mathcal{L}_{\text{dis}} = -l_G \cdot \log(p_G) - (1 - l_G) \cdot \log(1 - p_G). \tag{18}$$

3.2.3 Model Training. We here present the model training procedure of the GAN-based graph optimizer. The graph generator and the discriminator are optimized sequentially and iteratively.

For the generator, in each iteration, an augmented graph $I_{\rm gen}$ is generated and then the regularization loss (cf. Equation 13) is computed. In consideration of generating high-quality graphs, an adversarial classification loss is incorporated to cheat the graph discriminator by labeling $I_{\rm gen}$ with 1. According to the first term in Equation 18 and the regularization loss (cf. Equation 13), we have the final loss of the generator:

$$\mathcal{L}_{\text{gen}} = -\log(p_{I_{\text{gen}}}) - \mathcal{L}_{req}. \tag{19}$$

For the discriminator, their parameters are optimized by classifying the realistic graphs and the generated graphs using the loss in Equation 18. For the computational complexity, since our MotifCAR is designed based on GANs, it has the similar complexity and overhead with GAN-based baselines (e.g., CFAD [41]).

3.3 Graph-level Anomaly Detection

After training, the model is used to generate a number of counterfactual graphs. Both the generated counterfactual graphs \mathcal{I}_{gen} and realistic graphs \mathcal{G}_{rel} are adopted to train a robust classifier to identify anomalies from normal graphs. We adopt the discriminator as the classifier after changing its loss. Specifically, for each $G \in \mathcal{G}_{\text{rel}} \cup \mathcal{I}_{\text{gen}}$, we first adopt Equation 17 to obtain its graph representation. Then we label the normal graph with 1 and the anomalous graph with 0. Supposing that y_G denotes the label, the loss becomes:

$$\mathcal{L}_{class} = -y_G \cdot \log(p_G) - (1 - y_G) \cdot \log(1 - p_G). \tag{20}$$

Moreover, to enhance detection performance, we refine the graph representations following the work [31] during the process of the model training.

Method IMDB-BINARY			Y	IMDB-MULTI		REDDIT-BINARY			REDDIT-MULTI-5K			
	PRECISION	RECALL	F1	PRECISION	RECALL	F1	PRECISION	RECALL	F1	PRECISION	RECALL	F1
g-U-Nets	0.87±0.01	0.91±0.01	0.88±0.01	0.80±0.01	0.86±0.01	0.80±0.01	0.86±0.01	0.75±0.02	0.78±0.01	0.76±0.01	0.71±0.02	0.73±0.01
DiffPool	0.66±0.01	0.69 ± 0.01	0.81 ± 0.02	0.65 ± 0.01	0.70 ± 0.01	0.83 ± 0.02	0.71 ± 0.01	0.63 ± 0.01	0.64 ± 0.01	0.71 ± 0.01	0.61±0.01	0.71 ± 0.01
SAGPool	0.80±0.01	0.88±0.01	0.84 ± 0.01	0.73±0.01	0.90±0.01	0.82±0.01	0.81±0.01	0.81±0.01	0.83 ± 0.01	0.76 ± 0.01	0.61±0.01	0.74 ± 0.01
GMT	0.83 ± 0.03	0.88 ± 0.03	0.88 ± 0.02	0.82 ± 0.02	0.84 ± 0.02	0.83 ± 0.02	0.83 ± 0.02	0.83 ± 0.01	0.83 ± 0.02	0.78 ± 0.02	0.80 ± 0.01	0.67 ± 0.02
CFGL-LCR	0.85±0.03	0.87±0.03	0.84±0.03	0.83±0.04	0.83±0.02	0.76±0.02	086±0.02	0.88±0.02	0.82±0.03	0.82±0.03	0.80±0.03	0.78±0.03
CFAD	0.87±0.03	0.88±0.02	0.86 ± 0.03	0.86±0.03	0.88±0.02	0.78±0.03	0.87±0.03	0.86 ± 0.02	0.84 ± 0.03	0.83 ± 0.02	0.81±0.03	0.82±0.02
CGC	0.92±0.02	0.94±0.03	0.92±0.02	0.89 ± 0.02	0.93±0.02	0.92±0.02	0.90±0.03	0.93±0.02	0.91±0.02	0.84 ± 0.03	0.95±0.02	0.84 ± 0.02
CF-HGExp	0.93 ± 0.01	0.94 ± 0.01	0.94±0.01	0.90 ± 0.01	0.95 ± 0.01	0.94±0.01	0.92 ± 0.01	0.94±0.01	0.93±0.01	0.85 ± 0.01	0.95 ± 0.01	0.85 ± 0.01
OCGTL	0.88±0.01	0.89±0.01	0.87±0.02	0.82±0.01	0.84±0.01	0.80±0.02	0.89±0.01	0.89±0.01	0.88±0.01	0.86±0.01	0.84±0.01	0.82±0.01
GLocalKD	0.65±0.01	0.66±0.01	0.66±0.01	0.60±0.01	0.63±0.01	0.61±0.01	0.63±0.01	0.64±0.01	0.63±0.01	0.58±0.01	0.62±0.01	0.73±0.01
iGAD	0.88±0.03	0.78±0.02	0.87±0.02	0.86±0.02	0.75±0.02	0.82±0.02	0.57±0.02	0.58±0.03	0.55±0.02	0.53 ± 0.02	0.54±0.03	0.67±0.02
GmapAD	0.92±0.02	0.94±0.01	0.94 ± 0.01	0.90 ± 0.01	0.95±0.01	0.93±0.01	0.91 ± 0.02	0.96±0.01	0.95 ± 0.01	0.87 ± 0.02	0.95±0.01	0.86 ± 0.01
MotifCAR	0.94±0.01	0.95±0.01	0.96±0.01	0.91±0.01	0.96±0.01	0.94±0.01	0.93±0.01	0.96±0.01	0.96±0.01	0.89±0.01	0.97±0.01	0.87±0.01

Table 2: Performance comparison between MotifCAR and baselines on four datasets.

4 Expriments

In this section, we evaluate the MotifCAR performance, including the effectiveness, ablation study, hyper-parameter analysis and the quality of the generated counterfactual graphs.

4.1 Experiments Setup.

Datasets. We adopt four public datasets for our experiments, whose statistics are presented in Table 1. *IMDB-BINARY (IMDB-B)* and *IMDB-MULTI (IMDB-M)* are movie collaboration datasets. Each graph corresponds to an ego-network for each actor/actress, where nodes correspond to actors/actresses and an edge is drawn between two actors/actresses if they appear in the same movie. Each graph is derived from a pre-specified genre of movies, which is regarded as the category label. *REDDIT-BINARY (REDDIT-B)* and *REDDIT-MULTI-5K (REDDIT-M)* are balanced datasets where each graph corresponds to an online discussion thread and nodes correspond to users. An edge is drawn between two nodes if at least one of them responds to another's comment. The community or subreddit is considered as the label of the graph. Following previous works [31, 36], we downsample one of the categories as the anomalous one, and others as the normal data.

Baselines. We conduct a comparison between our developed framework MotifCAR and three categories of baselines: (1) The state-of-the-art GNN models (*g-U-Nets* [13], *SAGPool* [20], *DIFFPOOL* [45], and *GMT* [3]) leverage different pooling strategies and specially designed pooling layers for learning the graph-level representations. (2) The counterfactual graph augmentation methods (*CFGL-LCR* [49], *CFAD* [41], *CGC* [43], and *CF-HGExp* [44]) generate counterfactual graphs to enhance the node or graph representations and further improve the classification performance. (3) The graph-level anomaly detection approaches (*OCGTL* [36], *GLocalKD* [29], *iGAD* [48] and *GmapAD* [31]) explore the tailored classification models or graph-level anomaly patterns to the anomalies.

Experiments Settings. We use the Adam algorithm to optimize the model with learning rate 0.001. For hyper-parameters, we set τ_g = 0.0001 in Equation 8, γ = 0.75 in Equation 14, λ_1 = 1, λ_2 = 0.9, λ_3 = 0.6 in Equation 5, λ_g = 0.5 for IMDB-M and REDDIT-M and λ_g = 0.8 for IMDB-B and REDDIT-B in Equation 12. We split the dataset into train/validation/test data by 2 : 4 : 4. The best test epoch is selected

on the validation set, and we report the test accuracy on ten runs. All the experiments are conducted on an Ubuntu 20.04 server with a 12-core CPU, 1 Nvidia RTX 3090 GPU and 64Gb RAM.

4.2 Anomaly Detection Results

We report the anomaly detection performance of MotifCAR and the baselines in Table 2 and have the following observations. First, in terms of overall detection results, our model MotifCAR consistently outperforms all baseline models across the four datasets. In particular, our method achieves more gains on IMDB-Binary and REDDIT-Binary. This is because graphon estimation can perform better on dense graphs [25] and the graphs in these two datasets are denser. Correspondingly, our model exhibits more significant advantages on these two datasets. This result verifies that our model can generate effective counterfactual graphs and enhance detection performance. Besides, the t-test indicates that MotifCAR's performance is statistically significant compared to the baselines.

Second, some of the counterfactual graph augmentation methods, such as CGC and CF-HGExp, acquire relatively good performance, compared with the GNN models. This indicates that the appropriate counterfactual data can benefit graph-level anomaly detection. However, our model outperforms these state-of-the-art models, which validates the efficacy of the proposed GAN-based graph optimizer with the tailored losses for graph-level anomaly detection.

Third, the method GmapAD, which is specially designed for graph-level anomaly detection, exhibits remarkable performance. It can explore both the intra- and inter-graph node information to enhance the graph representations and detection performance. However, when faced with a smaller amount of training data, the trained model cannot generalize well to the test data, leading to a performance decline. However, our model can generate counterfactual data to handle varying environments and hence, outperform this method.

4.3 Ablation Study

In this section, we investigate the contributions of the key components in MotifCAR. We mainly observe the following variants: (1) In *MotifCAR-Opti*, we remove the GAN-based graph optimizer.

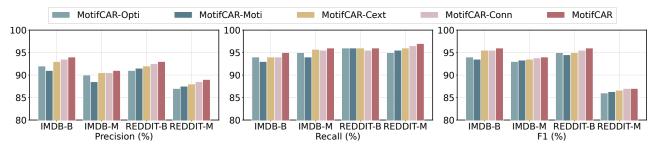


Figure 3: Ablation study: Variants of MotifCAR.

(2) In *MotifCAR-Moti*, we remove the motif consistency loss in the GAN-based optimizer. (3) In *MotifCAR-Cext*, we remove the contextual loss in the optimizer. (4) In *MotifCAR-Conn*, we remove the connection loss in the optimizer.

As shown in Figure 3, no matter which part we remove in MotifCAR, the model's performance degrades. This suggests the efficacy of our designed GAN-based optimizer as well as the losses in this optimizer. Moreover, when removing the motif consistency loss, the performance of *MotifCAR-Moti* delines more dramatically than *MotifCAR-Cext* and *MotifCAR-Conn*. This indicates that the motif consistency loss is more important when refining the counterfactual graphs. The reason is that this loss can ensure the generated graphs to have the corresponding identification information, which is vital for training a robust classification model.

4.4 Qualitative Analysis of the Counterfactuals

We here investigate the quality of the generated counterfactual graphs from the aspects of the counterfactual properties: *Realism*, *Proximity, Validity*, and *Sparsity*. According to the works [18, 44], the *Realism* score reflects the change of the shift detection result between the realistic graphs and the counterfactual graphs. The *Proximity* score estimates the mean of feature distances between the counterfactual graph and the realistic graph. The *Validity* score gauges the fraction of the generated counterfactual examples that are correctly predicted by the classifier to the corresponding class. The *Sparsity* score measures the difference between the edges of the counterfactual graph and the realistic graph. We compare MotifCAR with counterfactual data generation baselines: *CFGL-LCR* (CFGL) [49], *CFAD* [41], *CGC*[43], and *CF-HGExp* (CF-HG) [44]).

Table 3 demonstrates the *Realism* and *Proximity* scores. For the Realism score (the lower the better), our method achieves the best result on the four datasets. The reason is that adversarial training is beneficial to generate counterfactual graphs that are close to realistic ones. Also, the motif consistency and contextual loss also facilitate aligning the motif and contextual subgraph well. For the Proximity score (the lower the better), our model MotifCAR does not obtain the best result. This is because MotifCAR replaces the contextual subgraph in the counterfactual graphs, which results in a few differences in the feature space. While our result is competitive to the best baseline, CF-HGExp. In fact, both of their scores are very close. However, our model can achieve better detection performance than CF-HGExp.

Table 3: Realism and Proximity.

		CFGL	CFAD	CGC	CF-HG	MotifCAR
Realism	IMDB-B	1.980	1.84	0.844	0.688	0.622
	IMDB-M	1.946	1.850	0.752	0.562	0.514
	REDDIT-B	1.910	1.844	1.066	0.850	0.784
	REDDIT-M	1.896	1.860	1.142	0.754	0.462
Proximity	IMDB-B	0.532	0.519	0.179	0.098	0.076
	IMDB-M	0.548	0.558	0.118	0.075	0.078
	REDDIT-B	0.613	0.578	0.189	0.057	0.065
	REDDIT-M	0.656	0.623	0.185	0.129	0.149

The *Validity* and *Sparsity* scores are presented in Figure 4, where the misalignment of points on the x-axis is caused by discontinuous Sparsity values. As shown, our proposed MotifCAR achieves the best Validity performance at all levels of sparsity on the four datasets. This is because the raw counterfactual graphs contain the motif, the core subgraph of the graph, which can decide their category. Further, during the graph refinement procedure, the motif consistency provides strong consistent power for the counterfactual graphs keeping the identification information. As a result, the generated counterfactual graphs can possess the identification information and obtain higher Validity scores.

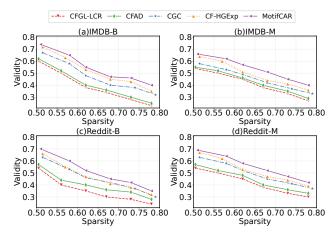


Figure 4: Validity and Sparsity.

4.5 Hyper-Parameter Sensitivity

In this section, we analyze the sensitivity of the hyper-parameters, λ_g , λ_1 , λ_2 , and λ_3 in Equation 12 and 13. The results are reported in Figure 5. For λ_g , the performance increases with the rise of the λ_g value. A smaller value leads to a lower F1-score. This is because λ_g determines the number of edges between the motif and the contextual subgraph, and a smaller value might cause the two subgraphs to disconnect. Correspondingly, the generated counterfactual graphs mainly contain the motif, which has little effect on improving the model's generalization ability.

For λ_1 , λ_2 , and λ_3 , these three parameters determine the weights of the motif consistency loss, the contextual loss, and the connection loss, individually. These figures show that F1-score increases with the rise of λ_1 , which suggests that more attention should be paid to the motif consistency loss. The reason is that this loss determines the identification information of the generated counterfactual samples, which is quite important for training a robust model. Besides, the detection performance increases with the rise of λ_2 and λ_3 , and then starts to decrease when $\lambda_2 > 0.9$ and $\lambda_3 > 0.6$. When they are too small, the contextual subgraph cannot effectively be incorporated into the generated counterfactual graphs, resulting in a lower performance. On the contrary, the higher λ_2 and λ_3 might attenuate the effectiveness of other losses, leading to inferior performance.

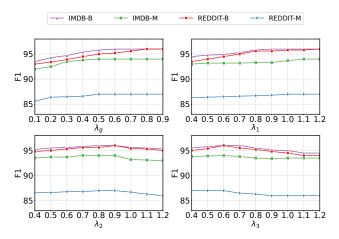


Figure 5: The effect of the hyper-parameters.

5 Related Work

Graph-level Anomaly Detection. Graph-level Anomaly Detection focuses on detecting entire abnormal graphs, rather than localizing anomalies in graphs. Most graph anomaly detection work is devoted to detecting irregular nodes and edges within graphs [12, 21, 30, 38]. Recently, graph-level anomaly detection has started to receive in-depth exploration. Initially, researchers adopt shallow learning techniques to detect graph-level anomalies. They mainly exploit graph kernels and graph signals to detect graph-level anomalies. For example, graph kernels (e.g., Weisfeiler-Leman Kernel and propagation kernels) are explored to measure the pairwise similarity of nodes in the graph, and the similarity score is used to identify anomalies based on the structural characteristics of the graph [32]. Changedar [16] outlines anomaly detection using graph signals generated by anomalous node sets. By examining patterns

of signal changes, this approach can effectively pinpoint anomalies within the graph. Also, a number of works [17, 35, 47] leverage frequent graph motifs to model the network topology among nodes and motifs. These frequently occurring subgraphs capture crucial high-order structural information, which enables models to detect anomalies more comprehensively. However, the shallow learning-based methods may achieve sub-optimal performance due to the low coupling of the detector and graph representation learning [34].

The deep learning-based methods, in contrast, are end-to-end strategies that have been successfully applied to both static and dynamic graphs for anomaly detection [10, 46, 54]. The remarkable development of GNNs [6, 7] has led to great progress in the field of graph-level anomaly detection. One approach involves utilizing GNNs in conjunction with classification loss functions to train a graph-level anomaly detection framework, exemplified by the works such as OCGIN [51] and OCGTL [36]. However, these classification methods do not cope well with imbalanced datasets, resulting in an underfit for anomalous graphs. iGAD [48] tackles this problem by modeling Point Mutual Information. And the out-ofdistribution problem of graph data is addressed for better anomaly detection [22, 24]. Another method centers around the identification of anomalies by scrutinizing the irregular attributes within each graph concerning the overall graph structure, as evidenced in GmapAD [31] and GLADST [23]. Moreover, GLocalKD [29] and GLADC [27] capture anomalies from global and local graph perspectives. Additionally, some approaches concurrently consider node-level anomalies and substructure anomalies in anomaly detection, such as iGAD [48], GLAM [52] and HO-GAT [17].

Counterfactual Graph Learning. The works about counterfactual graph learning primarily fall into two folds: counterfactual explanations and counterfactual data augmentation. The former aims to identify the necessary changes to the input graph that can alter the prediction outcome, which can help to filter out spurious explanations [14]. The related methods try to find a counterfactual graph by conducting minimal perturbations (i.e., adding or removing the minimum number of edges) that could lead to counterfactual predictions [4, 43, 44]. Counterfactual data augmentation is a promising technique used to augment training data to reduce model reliance on spurious correlations and aid in learning causal representations [8, 14, 50]. To alleviate the problem of spurious correlations and enhance models' generalization capacity, some researchers try to inject interventions (perturbations) on the node attributes and the graph structure to generate counterfactual data [28, 41]. While, others attend to match counterfactual examples which are the most similar items with different treatments [5, 53] to boost models' robustness.

6 Conclusion

In this paper, we proposed a novel framework, MotifCAR, for graph-level anomaly detection. We designed a counterfactual graph producer to produce raw counterfactual graphs by combining the discriminative motif and the contextual subgraph. It can generate high-quality counterfactual graphs and effectively alleviate the performance decline issue under varying environments. We also proposed a GAN-based graph optimizer to refine the raw graphs

through capturing the motif consistency and controlling the contextual subgraph structures. Extensive experiments demonstrate the superiority of our proposed approach over state-of-the-art baselines on graph benchmarks. Our future work will focus on further improving the efficiency of MotifCAR and extending it to more complicated scenarios such as evolving graph anomaly detection.

References

- Edo M Airoldi, Thiago B Costa, and Stanley H Chan. 2013. Stochastic blockmodel approximation of a graphon: Theory and consistent estimation. NIPS (2013).
- [2] Leman Akoglu, Hanghang Tong, and Danai Koutra. 2015. Graph based anomaly detection and description: a survey. Data mining and knowledge discovery 29 (2015), 626–688.
- [3] Jinheon Baek, Minki Kang, and Sung Ju Hwang. 2021. Accurate learning of graph representations with graph multiset pooling. In ICLR.
- [4] Mohit Bajaj, Lingyang Chu, Zi Yu Xue, Jian Pei, Lanjun Wang, Peter Cho-Ho Lam, and Yong Zhang. 2021. Robust counterfactual explanations on graph neural networks. NIPS (2021), 5644–5655.
- [5] Heng Chang, Jie Cai, and Jia Li. 2023. Knowledge Graph Completion with Counterfactual Augmentation. In WWW. 2611–2620.
- [6] Hao Chen, Yuanchen Bei, Qijie Shen, Yue Xu, Sheng Zhou, Wenbing Huang, Feiran Huang, Senzhang Wang, and Xiao Huang. 2024. Macro graph neural networks for online billion-scale recommender systems. In WWW. 3598–3608.
- [7] Hao Chen, Yue Xu, Feiran Huang, Zengde Deng, Wenbing Huang, Senzhang Wang, Peng He, and Zhoujun Li. 2020. Label-Aware Graph Convolutional Networks. In CIKM. 1977–1980.
- [8] Zeming Chen, Qiyue Gao, Antoine Bosselut, Ashish Sabharwal, and Kyle Richardson. 2023. Disco: distilling counterfactuals with large language models. In ACL. 5514–5528.
- [9] Yu-Liang Chou, Catarina Moreira, Peter Bruza, Chun Ouyang, and Joaquim Jorge. 2022. Counterfactuals and causability in explainable artificial intelligence: Theory, algorithms, and applications. *Information Fusion* 81 (2022), 59–83.
- [10] Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. 2019. Deep anomaly detection on attributed networks. In ICDM. 594–602.
- [11] Kaize Ding, Zhe Xu, Hanghang Tong, and Huan Liu. 2022. Data augmentation for deep graph learning: A survey. ACM SIGKDD Explorations Newsletter 24, 2 (2022), 61–77.
- [12] Jingcan Duan, Siwei Wang, Pei Zhang, En Zhu, Jingtao Hu, Hu Jin, Yue Liu, and Zhibin Dong. 2023. Graph anomaly detection via multi-scale contrastive learning networks with augmented view. In AAAI. 7459–7467.
- [13] Hongyang Gao and Shuiwang Ji. 2019. Graph u-nets. In ICML. 2083–2092.
- [14] Zhimeng Guo, Teng Xiao, Charu Aggarwal, Hui Liu, and Suhang Wang. 2023. Counterfactual Learning on Graphs: A Survey. arXiv:2304.01391 (2023).
- [15] Xiaotian Han, Zhimeng Jiang, Ninghao Liu, and Xia Hu. 2022. G-mixup: Graph data augmentation for graph classification. In ICML. 8230–8248.
- [16] Bryan Hooi, Leman Akoglu, Dhivya Eswaran, Amritanshu Pandey, Marko Jereminov, Larry Pileggi, and Christos Faloutsos. 2018. Changedar: Online localized change detection for sensor data on a graph. In CIKM. 507–516.
- [17] Ling Huang, Ye Zhu, Yuefang Gao, Tuo Liu, Chao Chang, Caixing Liu, Yong Tang, and Chang-Dong Wang. 2021. Hybrid-order anomaly detection on attributed networks. IEEE Transactions on Knowledge and Data Engineering (2021).
- [18] Saeed Khorram and Li Fuxin. 2022. Cycle-consistent counterfactuals by latent transformations. In CVPR. 10203–10212.
- [19] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. 2021. Out-ofdistribution generalization via risk extrapolation (rex). In ICML. 5815–5826.
- [20] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. 2019. Self-attention graph pooling. In ICML. 3734–3743.
- [21] Xuan Li, Chunjing Xiao, Ziliang Feng, Shikang Pang, Wenxin Tai, and Fan Zhou. 2023. Controlled graph neural networks with denoising diffusion for anomaly detection. Expert Systems with Applications (2023), 121533.
- [22] Zenan Li, Qitian Wu, Fan Nie, and Junchi Yan. 2022. Graphde: A generative framework for debiased learning and out-of-distribution detection on graphs. In NIPS. 30277–30290.
- [23] Fu Lin, Xuexiong Luo, Jia Wu, Jian Yang, Shan Xue, Zitong Wang, and Haonan Gong. 2023. Discriminative Graph-level Anomaly Detection via Dual-studentsteacher Model. In ADMA.
- [24] Yixin Liu, Kaize Ding, Huan Liu, and Shirui Pan. 2023. Good-d: On unsupervised graph out-of-distribution detection. In WSDM. 339–347.
- [25] László Lovász. 2012. Large networks and graph limits. Vol. 60. American Mathematical Society.
- [26] László Lovász and Balázs Szegedy. 2006. Limits of dense graph sequences. Journal of Combinatorial Theory, Series B 96, 6 (2006), 933–957.
- [27] X Luo, J Wu, J Yang, S Xue, H Peng, C Zhou, H Chen, Z Li, and QZ Sheng. 2022. Deep graph level anomaly detection with contrastive learning. Scientific Reports

- 12. 1 (2022), 19867-19867.
- [28] Jing Ma, Ruocheng Guo, Mengting Wan, Longqi Yang, Aidong Zhang, and Jundong Li. 2022. Learning fair node representations with graph counterfactual fairness. In WSDM. 695–703.
- [29] Rongrong Ma, Guansong Pang, Ling Chen, and Anton van den Hengel. 2022. Deep graph-level anomaly detection by glocal knowledge distillation. In WSDM. 704–714.
- [30] Xiaoxiao Ma, Jia Wu, Shan Xue, Jian Yang, Chuan Zhou, Quan Z Sheng, Hui Xiong, and Leman Akoglu. 2021. A comprehensive survey on graph anomaly detection with deep learning. IEEE Transactions on Knowledge and Data Engineering (2021).
- [31] Xiaoxiao Ma, Jia Wu, Jian Yang, and Quan Z Sheng. 2023. Towards graph-level anomaly detection via deep evolutionary mapping. In KDD. 1631–1642.
- [32] Emaad Manzoor, Sadegh M Milajerdi, and Leman Akoglu. 2016. Fast memoryefficient anomaly detection in streaming heterogeneous graphs. In KDD. 1035– 1044
- [33] Raha Moraffah, Mansooreh Karami, Ruocheng Guo, Adrienne Raglin, and Huan Liu. 2020. Causal interpretability for machine learning-problems, methods and evaluation. ACM SIGKDD Explorations Newsletter 22, 1 (2020), 18–33.
- [34] Chaoxi Niu, Guansong Pang, and Ling Chen. 2023. Graph-Level Anomaly Detection via Hierarchical Memory Networks. In ECML. 201–218.
- [35] Caleb C Noble and Diane J Cook. 2003. Graph-based anomaly detection. In KDD. 631–636.
- [36] Chen Qiu, Marius Kloft, Stephan Mandt, and Maja Rudolph. 2022. Raising the bar in graph-level anomaly detection. In IJCAI. 2196–2203.
- [37] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In ICLR.
- [38] Jianheng Tang, Jiajin Li, Ziqi Gao, and Jia Li. 2022. Rethinking graph neural networks for anomaly detection. In ICML. 21076–21089.
- [39] Cheng Wu, Chaokun Wang, Jingcao Xu, Ziyang Liu, Kai Zheng, Xiaowei Wang, Yang Song, and Kun Gai. 2023. Graph Contrastive Learning with Generative Adversarial Network. In KDD. 2721–2730.
- [40] Chunjing Xiao, Zehua Gou, Wenxin Tai, Kunpeng Zhang, and Fan Zhou. 2023. Imputation-based Time-Series Anomaly Detection with Conditional Weight-Incremental Diffusion Models. In KDD. 2742–2751.
- [41] Chunjing Xiao, Xovee Xu, Yue Lei, Kunpeng Zhang, Siyuan Liu, and Fan Zhou. 2023. Counterfactual Graph Learning for Anomaly Detection on Attributed Networks. *IEEE Transactions on Knowledge and Data Engineering* 35, 10 (2023), 10540–10553.
- [42] Hongteng Xu, Dixin Luo, Lawrence Carin, and Hongyuan Zha. 2021. Learning graphons via structured gromov-wasserstein barycenters. In AAAI. 10505–10513.
- [43] Haoran Yang, Hongxu Chen, Sixiao Zhang, Xiangguo Sun, Qian Li, Xiangyu Zhao, and Guandong Xu. 2023. Generating Counterfactual Hard Negative Samples for Graph Contrastive Learning. In WWW. 621–629.
- [44] Qiang Yang, Changsheng Ma, Qiannan Zhang, Xin Gao, Chuxu Zhang, and Xiangliang Zhang. 2023. Counterfactual Learning on Heterogeneous Graphs with Greedy Perturbation. In KDD. 2988–2998.
- [45] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. 2018. Hierarchical graph representation learning with differentiable pooling. NIPS (2018).
- [46] Minji Yoon, Bryan Hooi, Kijung Shin, and Christos Faloutsos. 2019. Fast and accurate anomaly detection in dynamic graphs with a two-pronged approach. In KDD. 647–657.
- [47] Zirui Yuan, Minglai Shao, and Qiben Yan. 2023. Motif-level Anomaly Detection in Dynamic Graphs. IEEE Transactions on Information Forensics and Security (2023).
- [48] Ge Zhang, Zhenyu Yang, Jia Wu, Jian Yang, Shan Xue, Hao Peng, Jianlin Su, Chuan Zhou, Quan Z Sheng, Leman Akoglu, et al. 2022. Dual-discriminative graph neural network for imbalanced graph-level anomaly detection. In NIPS. 24144–24157.
- [49] Kun Zhang, Chong Chen, Yuanzhuo Wang, Qi Tian, and Long Bai. 2023. CFGL-LCR: A Counterfactual Graph Learning Framework for Legal Case Retrieval. In KDD. 3332–3341.
- [50] Xiheng Zhang, Yongkang Wong, Xiaofei Wu, Juwei Lu, Mohan Kankanhalli, Xiangdong Li, and Weidong Geng. 2021. Learning causal representation for training cross-domain pose estimator via generative interventions. In ICCV. 11270–11280.
- [51] Lingxiao Zhao and Leman Akoglu. 2023. On using classification datasets to evaluate graph outlier detection: Peculiar observations and new insights. Big Data 11, 3 (2023), 151–180.
- [52] Lingxiao Zhao, Saurabh Sawlani, Arvind Srinivasan, and Leman Akoglu. 2022. Graph anomaly detection with unsupervised GNNs. In ICDM.
- [53] Tong Zhao, Gang Liu, Daheng Wang, Wenhao Yu, and Meng Jiang. 2022. Learning from counterfactual links for link prediction. In ICML. 26911–26926.
- [54] Li Zheng, Zhenpeng Li, Jian Li, Zhao Li, and Jun Gao. 2019. AddGraph: Anomaly Detection in Dynamic Graph Using Attention-based Temporal GCN. In IJCAI. 4419–4425.
- [55] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2020. Deep graph contrastive representation learning. arXiv:2006.04131 (2020).