# TOWARDS A UNIFIED FRAMEWORK OF CLUSTERING BASED ANOMALY DETECTION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Unsupervised Anomaly Detection (UAD) plays a crucial role in identifying abnormal patterns within data without labeled examples, holding significant practical implications across various domains. Although the individual contributions of representation learning and clustering to anomaly detection are well-established, their interdependencies remain under-explored due to the absence of a unified theoretical framework. Consequently, their collective potential to enhance anomaly detection performance remains largely untapped. To bridge this gap, in this paper, we propose a novel probabilistic mixture model for anomaly detection to establish a theoretical connection among representation learning, clustering, and anomaly detection. By maximizing a novel anomaly-aware data likelihood, representation learning and clustering can effectively reduce the adverse impact of anomalous data and collaboratively benefit anomaly detection. Meanwhile, a theoretically substantiated anomaly score is naturally derived from this framework. Lastly, drawing inspiration from gravitational analysis in physics, we have devised an improved anomaly score that more effectively harnesses the combined power of representation learning and clustering. Extensive experiments, involving 17 baseline methods across 30 diverse datasets, validate the effectiveness and generalization capability of the proposed method, surpassing state-of-the-art methods.

#### 028 029

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

031

#### 1 INTRODUCTION

Unsupervised Anomaly Detection (UAD) refers to the task dedicated to identifying abnormal patterns 033 or instances within data in the absence of labeled examples Chandola et al. (2009). It has long 034 received extensive attention in the past decades for its wide-ranging applications in numerous practical scenarios, including financial auditing Bakumenko & Elragal (2022), healthcare monitoring Salem 035 et al. (2014) and e-commerce sector Kou et al. (2004). Due to the lack of explicit label guidance, the key to UAD is to uncover the dominant patterns that widely exist in the dataset so that samples do not 037 conform to these patterns can be recognized as anomalies. To achieve this, early works Chalapathy & Chawla (2019) have heavily relied on powerful unsupervised *representation learning* methods to extract the normal patterns from high-dimensional and complex data such as images, text, and 040 graphs. More recent works Song et al. (2021); Aytekin et al. (2018) have utilized *clustering*, a 041 widely observed natural pattern in real-world data, to provide critical global information for anomaly 042 detection and achieved tremendous success. 043

While the individual contributions of representation learning and clustering to anomaly detection 044 are well-established, their interrelationships remain largely unexplored. Intuitively, discriminative representation learning can leverage accurate clustering results to differentiate samples from distinct 046 clusters in the embedding space (i.e., ①). Similarly, it can utilize accurate anomaly detection to 047 avoid preserving abnormal patterns (i.e., 2). For accurate clustering, it can gain advantages from 048 representation learning by operating in the discriminative embedding space (i.e., ③). Meanwhile, it can potentially benefit from accurate anomaly detection by excluding anomalies when formulating clusters (i.e., 4). Anomaly detection can greatly benefit from both discriminative representation 051 learning and accurate clustering (i.e., (5) & (6)). However, these benefits hinge on the successful identification of anomalies and the reduction of their detrimental impact on the aforementioned 052 tasks. As depicted in Figure 1, the integration of these three elements exhibits a significant reciprocal nature. In summary, representation learning, clustering, and anomaly detection are interdependent and



Figure 1: Interdependent relationships among representation learning, clustering, and anomaly detection.

064 065 066

067

068

063

intricately intertwined. Therefore, it is crucial for anomaly detection to *fully leverage and mutually enhance the relationships among these three components*.

Despite the intuitive significance of the interactions among representation learning, clustering, and 069 anomaly detection, existing methods have only made limited attempts to exploit them and fall short of expectations. On one hand, some methods Zong et al. (2018) have acknowledged the interplay among 071 these three factors, but their focus remains primarily on the interactions between two factors at a time, 072 making only targeted improvements. For instance, some strategies include explicitly removing outlier 073 samples during the clustering process Chawla & Gionis (2013) or designing robust representation 074 learning methods Cho et al. (2021) to mitigate the influence of anomalies. On the other hand, recent 075 methods Song et al. (2021) have begun to explore the simultaneous optimization of these three factors 076 within a single framework. However, these attempts are still in the stage of merely superimposing 077 the objectives of the three factors without a unified theoretical framework. This lack of a guiding framework prevents the adequate modeling of the interdependencies among these factors, thereby 079 limiting their collective contribution to a unified anomaly detection objective. Consequently, we aim to address the following question: Is it possible to employ a unified theoretical framework to jointly model these three interdependent objectives, thereby leveraging their respective strengths to enhance 081 anomaly detection?

083 In this paper, we try to answer this question and propose a novel model named UniCAD for anomaly 084 detection. The proposed UniCAD integrates representation learning, clustering, and anomaly de-085 tection into a unified framework, achieved through the theoretical guidance of maximizing the anomaly-aware data likelihood. Specifically, we explicitly model the relationships between samples 086 and multiple clusters in the representation space using the probabilistic mixture models for the 087 likelihood estimation. Moreover, we creatively introduce a learnable indicator function into the 088 objective of maximum likelihood to explicitly attenuate the influence of anomalies on representation learning and clustering. Under this framework, we can theoretically derive an anomaly score that 090 indicates the abnormality of samples, rather than heuristically designing it based on clustering results 091 as existing works do. Furthermore, building upon this theoretically supported anomaly score and 092 inspired by the theory of universal gravitation, we propose a more comprehensive anomaly metric 093 that considers the complex relationships between samples and multiple clusters. This allows us to 094 better utilize the learned representations and clustering results from this framework for anomaly detection. We conduct extensive experiments with 15 baselines on 30 datasets from different data 096 domains to evaluate the effectiveness of the proposed method. The results verify the effectiveness and generalization capability in detecting anomalies in real-world applications.

098 099

100

101

102

103

104

105

To sum up, we underline our contributions as follows:

- We propose a unified theoretical framework to jointly optimize representation learning, clustering, and anomaly detection, allowing their mutual enhancement and aid in anomaly detection.
- Based on the proposed framework, we derive a theoretically grounded anomaly score and further introduce a more comprehensive score with the vector summation, which fully releases the power of the framework for effective anomaly detection.
- Extensive experiments have been conducted on 30 datasets to validate the superior unsupervised anomaly detection performance of our approach, which surpassed the state-of-the-art through comparative evaluations with 17 baseline methods.

#### 108 **RELATED WORK** 2

109 110

Typical unsupervised anomaly detection (UAD) methods calculate a continuous score for each sample 111 to measure its anomaly degree. Various UAD methods have been proposed based on different assump-112 tions, making them suitable for detecting various types of anomaly patterns, including subspace-based 113 models Kriegel et al. (2009), statistical models Goldstein & Dengel (2012), linear models Wold et al. 114 (1987); Manevitz & Yousef (2001), density-based models Breunig et al. (2000); Peterson (2009), 115 ensemble-based models Pevný (2016); Liu et al. (2008), probability-based models Reynolds et al. 116 (2009); Zong et al. (2018); Li et al. (2022; 2020), representation-based models Ruff et al. (2018); Xu et al. (2023), and cluster-based models He et al. (2003); Chawla & Gionis (2013). Considering the 117 field of anomaly detection has progressed by integrating clustering information to enhance detection 118 accuracy Li et al. (2021); Zhou et al. (2022), we primarily focus on and analyze anomaly patterns 119 related to clustering, incorporating a global clustering perspective to assess the degree of anomaly. 120 Notable methods in this context include CBLOF He et al. (2003), which evaluates anomalies based 121 on the size of the nearest cluster and the distance to the nearest large cluster. Similarly, DCFOD Song 122 et al. (2021) introduces innovation by applying the self-training architecture of the deep cluster-123 ing Xie et al. (2016) to outlier detection. Meanwhile, DAGMM Zong et al. (2018) combines deep 124 autoencoders with Gaussian mixture models, utilizing sample energy as a metric to quantify the 125 anomaly degree. In contrast, our approach introduces a unified theoretical framework that integrates 126 representation learning, clustering, and anomaly detection, overcoming the limitations of heuristic 127 designs and the overlooked anomaly influence in existing methods.

128 129

#### 3 METHODOLOGY

130 131

In this section, we first define the problem we studied and the notations used in this paper. Then we 132 elaborate on the proposed method UniCAD. More specifically, we first introduce a novel learning 133 objective that optimizes representation learning, clustering, and anomaly detection within a unified 134 theoretical framework by maximizing the data likelihood. A novel anomaly score with theoretical 135 support is also naturally derived from this framework. Then, inspired by the concept of universal 136 gravitation, we further propose an enhanced anomaly scoring approach that leverages the intricate 137 relationship between samples and clustering to detect anomalies effectively. Finally, we present an 138 efficient iterative optimization strategy to optimize this model and provide a complexity analysis for 139 the proposed model.

140 **Definition 1** (Unsupervised Anomaly Detection). Given a dataset  $\mathbf{X} \in \mathbb{R}^{N \times D}$  comprising N 141 instances with D-dimensional features, unsupervised anomaly detection aims to learn an anomaly 142 score  $o_i$  for each instance  $\mathbf{x}_i$  in an unsupervised manner so that the abnormal ones have higher 143 scores than the normal ones.

144 145 146

## 3.1 MAXIMIZING ANOMALY-AWARE LIKELIHOOD

147 Previous research has demonstrated the importance of discriminative representation and accurate 148 clustering in anomaly detection Song et al. (2021). However, the presence of anomalous samples can 149 significantly disrupt the effectiveness of both representation learning and clustering Duan et al. (2009). 150 While some existing studies have attempted to integrate these three separate learning objectives, the 151 lack of a unified theoretical framework has hindered their mutual enhancement, leading to suboptimal 152 results.

153 To tackle this issue, in this paper, we propose a unified and coherent approach that considers 154 representation learning, clustering, and anomaly detection by maximizing the likelihood of the 155 observed data. Specifically, we denote the parameters of representation learning as  $\Theta$ , the clustering 156 parameter as  $\Phi$ , and the dynamic indicator function for anomaly detection as  $\delta(\cdot)$ . These parameters 157 are optimized simultaneously by maximizing the likelihood of the observed data X: 158

1

$$\max \log p(\mathbf{X}|\Theta, \Phi) = \max \sum_{i=1}^{N} \delta(\mathbf{x}_i) \log p(\mathbf{x}_i|\Theta, \Phi) = \max \sum_{i=1}^{N} \delta(\mathbf{x}_i) \log \sum_{k=1}^{K} p(\mathbf{x}_i, c_i = k|\Theta, \Phi),$$
(1)

where  $c_i$  represents the latent cluster variable associated with  $\mathbf{x}_i$ , and  $c_i = k$  denotes the probabilistic event that  $\mathbf{x}_i$  belongs to the k-th cluster. The  $\delta(\mathbf{x}_i)$  is an indicator function that determines whether a sample  $\mathbf{x}_i$  is an anomaly of value 0 or a normal sample of value 1.

3.1.1 JOINT REPRESENTATION LEARNING AND CLUSTERING WITH  $p(\mathbf{x}_i | \Theta, \Phi)$ 

Based on the aforementioned advantages of MMs, we estimate the likelihood  $p(\mathbf{x}_i | \Theta, \Phi)$  with mixture models defined as:

167

168

$$p(\mathbf{x}_{i}|\Theta, \Phi) = \sum_{k=1}^{K} p(\mathbf{x}_{i}, c_{i} = k|\Theta, \Phi) = \sum_{k=1}^{K} p(c_{i} = k) \cdot p(\mathbf{x}_{i}|c_{i} = k, \Theta, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

$$= \sum_{k=1}^{K} \omega_{k} \cdot p(\mathbf{x}_{i}|c_{i} = k, \Theta, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}),$$
(2)

where  $\Phi = \{\omega_k, \mu_k, \Sigma_k\}$ . The mixture model is parameterized by the prototypes  $\mu_k$ , covariance matrices  $\Sigma_k$ , and mixture weights  $\omega_k$  from all clusters.  $\sum_{k=1}^{K} \omega_k = 1$ , and  $k = 1, 2, \cdots, K$ .

In practice, the samples are usually attributed to high-dimensional features and it is challenging to detect anomalies from the raw feature space Ruff et al. (2021). Therefore, modern anomaly detection methods Ruff et al. (2018); Zong et al. (2018) often map raw data samples  $\mathbf{X} = {\mathbf{x}_i} \in \mathbb{R}^{N \times D}$  into a low-dimensional representation space  $\mathbf{Z} = {\mathbf{z}_i} \in \mathbb{R}^{N \times d}$  with a representation learning function  $\mathbf{z}_i = f_{\Theta}(\mathbf{x}_i)$  and detect anomalies within this latent representation space.

Following this widely adopted practice, we model the distribution of samples in the latent representation space with a multivariate Student's-*t* distribution giving its cluster  $c_i = k$ . The Student's-*t* distribution is robust against outliers due to its heavy tails. Bayesian robustness theory leverages such distributions to dismiss outlier data, favoring reliable sources, making the Student's-*t* process preferable over Gaussian processes for data with atypical information Andrade (2023). Thus the probability distribution of generating  $\mathbf{x}_i$  with latent representation  $\mathbf{z}_i$  given its cluster  $c_i = k$  can be expressed as:

191 192

193 194

$$p(\mathbf{x}_{i}|c_{i}=k,\Theta,\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k}) = \frac{\Gamma(\frac{\nu+1}{2})|\boldsymbol{\Sigma}_{k}|^{-1/2}}{\Gamma(\frac{\nu}{2})\sqrt{\nu\pi}} \left(1 + \frac{1}{\nu}D_{M}(\mathbf{z}_{i},\boldsymbol{\mu}_{k})^{2}\right)^{-\frac{\nu+1}{2}},$$
(3)

where  $\mathbf{z}_i = f_{\Theta}(\mathbf{x}_i)$  denotes the representation obtained from the data mapped through the neural network parameterized by  $\Theta$ .  $\Gamma$  denotes the gamma function while  $\nu$  is the degree of freedom.  $\Sigma_k$  is the scale parameter.  $D_M(\mathbf{z}_i, \boldsymbol{\mu}_k) = \sqrt{(\mathbf{z}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_k)}$  represents the Mahalanobis distance McLachlan (1999). In the unsupervised setting, as cross-validating  $\nu$  on a validation set or learning it is unnecessary,  $\nu$  is set as 1 for all experiments Xie et al. (2016); Van Der Maaten (2009). The overall marginal likelihood of the observed data  $\mathbf{x}_i$  can be simplified as:

201 202 203

204 205 206

207

$$p(\mathbf{x}_i|\Theta, \Phi) = \sum_{k=1}^{K} \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}.$$
(4)

#### 3.1.2 ANOMALY INDICATOR $\delta(\mathbf{x}_i)$ and Score $o_i$

As we have discussed, the indicator function  $\delta(\mathbf{x}_i)$  not only benefits both representation and clustering but also directly serves as the output of anomaly detection. Ideally, with the percentage of outliers denoted as l, an optimal solution for  $\delta(\mathbf{x}_i)$  that maximizes the objective function  $J(\Theta, \Phi)$  entails setting all  $\delta(\mathbf{x}_i) = 0$  for  $\mathbf{x}_i$  among the l percent of outliers with lowest generation possibility  $p(\mathbf{x}_i | \Theta, \Phi)$ , and otherwise  $\delta(\mathbf{x}_i) = 1$  is set for the remaining normal samples. Therefore, the indicator function is determined as:

214 215

 $\delta(\mathbf{x}_i) = \begin{cases} 0, & \text{if } p(\mathbf{x}_i | \Theta, \Phi) \text{ is among the } l \text{ lowest,} \\ 1, & \text{otherwise.} \end{cases}$ (5)

As this method involves sorting the samples based on the generation probability as being anomalous, the values of  $p(\mathbf{x}_i | \Theta, \Phi)$  can serve as a form of anomaly score, a classic approach within the mixture model framework Reynolds et al. (2009); Zong et al. (2018). This suggests that the likelihood of a sample being anomalous is inversely related to its generative probability since a lower generative probability indicates a higher chance of the sample being an outlier. Thus the anomaly score of sample  $\mathbf{x}_i$  can be defined as:

$$p_i = \frac{1}{p(\mathbf{x}_i | \Theta, \Phi)} = \frac{1}{\sum_{k=1}^{K} \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}}.$$
(6)

#### 3.2 GRAVITY-INSPIRED ANOMALY SCORING

In practical applications, it is proved that anomaly scores derived from generation probabilities often yield suboptimal performance Han et al. (2022). This observation prompts a reconsideration of *how to fully leverage the complex relationships among samples or even across multiple clusters for anomaly detection*. In this section, we first provide a brief introduction to the concept of Newton's Law of Universal Gravitation Newton (1833) and then demonstrate how the anomaly score is intriguingly similar to this cross-field principle. Finally, we discuss the advantages of introducing the vector sum operation into the anomaly score inspired by the analogy.

#### 3.2.1 ANALOG ANOMALY SCORING AND FORCE ANALYSIS

To begin with, Newton's Law of Universal Gravitation Newton (1833) stands as a fundamental framework for describing the interactions among entities in the physical world. According to this law, every object in the universe experiences an attractive force from another object. In classical mechanics, force analysis involves calculating the vector sum of all forces acting on an object, known as the **resultant force**, which is crucial in determining an object's acceleration or change in motion:

$$\vec{\mathbf{F}}_{i,\text{total}} = \sum_{k=1}^{K} \vec{\mathbf{F}}_{ik}, \text{ with } \vec{\mathbf{F}}_{ik} = \frac{G \cdot m_i m_k}{r_{ik}^2} \cdot \vec{\mathbf{r}}_{ik}, \tag{7}$$

where  $\vec{\mathbf{F}}_{ik}$  represents the k-th force acting on the object i. This force is proportional to the product of their masses,  $(m_i \text{ and } m_k)$ , and inversely proportional to the square of the distance  $r_{ik}$  between them. G represents the gravitational constant, and  $\vec{\mathbf{r}}_{ij}$  is the unit direction vector.

248 249 Similarly, if denoting:  $\widetilde{\mathbf{F}}_{ik} = p(\mathbf{x}_i, c_i = k | \Theta, \Phi) = \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M (\mathbf{z}_i, \boldsymbol{\mu}_k)^2}$ , the score of Equation equa-250 tion 6 bears analogies to the summation of the magnitudes of forces as:

222

224 225

226

227 228

229

230

231

232

233

234 235

236

242 243 244

$$o_i = \frac{1}{\sum_{k=1}^{K} \widetilde{\mathbf{F}}_{ik}}, \text{ with } \widetilde{\mathbf{F}}_{ik} = \frac{\widetilde{G} \cdot \widetilde{m}_i \widetilde{m}_k}{\widetilde{r}_{ik}^2}, \tag{8}$$

where  $\tilde{G} = \pi^{-1}$ ,  $\tilde{m}_k = \omega_k |\Sigma_k|^{-1/2}$ ,  $\tilde{m}_i = 1$ , and  $\tilde{r}_{ik} = \sqrt{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}$ . Here,  $\tilde{r}_{ik}$  is taken as the measure of distance within the representation space, modified slightly by an additional term for smoothness. The constant  $\tilde{G}$  serves a role akin to the gravitational constant in this analogy, whereas  $\tilde{m}_k$  resembles the concept of mass for the cluster. The notation  $\tilde{m}_i$  suggests a standardization where the mass of each data point is considered uniform and not differentiated.

# 260 3.2.2 ANOMALY SCORING WITH VECTOR SUM261

Comparing Equation equation 7 with Equation equation 8, what still differs is that, unlike a simple sum of the scalar value, the resultant force  $\vec{\mathbf{F}}_{i,\text{total}}$  employs the vector sum and incorporates both the magnitude and direction  $\hat{\mathbf{r}}_{ik}$  of each force. This distinction is crucial because forces in different directions can neutralize each other with a large angle between them or enhance each other's effects with a small angle. Inspired by this difference, we consider modeling the relationship between samples and clusters as a vector, and aggregating them through vector summation. The vector-formed anomaly score  $o_i^V$  is defined as:

$$o_i^V = \frac{1}{\|\sum_{k=1}^K \widetilde{\mathbf{F}}_{ik} \cdot \vec{\mathbf{r}}_{ik}\|},\tag{9}$$



where  $\vec{\mathbf{r}}_{ik}$  represents the unit direction vector in the representation space from the sample  $\mathbf{z}_i$  to the cluster prototype  $\boldsymbol{\mu}_k$ , and  $\|\cdot\|$  represents the  $L_2$  norm.

#### 284 285 3.2.3 Advantages of Vector Sum

The application of the vector sum principle extends beyond physical mechanics and finds relevance
 in various domains. In relational embedding Bordes et al. (2013), for example, relationships can be
 represented as vectors. Aggregating these vectors allows for capturing complexities like transitivity,
 symmetry, and antisymmetry.

290 Similarly, in our context, the vector sum can help capture more complex relationships along clusters. 291 In Figure 2, a sample v is attracted to two groups of cluster prototypes,  $\{\mu_1, \mu_2\}$  and  $\{\mu_3, \mu_4\}$ , 292 with equal mass and distances. While both groups exert equal forces, we argue that their influences 293 differ: a sample near two clusters with a large difference is more likely to be an anomaly than one near similar clusters. For instance, a user liking both money-saving tips and luxury items is more 294 295 anomalous than one liking two similar luxury items. The vector sum shows that the total force from  $\{\mu_1, \mu_2\}$  is smaller, leading to a higher anomaly score, thus demonstrating its effectiveness in 296 identifying subtle distinctions among clusters. 297

#### 298 299

281

282

283

#### 3.3 ITERATIVE OPTIMIZATION

Given the challenge posed by the interdependence of the parameters of the network  $\Theta$  and those of the mixture model { $\omega_k, \mu_k, \Sigma_k$ } in joint optimization, we propose an iterative optimization procedure. The pseudocode for training the model is presented in Algorithm 1.

**304** 3.3.1 UPDATE Φ

To update the parameters of the mixture model  $\Phi = \{\omega_k, \mu_k, \Sigma_k\}$ , we use the Expectation-Maximization (EM) algorithm to maximize equation equation 1 Peel & McLachlan (2000). The detailed derivation is included in Appendix B.

**E-step.** During the E-step of iteration (t + 1), our goal is to compute the posterior probabilities of each data point belonging to the *k*-th cluster within the mixture model. Given the observed sample  $\mathbf{x}_i$  and the current estimates of the parameters  $\Theta^{(t)}$  and  $\Phi^{(t)}$ , the expected value of the likelihood function of latent variable  $c_k$ , or the posterior possibilities, can be expressed as:

313 314 315

316 317

318

323

$$\boldsymbol{\tau}_{ik}^{(t+1)} = p(c_i = k | \mathbf{x}_i, \Theta, \Phi^{(t)}) = \frac{p(\mathbf{x}_i, c_i = k | \Theta, \Phi^{(t)})}{\sum_{j=1}^{K} p(\mathbf{x}_i, c_i = j | \Theta, \Phi^{(t)})} = \frac{\widetilde{\mathbf{F}}_{ik}^{(t)}}{\sum_{j=1}^{K} \widetilde{\mathbf{F}}_{ij}^{(t)}}.$$
 (10)

The scale factorPeel & McLachlan (2000) serving as an intermediate result for subsequent updates in the M-step is :

$$\mathbf{u}_{ik}^{(t+1)} = \frac{2}{1 + D_M(\mathbf{z}_i^{(t)}, \boldsymbol{\mu}_k^{(t)})}.$$
(11)

**M-step.** In the M-step of iteration (t + 1), given the gradients  $\frac{\partial J(\Theta, \Phi)}{\partial \omega_k} = 0$ ,  $\frac{\partial J(\Theta, \Phi)}{\partial \mu_k} = 0$ , and  $\frac{\partial J(\Theta, \Phi)}{\partial \Sigma_k} = 0$ , we derive the analytical solutions for the mixture model parameters  $\omega_k$ ,  $\mu_k$ , and  $\Sigma_k$ .

Algorithm 1 Model training for UniCAD 325 **Input:** data points **X**, cluster number K, outlier ratio l, tolerance  $\lambda$ , iterations t 326 **Output:** network parameters  $\Theta$ , mixture parameters  $\{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$ 327 1: Initialize  $\Theta$  and  $\{\mu_k, \omega_k, \Sigma_k\}$ ; 328 2: **for** i = 1 to t **do** if i = 1 then 3: 330 4:  $\mathbf{X}_i \leftarrow \mathbf{X};$ 331 5: else 332 6: Re-order the point in X such that  $o_1 \ge \cdots \ge o_n$ ;  $L_i \leftarrow \{x_1, \dots, x_{\lfloor N * l \rfloor}\}; \\ \mathbf{X}_i \leftarrow \mathbf{X} \setminus L_i;$ 7: 333 8: 334 9: end if 335 10: Update  $\Theta$  with equation 15; 336 while  $|J(\Theta, \Phi) - J^{old}(\Theta, \Phi)| > \lambda$  do 11: 337  $J^{old}(\Theta, \Phi) = J(\Theta, \Phi);$ 12: 338 Calculate  $\tau$  with equation 10; 13: 339 14: Update  $\{\omega_k, \mu_k, \Sigma_k\}$  with equation 12, equation 13 and equation 14; 340 15: end while 341 16: Calculate  $o_i$  with equation 9; 342 17: end for 18: return  $\Theta$  and  $\{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$ 343 344

Assume the anomalous ratio is  $l \in [0, 1]$ , the number of the normal samples is n = int(l \* N). The updating process for  $\{\omega_k^{(t+1)}, \boldsymbol{\mu}_k^{(t+1)}, \boldsymbol{\Sigma}_k^{(t+1)}\}\$  is as follows:

• The mixture weights  $\omega_k$  are updated by averaging the posterior probabilities over all data points with the number of samples, reflecting the relative presence of each component in the mixture:

$$\omega_k^{(t+1)} = \sum_{i=1}^n \tau_{ik}^{(t+1)} / n.$$
(12)

• The prototypes  $\mu_k$  are updated to be the weighted average of the data points, where weights are the posterior probabilities:

$$\boldsymbol{\mu}_{k}^{(t+1)} = \sum_{i=1}^{n} \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \mathbf{z}_{i} \right) / \sum_{i=1}^{n} \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \right).$$
(13)

• The covariance matrices  $\Sigma_k$  are updated by considering the dispersion of the data around the newly computed prototypes:

$$\boldsymbol{\Sigma}_{k}^{(t+1)} = \frac{\sum_{i=1}^{n} \boldsymbol{\tau}_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} (\mathbf{z}_{i} - \boldsymbol{\mu}_{k}^{(t+1)}) (\mathbf{z}_{i} - \boldsymbol{\mu}_{k}^{(t+1)})^{\mathsf{T}}}{\sum_{j=1}^{K} \boldsymbol{\tau}_{ij}^{(t+1)}}.$$
(14)

## 3.3.2 Update $\Theta$

324

345 346

347

348 349

350

351 352

353 354 355

356

362

368 369

370

371 372 We focus on anomaly-aware representation learning and use stochastic gradient descent to optimize the network parameters  $\Theta$ , by minimizing the following joint loss:

$$\mathcal{L} = -J(\Theta, \Phi) + g(\Theta), \tag{15}$$

373 where  $J(\Theta, \Phi) = \log p(\mathbf{X}|\Theta, \Phi)$ . An additional constraint term  $g(\Theta)$  is introduced to prevent 374 shortcut solution Geirhos et al. (2020). In practice, an autoencoder architecture is implemented, 375 utilizing a reconstruction loss  $g(\Theta) = ||x - \hat{x}||^2$  as the constraint. 376

These updates are iteratively performed until convergence, resulting in optimized model parameters 377 that best fit the given data according to the mixture model framework.

## 378 4 EXPERIMENTS

## 380<br/>3814.1DATASETS & BASELINES

We evaluated UniCAD on an extensive collection of datasets, comprising 30 tabular datasets that span 16 diverse fields. We specifically focused on naturally occurring anomaly patterns, rather than synthetically generated or injected anomalies, as this aligns more closely with real-world scenarios. The detailed descriptions are provided in Table 4 of Appendix D.1. Following the setup in ADBench Han et al. (2022), we adopt an inductive setting to predict newly emerging data, a highly beneficial approach for practical applications.

388 To assess the effectiveness of UniCAD, we compared it with 17 advanced unsupervised anomaly 389 detection methods, including: (1) traditional methods: SOD Kriegel et al. (2009) and HBOS Goldstein 390 & Dengel (2012); (2) linear methods: PCA Wold et al. (1987) and OCSVM Manevitz & Yousef 391 (2001); (3) density-based methods: LOF Breunig et al. (2000) and KNN Peterson (2009); (4) ensemblebased methods: LODA Pevný (2016) and IForest Liu et al. (2008); (5) probability-based methods: 392 DAGMM Zong et al. (2018), ECOD Li et al. (2022), and COPOD Li et al. (2020); (6) cluster-based 393 methods: DBSCAN Ester et al. (1996), CBLOF He et al. (2003), DCOD Song et al. (2021) and 394 KMeans-- Chawla & Gionis (2013); and (7) representation-based methods: DeepSVDD Ruff et al. 395 (2018) and DIF Xu et al. (2023). These baselines encompass the majority of the latest methods, 396 providing a comprehensive overview of the state-of-the-art. For a detailed description, please refer to 397 Appendix D.2.

398 399

400

#### 4.2 EXPERIMENT SETTINGS

401 In the unsupervised setting, we employ the default hyperparameters from the original papers for all 402 comparison methods. Similarly, the UniCAD also utilizes a fixed set of parameters to ensure a fair 403 comparison. For all datasets, we employ a two-layer MLP with a hidden dimension of d = 128 and 404 ReLU activation function as both encoder and decoder. We utilize the Adam optimizer Kingma & Ba 405 (2014) with a learning rate of  $1e^{-4}$  for 100 epochs. For the EM process, we set the maximum iteration number to 100 and a tolerance of  $1e^{-3}$  for stopping training when the objectives converge. The 406 407 number of components in the mixture model is set as k = 10, and the proportion of the outlier is set as l = 1%. We evaluate the methods using Area Under the Receiver Operating Characteristic (AUC-408 ROC) and Area Under the Precision-Recall Curve (AUC-PR) metrics Han et al. (2022), reporting the 409 average ranking (Avg. Rank) across all datasets. All experiments are run 3 times with different seeds, 410 and the mean results are reported. 411

412

#### 4.3 PERFORMANCE AND ANALYSIS

414 **Performance Comparison**. Table 1 presents a comparison of UniCAD with 10 unsupervised 415 baseline methods across 30 tabular datasets using the AUC-ROC metric. The experimental results, 416 which encompass 17 baselines, are included in Tables 5 and 6 of Appendix D.3, with additional 417 experiments on other data domains presented in Appendix E. Our proposed UniCAD achieves the 418 top average ranking, exhibiting the best or near-best performance on a larger number of datasets 419 and confirming advanced capabilities. It is noteworthy that there is no one-size-fits-all unsupervised 420 anomaly detection method suitable for every type of dataset, as demonstrated by the observation that 421 other methods have also achieved some of the best results on certain datasets. However, our model 422 showcased a remarkable ability to generalize across most datasets featuring natural anomalies, as evidenced by statistical average ranking. As for clustering-based methods such as KMeans--, DCOD, 423 and CBLOF, they mostly rank in the top tier among all baseline methods, supporting the advantage of 424 combining deep clustering with anomaly detection. However, our method significantly outperformed 425 these methods by mitigating their limitations and further providing a unified framework for joint 426 representation learning, clustering, and anomaly detection. 427

428 **Effectiveness of Vector Sum in Anomaly Scoring.** As demonstrated in Table 1, we compare the 429 anomaly score  $o_i$  derived directly from the generation possibility with its vector summation form  $o_i^V$ . 430 According to our statistical findings, we observe that vector scores  $o_i^V$  consistently outperform scalar 431 scores  $o_i$ . This indicates that the introduction of the vector summation, analogous to the concept of resultant force, makes a substantial difference in anomaly detection scenarios involving multiple Table 1: AUCROC of 10 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCROC is marked in red, the second highest in blue, and the third highest in green.

435													
436	Dataset	OC SVM	LOF	IForest	DA GMM	ECOD	DB SCAN	CBLOF	DCOD	KMeans	DIF	UniCAD w/ $o_i^S$	UniCAD w/ $o_i^V$
437	annthyroid	57.23	70.20	82.01	56.53	78.66	50.08	62.28	55.01	64.99	66.76	75.27	72.72
438	backdoor	85.04	85.79	72.15	55.98	86.08	76.55	81.91	79.57	89.11	92.87	87.28	89.24
400	breastw	80.30	40.61	98.32	N/A	99.17	85.20	96.86	99.02	97.05	77.45	98.15	98.56
439	campaign	65.70	59.04	71.71	56.03	76.10	50.60	64.34	63.16	63.51	67.53	73.52	73.64
440	celeba	70.70	38.95	70.41	44.74	76.48	50.36	73.99	91.41	56.76	65.29	81.38	82.00
	census	54.90	47.46	59.52	59.65	67.63	58.50	60.17	72.84	63.33	59.66	67.90	67.84
441	glass	35.36	69.20	77.13	76.09	65.83	54.55	78.30	78.07	77.30	84.57	79.52	82.17
442	Hepatitis	67.75	38.06	69.75	54.80	75.22	68.12	73.05	48.38	64.64	74.24	75.53	80.62
-1-12	http	99.59	27.46	99.96	N/A	98.10	49.97	99.60	99.53	99.55	99.49	99.53	99.52
443	Ionosphere	75.92	90.59	84.50	73.41	73.15	81.12	90.79	57.78	91.36	89.74	92.04	90.37
444	landsat	36.15	53.90	47.64	43.92	36.10	50.17	63.69	33.40	55.31	54.84	49.60	57.37
	Lymphography	99.54	89.86	99.81	72.11	99.52	74.16	99.81	81.19	100.00	83.67	99.29	99.73
445	mnist	82.95	67.13	80.98	67.23	74.61	50.00	79.96	65.23	82.45	88.16	86.00	86.64
440	musk	80.58	41.18	99.99	76.85	95.40	50.00	100.00	42.19	72.16	98.22	99.92	100.00
446	pendigits	93.75	47.99	94.76	64.22	93.01	55.33	96.93	94.33	94.37	93.79	95.12	95.52
447	Pima	66.92	65.71	72.87	55.93	63.05	51.39	71.49	72.16	70.44	67.28	75.16	74.87
	satellite	59.02	55.88	70.43	62.33	58.09	55.52	71.32	55.97	67.71	74.52	72.46	77.65
448	satimage-2	97.35	47.36	99.16	96.29	96.28	75.74	99.84	86.01	99.88	99.63	99.87	99.88
449	shuttle	97.40	57.11	99.56	97.92	99.13	50.40	93.07	97.20	69.97	97.00	99.15	98.75
445	skin	49.45	46.47	68.21	N/A	49.08	50.00	68.03	64.34	65.47	66.36	72.26	69.69
450	Stamps	83.86	51.26	91.21	88.89	87.87	52.08	69.89	93.41	79.78	87.95	91.37	94.18
454	thyroid	87.92	86.86	98.30	79.75	97.94	53.57	94.74	78.55	92.26	96.26	97.66	97.48
431	vertebral	37.99	49.29	36.66	53.20	40.66	49.74	41.01	38.13	38.14	47.20	33.11	47.37
452	vowels	61.59	93.12	73.94	60.58	62.24	57.50	92.12	51.56	93.45	81.02	88.38	92.09
450	Waveform	56.29	73.32	71.47	49.35	62.36	66.41	71.27	63.47	74.35	75.33	71.81	74.29
453	WBC	99.03	54.17	99.01	N/A	99.11	87.43	96.88	94.92	97.45	81.27	97.68	98.93
454	Wilt	31.28	50.65	41.94	37.29	36.30	49.96	34.50	44.71	34.91	39.46	48.95	52.56
	wine	73.07	37.74	80.37	61.70	77.22	40.33	27.14	82.18	27.36	41.69	82.72	95.25
455	WPBC	45.35	41.41	46.63	47.80	46.65	52.22	45.32	49.67	45.01	44.69	48.02	49.90
456	Avg. Rank	7.8	8.9	5.1	8.7	6.4	9.3	5.7	7.4	6.0	5.8	3.7	2.6



Figure 3: (a) demonstrates the performance variations during the optimization process on the satimage-2 dataset. (b) & (c) Analysis of cluster count k, anomaly ratio l.

clusters. The performance gains of the vector sum scores strongly demonstrate the effectiveness of the UniCAD in capturing the subtle differences in the distinctions among multiple clusters and underscore the utility of this factor in the context of anomaly detection based on clustering. 

Runtime Comparison. We present a analysis of the runtime performance of various methods, including our proposed approach, as detailed in Table 2. Our experiments, conducted on the backdoor dataset, reveal that while non-deep learning methods exhibit lower runtime, they often simplify the problem space excessively, failing to capture the complex non-linear relationships present in the data. In contrast, our method, when compared to existing deep learning techniques, demonstrates a significant reduction in computational time. This indicates that our approach not only manages to efficiently model complex patterns but also achieves an optimal balance between computational efficiency and modeling capability.



486 487 488

Table 2: Runtime Comparison. The runtime is reported in seconds (s).

Phase	IForest	KMeans	DAGMM	DCOD	UniCAD
Fit	0.256	103.697	795.004	4548.634	246.113
Infer	0.018	0.059	4.190	16.190	0.079

#### Table 3: Ablation study on AUC-ROC scores, calculated across 30 datasets.

Metric	w/ Gauss.	w/o $J(\Theta,\Phi)$	w/o $\delta(\mathbf{x}_i)$	Full Model
Avg. Rank (w/ baselines & variants)	6.2	6.6	5.0	4.2

#### 497 498 499

500

#### 4.4 ABLATION STUDIES

In this section, we examine the contributions of different components in UniCAD. Tables 3 reports the 501 results. We make three major observations. Firstly, the anomaly detection performance experiences a 502 significant drop when replacing the Student's t distribution with a Gaussian distribution for the Mixture Model, highlighting the robustness of the Student's t distribution in unsupervised anomaly detection. 504 **Secondly**, omitting the likelihood maximization loss (w/o  $J(\Theta, \Phi)$ ) also results in a considerable 505 decrease in overall performance. This observation underscores the importance of deriving both 506 the optimization objectives and anomaly scores from the likelihood generation probability through 507 a theoretical framework, which allows for unified joint optimization of anomaly detection and 508 clustering in the representation space. Furthermore, the indicator function  $\delta(\mathbf{x}_i)$  also contributes to a 509 performance increase. These results further confirm the effectiveness of our UniCAD in mitigating the 510 negative influence of anomalies in the clustering process, as the existence of outliers may significantly 511 degrade the performance of clustering. In summary, all these ablation studies clearly demonstrate 512 the effectiveness of our theoretical framework in simultaneously considering representation learning, clustering, and anomaly detection. 513

514 515

516

#### 4.5 HYPERPARAMETERS ANALYSIS

This section analyzes how hyperparameters affect our model's performance during the iterative 517 training process. As shown in Figure 3a, we tracked iteration counts from 0 to 10 for the satimage-2 518 dataset, keeping other parameters constant. The AUC-ROC and AUC-PR curves demonstrated stable 519 performance with only minor fluctuations initially, highlighting the convergence of the iterative EM 520 optimization. We also conducted a sensitivity analysis on key hyperparameters for the donors dataset, 521 focusing on the number of clusters k and the outlier set proportion l. The results, shown in Figure 3, 522 reveal that the optimal l is generally lower than the actual anomaly proportion. Furthermore, a pattern 523 was observed with the number of clusters k, where the model performance initially improved with an 524 increase in k, followed by a subsequent decline. This suggests the existence of an optimal range for 525 the number of clusters, which should be carefully selected based on the specific application context.

526 527

528

#### 5 CONCLUSION

529 This paper presents UniCAD, a novel model for Unsupervised Anomaly Detection (UAD) that 530 seamlessly integrates representation learning, clustering, and anomaly detection within a unified 531 theoretical framework. Specifically, UniCAD introduces an anomaly-aware data likelihood based on 532 the mixture model with the Student-t distribution to guide the joint optimization process, effectively 533 mitigating the impact of anomalies on representation learning and clustering. This framework 534 enables a theoretically grounded anomaly score inspired by universal gravitation, which considers complex relationships between samples and multiple clusters. Extensive experiments on 30 datasets 536 across various domains demonstrate the effectiveness and generalization capability of UniCAD, surpassing 15 baseline methods and establishing it as a state-of-the-art solution in unsupervised anomaly detection. Despite its potential, the proposed method's applicability to broader fields like 538 time series and multimodal anomaly detection requires further exploration and validation, highlighting a significant area for future work.

## 540 REFERENCES

547

553

565

566

567

568

577

578

- J Ailton A Andrade. On the robustness to outliers of the student-t process. Scandinavian Journal of Statistics, 50(2):725–749, 2023.
- Caglar Aytekin, Xingyang Ni, Francesco Cricri, and Emre Aksu. Clustering and unsupervised anomaly detection with 1 2 normalized deep auto-encoder representations. In 2018 International Joint Conference on Neural Networks (IJCNN), pp. 1–6. IEEE, 2018.
- Alexander Bakumenko and Ahmed Elragal. Detecting anomalies in financial data using machine learning algorithms. *Systems*, 10(5):130, 2022.
- Sambaran Bandyopadhyay, Saley Vishal Vivek, and MN Murty. Outlier resistant unsupervised deep architectures for attributed network embedding. In *Proceedings of the 13th international conference on web search and data mining*, pp. 25–33, 2020.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.
   Translating embeddings for modeling multi-relational data. Advances in neural information processing systems, 26, 2013.
- Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying density based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pp. 93–104, 2000.
- Raghavendra Chalapathy and Sanjay Chawla. Deep learning for anomaly detection: A survey. *arXiv* preprint arXiv:1901.03407, 2019.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. ACM computing
   *surveys (CSUR)*, 41(3):1–58, 2009.
  - Sanjay Chawla and Aristides Gionis. k-means-: A unified approach to clustering and outlier detection. In Proceedings of the 2013 SIAM international conference on data mining, pp. 189–197. SIAM, 2013.
- Hyunsoo Cho, Jinseok Seol, and Sang-goo Lee. Masked contrastive learning for anomaly detection.
   *arXiv preprint arXiv:2105.08793*, 2021.
- Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. Deep anomaly detection on attributed networks. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pp. 594–602. SIAM, 2019.
- Lian Duan, Lida Xu, Ying Liu, and Jun Lee. Cluster-based outlier detection. *Annals of Operations Research*, 168:151–168, 2009.
  - Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pp. 226–231, 1996.
- Haoyi Fan, Fengbin Zhang, and Zuoyong Li. Anomalydae: Dual autoencoder for anomaly detection
   on attributed networks. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5685–5689. IEEE, 2020.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
  Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track*, 1:59–63, 2012.
- Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. Adbench: Anomaly detection benchmark. *Advances in Neural Information Processing Systems*, 35:32142–32159, 2022.
- 593 Zengyou He, Xiaofei Xu, and Shengchun Deng. Discovering cluster-based local outliers. Pattern recognition letters, 24(9-10):1641–1650, 2003.

594 Meng Jiang. Catching social media advertisers with strategy analysis. In Proceedings of the First 595 International Workshop on Computational Methods for CyberSafety, pp. 5–10, 2016. 596 Ming Jin, Yixin Liu, Yu Zheng, Lianhua Chi, Yuan-Fang Li, and Shirui Pan. Anemone: Graph 597 anomaly detection with multi-scale contrastive learning. In Proceedings of the 30th ACM Interna-598 tional Conference on Information & Knowledge Management, pp. 3122–3126, 2021. 600 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint 601 arXiv:1412.6980, 2014. 602 Thomas N Kipf and Max Welling. Variational graph auto-encoders. arXiv preprint arXiv:1611.07308, 603 2016. 604 605 Yufeng Kou, Chang-Tien Lu, Sirirat Sirwongwattana, and Yo-Ping Huang. Survey of fraud detection 606 techniques. In IEEE International Conference on Networking, Sensing and Control, 2004, volume 2, pp. 749–754. IEEE, 2004. 607 608 Hans-Peter Kriegel, Peer Kröger, Erich Schubert, and Arthur Zimek. Outlier detection in axis-parallel 609 subspaces of high dimensional data. In Advances in Knowledge Discovery and Data Mining: 13th 610 Pacific-Asia Conference, PAKDD 2009 Bangkok, Thailand, April 27-30, 2009 Proceedings 13, pp. 611 831-838. Springer, 2009. 612 Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in 613 temporal interaction networks. In Proceedings of the 25th ACM SIGKDD international conference 614 on knowledge discovery & data mining, pp. 1269–1278, 2019. 615 616 Jinbo Li, Hesam Izakian, Witold Pedrycz, and Iqbal Jamal. Clustering-based anomaly detection in 617 multivariate time series data. Applied Soft Computing, 100:106919, 2021. 618 Zheng Li, Yue Zhao, Nicola Botta, Cezar Ionescu, and Xiyang Hu. Copod: copula-based outlier 619 detection. In 2020 IEEE international conference on data mining (ICDM), pp. 1118–1123. IEEE, 620 2020. 621 622 Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George Chen. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. IEEE Transactions on 623 Knowledge and Data Engineering, 2022. 624 625 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In 2008 eighth ieee international 626 conference on data mining, pp. 413-422. IEEE, 2008. 627 Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly detection on 628 attributed networks via contrastive self-supervised learning. IEEE transactions on neural networks 629 and learning systems, 33(6):2378-2392, 2021. 630 631 Xuexiong Luo, Jia Wu, Amin Beheshti, Jian Yang, Xiankun Zhang, Yuan Wang, and Shan Xue. 632 Comga: Community-aware attributed graph anomaly detection. In Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining, pp. 657–665, 2022. 633 634 Larry M Manevitz and Malik Yousef. One-class syms for document classification. Journal of machine 635 *Learning research*, 2(Dec):139–154, 2001. 636 637 Goeffrey J McLachlan. Mahalanobis distance. *Resonance*, 4(6):20–26, 1999. 638 Emmanuel Müller, Patricia Iglesias Sánchez, Yvonne Mülle, and Klemens Böhm. Ranking outlier 639 nodes in subspaces of attributed graphs. In 2013 IEEE 29th international conference on data 640 engineering workshops (ICDEW), pp. 216–222. IEEE, 2013. 641 Isaac Newton. *Philosophiae naturalis principia mathematica*, volume 1. G. Brookman, 1833. 642 643 David Peel and Geoffrey J McLachlan. Robust mixture modelling using the t distribution. Statistics 644 and computing, 10:339-348, 2000. 645 Zhen Peng, Minnan Luo, Jundong Li, Luguo Xue, and Qinghua Zheng. A deep multi-view framework 646 for anomaly detection on attributed networks. IEEE Transactions on Knowledge and Data 647

Engineering, 34(6):2539-2552, 2020.

648 649	Leif E Peterson. K-nearest neighbor. Scholarpedia, 4(2):1883, 2009.
650 651	Tomáš Pevny. Loda: Lightweight on-line detector of anomalies. <i>Machine Learning</i> , 102:275–304, 2016.
652 653 654	Douglas A Reynolds et al. Gaussian mixture models. <i>Encyclopedia of biometrics</i> , 741(659-663), 2009.
655 656 657	Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In <i>International</i> <i>conference on machine learning</i> , pp. 4393–4402. PMLR, 2018.
658 659 660	Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. <i>Proceedings of the IEEE</i> , 109(5):756–795, 2021.
662 663 664	Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In <i>Proceedings of the MLSDA 2014 2nd workshop on machine learning for sensory data analysis</i> , pp. 4–11, 2014.
665 666 667	Osman Salem, Yaning Liu, Ahmed Mehaoua, and Raouf Boutaba. Online anomaly detection in wireless body area networks for reliable healthcare monitoring. <i>IEEE journal of biomedical and health informatics</i> , 18(5):1541–1551, 2014.
669 670 671	Hanyu Song, Peizhao Li, and Hongfu Liu. Deep clustering based fair outlier detection. In <i>Proceedings</i> of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, pp. 1481–1489, 2021.
672 673 674	Jianheng Tang, Jiajin Li, Ziqi Gao, and Jia Li. Rethinking graph neural networks for anomaly detection. In <i>International Conference on Machine Learning</i> , pp. 21076–21089. PMLR, 2022.
675 676 677	Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Mehdi Azabou, Eva L Dyer, Remi Munos, Petar Veličković, and Michal Valko. Large-scale representation learning on graphs via bootstrapping. <i>arXiv preprint arXiv:2102.06514</i> , 2021.
678 679 680	Laurens Van Der Maaten. Learning a parametric embedding by preserving local structure. In <i>Artificial intelligence and statistics</i> , pp. 384–391. PMLR, 2009.
681 682	Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. <i>Chemometrics and intelligent laboratory systems</i> , 2(1-3):37–52, 1987.
683 684 685	Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In <i>International conference on machine learning</i> , pp. 478–487. PMLR, 2016.
686 687	Hongzuo Xu, Guansong Pang, Yijie Wang, and Yongjun Wang. Deep isolation forest for anomaly detection. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2023.
688 689 690 691	Xiaowei Xu, Nurcan Yuruk, Zhidan Feng, and Thomas AJ Schweiger. Scan: a structural clustering algorithm for networks. In <i>Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining</i> , pp. 824–833, 2007.
692 693 694 695	Zhiming Xu, Xiao Huang, Yue Zhao, Yushun Dong, and Jundong Li. Contrastive attributed network anomaly detection with data augmentation. In <i>Advances in Knowledge Discovery and Data Mining:</i> 26th Pacific-Asia Conference, PAKDD 2022, Chengdu, China, May 16–19, 2022, Proceedings, Part II, pp. 444–457. Springer, 2022.
696 697 698 699	Xu Yuan, Na Zhou, Shuo Yu, Huafei Huang, Zhikui Chen, and Feng Xia. Higher-order structure based anomaly detection on attributed networks. In <i>2021 IEEE International Conference on Big Data (Big Data)</i> , pp. 2691–2700. IEEE, 2021.
700 701	Yu Zheng, Ming Jin, Yixin Liu, Lianhua Chi, Khoa T Phan, and Yi-Ping Phoebe Chen. Generative and contrastive self-supervised learning for graph anomaly detection. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2021.

702 703 704	Shuang Zhou, Qiaoyu Tan, Zhiming Xu, Xiao Huang, and Fu-lai Chung. Subtractive aggregation for attributed network anomaly detection. In <i>Proceedings of the 30th ACM International Conference on Information &amp; Knowledge Management</i> , pp. 3672–3676, 2021.
705 706 707	Shuang Zhou, Xiao Huang, Ninghao Liu, Qiaoyu Tan, and Fu-Lai Chung. Unseen anomaly detection on networks via multi-hypersphere learning. In <i>Proceedings of the 2022 SIAM International</i>
707	Conference on Data Mining (SDM), pp. 262–270. SIAM, 2022.
709	Bo Zong, Oi Song, Martin Rengiang Min, Wai Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng
710	Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In
711	International conference on learning representations, 2018.
712	
713	
714	
715	
716	
717	
718	
719	
720	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
740	
747	
749	
750	
751	
752	
753	
754	
755	

## A ITERATIVE TRAINING ALGORITHM

#### B DERIVATION OF EM ALGORITHM

This appendix provides the detailed derivation of the Expectation-Maximization (EM) algorithm for optimizing the parameters of a mixture model based on Student's t-distribution. The focus is on deriving analytical solutions for the maximization of the parameters  $\Phi = {\mu_k, \Sigma_k, \omega_k}$  of the mixture components. The EM algorithm alternates between two steps:

In the E-step, we calculate the posterior probabilities  $\tau_{ik}$ , representing the probability of data point *i* belonging to cluster *k*, given the current parameters. The posterior probabilities for a Student's t-distribution mixture model are formulated as:

$$\boldsymbol{\tau}_{ik} = \frac{\omega_k \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \omega_j \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)},\tag{16}$$

where  $\tau(\mathbf{z}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  denotes the Student's t-distribution for data point *i* with respect to cluster *k*, and *K* is the number of mixture components.

The Student's t-distribution is depicted as a hierarchical conditional probability, resembling a Gaussian distribution with an accuracy scale factor **u**, where its latent variable follows a gamma distribution. Adopting a degree of freedom  $\nu = 1$ , the value of  $\mathbf{u}_{ik}$  is given by:

$$\mathbf{u}_{ik} = \frac{\nu + 1}{\nu + D_M(z_i, \boldsymbol{\mu}_k)} = \frac{2}{1 + D_M(z_i, \boldsymbol{\mu}_k)}$$
(17)

**In the M-step**, we update the parameters  $\Phi = \{\omega_k, \mu_k, \text{ and } \Sigma_k\}$  using the derivatives obtained in the previous steps. In our model, the likelihood function for a Student's-t Distribution Mixture Model (SMM) is represented as:

$$L(\omega, \mu, \Sigma) = \sum_{i=1}^{N} \sum_{k=1}^{K} \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-\frac{1}{2}}}{1 + (\mathbf{z}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{z}_i - \mu_k)},$$
(18)

where  $\omega_k$  are the mixture weights,  $\Sigma_k$  the covariance matrices,  $\mu_k$  the means, and  $\mathbf{z}_i$  the data points.

The derivative with respect to  $\omega_k$  must consider the constraint that the sum of the mixture weights equals 1, i.e.,  $\sum_k \omega_k = 1$ . Hence, we introduce a Lagrange multiplier  $\lambda$  to address this constraint and construct the Lagrangian L':

$$L'(\omega, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda}) = L(\omega, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \boldsymbol{\lambda} \left( 1 - \sum_{k=1}^{K} \omega_k \right),$$
(19)

The derivative with respect to  $\omega_k$  is:

$$\frac{\partial L'}{\partial \omega_k} = \frac{\partial L}{\partial \omega_k} - \lambda,\tag{20}$$

Substituting the definition of  $L(\omega, \mu, \Sigma)$ , we obtain:

$$\frac{\partial L}{\partial \omega_k} = \sum_i \frac{p(\mathbf{z}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \omega_j \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \sum_i \frac{\boldsymbol{\tau}_{ik}}{\omega_k},$$
(21)

To solve for  $\omega_k$ , we first multiply both sides of the equation by  $\omega_k$  and apply the constraint condition:

$$\sum_{k} \omega_k \left( \sum_{i} \frac{\tau_{ik}}{\omega_k} - \lambda \right) = 0, \tag{22}$$

<sup>808</sup> Upon further organization, we find that the Lagrange multiplier  $\lambda$  actually equals the total number of <sup>809</sup> data points N (since  $\sum_i \tau_{ik} = N_k$ , where  $N_k$  is the expected total number of data points belonging to the kth component, and the sum of all  $N_k$  equals the total number of data points N).

801 802 803

804 805

806 807

799 800

758

759 760

761

762

763

764

765

766

767 768

769

776

777 778

782 783

784 785

794

Nearest Cluster

Figure 4: Score comparison with other methods.

Finally, we can solve for  $\omega_k$ :

Group Anomalies

$$\omega_k = \frac{\sum_i \tau_{ik}}{N},\tag{23}$$

**Resultant Force** 

This result indicates that the weight  $\omega_k$  of each mixture component equals the proportion of the posterior probabilities of the data points it contains relative to all data points.

To update  $\mu_k$  and  $\Sigma_k$ , we consider the conditional expectation of the data log-likelihood function:

$$Q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \sum_{i=1}^{N} \boldsymbol{\tau}_{ik} \left( -\log(\pi) - \frac{1}{2} \log |\boldsymbol{\sigma}_k| + \frac{1}{2} \log u_{ik} -\frac{1}{2} \mathbf{u}_{ik} (\mathbf{z}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_k) \right)$$
(24)

Generative Probability

Maximizing  $Q(\mu_k, \Sigma_k)$  with respect to  $\mu_k$  leads to:

$$\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = \frac{1}{2} \sum_{i=1}^N \boldsymbol{\tau}_{ik} \mathbf{u}_{ik} (2\Sigma_k^{-1} \boldsymbol{\mu}_k - 2\Sigma_k^{-1} \mathbf{z}_{ik})$$
(25)

Setting  $\frac{\partial Q}{\partial \mu_k} = 0$  results in the updated mean  $\mu_k^{(t+1)}$ :

$$\boldsymbol{\mu}_{k}^{(t+1)} = \sum_{i=1}^{n} \left( \boldsymbol{\tau}_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \mathbf{z}_{i} \right) / \sum_{i=1}^{n} \left( \boldsymbol{\tau}_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \right).$$
(26)

Considering the derivative of  $Q(\mu_k, \Sigma_k)$  with respect to  $\Sigma_k^{-1}$ :

$$\frac{\partial Q}{\partial \Sigma_k^{-1}} = \frac{1}{2} \sum_{i=1}^N \boldsymbol{\tau}_{ik} \left( \Sigma_k - \mathbf{u}_{ik} (\mathbf{z}_i - \boldsymbol{\mu}_k) \times (\mathbf{z}_i - \boldsymbol{\mu}_k)^T \right).$$
(27)

Setting  $\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = 0$  yields the updated covariance matrix  $\boldsymbol{\Sigma}_k^{(t+1)}$ :

$$\Sigma_{k}^{(t+1)} = \frac{\sum_{i=1}^{n} \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} (\mathbf{z}_{i} - \boldsymbol{\mu}_{k}^{(t+1)}) (\mathbf{z}_{i} - \boldsymbol{\mu}_{k}^{(t+1)})^{T}}{\sum_{j=1}^{K} \tau_{ij}^{(t+1)}}.$$
(28)

### C ANOMALY SCORE WITH VECTOR SUM

C.1 TOY EXAMPLE

In the appendix, as illustrated in Figure 4, we investigated a toy example. We discussed a specific
 pattern of anomalies termed *group anomalies*, where a small number of anomalous samples cluster
 together. It is crucial to note that we do not claim this anomaly pattern is common in real-world data;
 our goal is merely to point out a specific anomaly pattern that is challenging for traditional cluster based anomaly detection methods to detect. Specifically, we utilize three Gaussian distributions with

864 high variance (each generating 300 data samples) and one with lower variance (generating 30 data 865 samples). Because the samples from the smaller Gaussian follow a different generative mechanism 866 and represent a minority in the dataset, we consider them anomalies.

867 We set the cluster number for KMeans-- and GMM at four, indicating that the Gaussian distribution 868 comprising anomalous samples was also recognized as a cluster. KMeans-- employs a cluster-based 869 approach, using the distance to the nearest cluster center as the anomaly score, while GMM uses 870 a probability-based approach, considering the samples' likelihood in the mixture model as the 871 anomaly score. However, both approaches are ineffective in this scenario. Rather than identifying the 872 small cluster as anomalous, they tend to misidentify samples on the peripheries of larger clusters as 873 anomalies.

874 By contrast, our scoring method views the entire small cluster as more likely anomalous, followed by 875 outlier samples on the margins of the larger clusters. This visualization provides a perspective that 876 distinguishes our method from previous efforts.

- 877
- 878

D

879 880 881

882 883

884 885

886

D.1 BENCHMARK DATASETS DETAILS

EXPERIMENTAL SUPPLEMENTARY

Due to space constraints in the main text, we utilized 30 public datasets from ADBench Han et al. (2022), covering all different types of data. The details of the 30 datasets are presented in Table 4.

#### Table 4: Statistics of tabular benchmark datasets.

887						
888	Data	# Samples	# Features	# Anomaly	% Anomaly	Category
889	annthyroid	7200	6	534	7 42	Healthcare
890	backdoor	95329	196	2329	7.42 2 44	Network
891	breastw	683	9	232	34.99	Healthcare
892	campaign	41188	62	4640	11.27	Finance
893	celeba	202599	39	4547	2.24	Image
894	census	299285	500	18568	6.20	Sociology
895	glass	214	7	9	4.21	Forensic
896	Hepaitis	80	19	13	16.25	Healthcare
897	http	567498	3	2211	0.39	Web
898	Ionosphere	351	33	126	35.90	Oryctognosy
899	landsat	6435	36	1333	20.71	Astronautics
000	Lymphography	148	18	6	4.05	Healthcare
001	magic.gamma	19020	10	6688	35.16	Physical
901	mnist	7603	100	700	9.21	Image
902	musk	3062	166	97	3.17	Chemistry
903	pendigits	6870	16	156	2.27	Image
904	Pima	768	8	268	34.90	Healthcare
905	satellite	6435	36	2036	31.64	Astronautics
906	satimage-2	5803	36	71	1.22	Astronautics
907	shuttle	49097	9	3511	7.15	Astronautics
908	skin	245057	3	50859	20.75	Image
909	Stamps	340	9	31	9.12	Document
910	thyroid	3772	6	93	2.47	Healthcare
011	vertebral	240	6	30	12.50	Biology
911	vowels	1456	12	50	3.43	Linguistics
912	Waveform	3443	21	100	2.90	Physics
913	WBC	223	9	10	4.48	Healthcare
914	Wilt	4819	5	257	5.33	Botany
915	wine	129	13	10	7.75	Chemistry
916	WPBC	198	33	47	23.74	Healthcare
917						

918 919	D.2 BASELINES DETAILS	
920 921	A comprehensive overview of the unsupervised anomaly detection methods is presented below.	
922	D.2.1 TRADITIONAL MODELS	
924 925 926	• Subspace Outlier Detection (SOD) Kriegel et al. (2009): Identifies outliers in varying sub- spaces of a high-dimensional feature space, targeting anomalies that emerge in lower-dimensional projections.	
927 928 929	• Histogram-based Outlier Detection (HBOS) Goldstein & Dengel (2012): Assumes feature independence and calculates outlyingness via histograms, offering scalability and efficiency.	
930	D.2.2 LINEAR MODELS	
931 932	• <b>Principal Component Analysis (PCA) Wold et al. (1987):</b> Utilizes singular value decomposition for dimensionality reduction, with anomalies indicated by reconstruction errors.	
933 934 935	• One-class SVM (OCSVM) Manevitz & Yousef (2001): Defines a decision boundary to separate normal samples from outliers, maximizing the margin from the data origin.	
936	D.2.3 DENSITY-BASED MODELS	
938 939	• Local Outlier Factor (LOF) Breunig et al. (2000) : Measures local density deviation, marking samples as outliers if they lie in less dense regions compared to their neighbors.	
940 941 942	• K-Nearest Neighbors (KNN) Peterson (2009): Anomaly scores are assigned based on the distance to the k-th nearest neighbor, embodying a simple yet effective approach.	
943	D.2.4 ENSEMBLE-BASED MODELS	
944 945 946	• Lightweight On-line Detector of Anomalies (LODA) Pevný (2016) : An ensemble method suitable for real-time processing and adaptable to concept drift through random projections and histograms.	
948 949	• Isolation Forest (IForest) Liu et al. (2008): Isolates anomalies by randomly selecting features and split values, leveraging the ease of isolating anomalies to identify them efficiently.	
950 951	D.2.5 PROBABILITY-BASED MODELS	
952 953 954	• Deep Autoencoding Gaussian Mixture Model (DAGMM) Zong et al. (2018): Combines a deep autoencoder with a GMM for anomaly scoring, utilizing both low-dimensional representation and reconstruction error.	
955 956	• Empirical-Cumulative-distribution-based Outlier Detection (ECOD) Li et al. (2022): Uses ECDFs to estimate feature densities independently, targeting outliers in distribution tails.	
957 958 959	• Copula Based Outlier Detector (COPOD) Li et al. (2020): A hyperparameter-free method leveraging empirical copula models for interpretable and efficient outlier detection.	
960 961	D.2.6 CLUSTER-BASED MODELS	
962 963 964	• <b>DBSCAN Ester et al. (1996):</b> A density-based clustering algorithm that identifies clusters based on the density of data points, effectively separating high-density clusters from low-density noise, and is widely used for anomaly detection in spatial data.	
965 966	• Clustering Based Local Outlier Factor (CBLOF) He et al. (2003): Calculates anomaly scores based on cluster distances, using global data distribution.	
968 969	• <b>KMeans Song et al. (2021):</b> Extends k-means to include outlier detection in the clustering process, offering an integrated approach to anomaly detection.	
970 971	• Deep Clustering-based Fair Outlier Detection (DCFOD) Chawla & Gionis (2013): Enhances outlier detection with a focus on fairness, combining deep clustering and adversarial training for representation learning.	

Table 5: AUCROC of 17 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCROC is marked in red, the second highest in blue, and the third highest in green.

975																				
976	Dataset	SOD	HBOS	PCA	OC SVM	LOF	KNN	LODA	IForest	DA GMM	ECOD	COPOD	DB SCAN	CBLOF	DCOD	KMeans	Deep SVDD	DIF	UniCAD (Scalar)	UniCAD (Vector)
977	annthyroid backdoor	77.38 68.77	60.15 71.56	66.24 80.16	57.23 85.04	70.20 85.79	71.69 80.58	41.02 66.38	82.01 72.15	56.53 55.98	78.66 86.08	76.80 80.97	50.08 76.55	62.28 81.91	55.01 79.57	64.99 89.11	76.09 78.83	66.76 92.87	75.27 87.28	72.72 89.24
978	breastw campaign	93.97 69.16	98.94 78.55 76.18	95.13 72.78 70.28	80.30 65.70 70.70	40.61 59.04	97.01 72.27 50.62	98.49 51.67	98.32 71.71 70.41	N/A 56.03	99.17 76.10 76.48	99.68 77.69 75.69	85.20 50.60	96.86 64.34 73.00	99.02 63.16	97.05 63.51 56.76	63.36 54.42	77.45 67.53	98.15 73.52	98.56 73.64
979	census glass	62.12 73.36	64.89 77.23	68.74 66.29	54.90 35.36	47.46 69.20	66.88 82.29	37.14 73.13	59.52 77.13	59.65 76.09	67.63 65.83	69.07 72.43	58.50 54.55	60.17 78.30	72.84 78.07	63.33 77.30	54.16 55.71	59.66 84.57	67.90 79.52	67.84 82.17
980	Hepatitis http	67.83 78.04	79.85 99.53	75.95 99.72	67.75 99.59	38.06 27.46	52.76 3.37	64.87 12.48	69.75 99.96	54.80 N/A	75.22 98.10	82.05 99.29	68.12 49.97	73.05 99.60	48.38 99.53	64.64 99.55	57.45 60.38	74.24 99.49	75.53 99.53	80.62 99.52
981	Ionosphere landsat	86.37 59.54 71.22	62.49 55.14 99.49	79.19	75.92 36.15 99.54	90.59 53.90	88.26 57.95 55.91	78.42 38.17 85.55	84.50 47.64 99.81	73.41 43.92 72.11	73.15 36.10 99.52	79.34 41.55 99.48	81.12 50.17 74.16	90.79 63.69 99.81	57.78 33.40 81.10	91.36 55.31	53.94 62.48 71.91	89.74 54.84 83.67	92.04 49.60 90.20	90.37 57.37 99.73
982	mnist musk	60.10 74.09	60.42 100.00	85.29 100.00	82.95 80.58	67.13 41.18	80.58 69.89	72.27 95.11	80.98 99.99	67.23 76.85	74.61 95.40	77.74 94.20	50.00 50.00	79.96 100.00	65.23 42.19	82.45 72.16	50.98 66.02	88.16 98.22	86.00 99.92	86.64 100.00
983	pendigits Pima	66.29 61.25	93.04 71.07	93.73 70.77	93.75 66.92	47.99 65.71	72.95 73.43	89.10 65.93	94.76 72.87	64.22 55.93	93.01 63.05	90.68 69.10	55.33 51.39	96.93 71.49	94.33 72.16	94.37 70.44	27.32 49.49	93.79 67.28	95.12 75.16	95.52 74.87
984	satellite satimage-2 shuttle	63.96 83.08 69.51	97.65 98.63	59.62 97.62 98.62	59.02 97.35 97.40	55.88 47.36 57.11	65.18 92.60 69.64	97.56 60.95	70.43 99.16 99.56	62.33 96.29 97.92	96.28 99.13	63.20 97.21 99.35	55.52 75.74 50.40	99.84 93.07	55.97 86.01 97.20	67.71 99.88 69.97	57.40 55.68 51.81	74.52 99.63 97.00	72.46 99.87 99.15	99.88 98.75
985	skin Stamps	60.35 73.26	60.15 90.73	45.26 91.47	49.45 83.86	46.47 51.26	71.46 68.61	45.75 87.18	68.21 91.21	N/A 88.89	49.08 87.87	47.55 93.40	50.00 52.08	68.03 69.89	64.34 93.41	65.47 79.78	45.69 59.48	66.36 87.95	72.26 91.37	69.69 94.18
986	thyroid vertebral	92.81 40.32 92.65	95.62 28.56 72.21	96.34 37.06 65.29	87.92 37.99 61.59	86.86 49.29 93.12	95.93 33.79 97.26	74.30 30.57 70.36	98.30 36.66 73.94	79.75 53.20 60.58	97.94 40.66 62.24	94.30 25.64 53.15	53.57 49.74 57.50	94.74 41.01 92.12	78.55 38.13 51.56	92.26 38.14 93.45	52.14 37.81 49.87	96.26 47.20 81.02	97.66 33.11 88.38	97.48 47.37 92.09
987	Waveform WBC	68.57 94.60	68.77 98.72	65.48 98.20	56.29 99.03	73.32 54.17	73.78 90.56	60.13 96.91	71.47 99.01	49.35 N/A	62.36 99.11	75.03 99.11	66.41 87.43	71.27 96.88	63.47 94.92	74.35 97.45	53.94 62.46	75.33 81.27	71.81 97.68	74.29 98.93
988	Wilt wine	53.25 46.11	32.49 91.36	20.39 84.37	31.28 73.07	50.65 37.74	48.42 44.98	26.42 90.12	41.94 80.37	37.29 61.70	36.30 77.22	33.40 88.65	49.96 40.33	34.50 27.14	44.71 82.18	34.91 27.36	45.90 64.26	39.46 41.69	48.95 82.72	52.56 95.25
989	Avg. Rank	51.28	51.24 8.26	46.01 8.98	45.35 11.59	41.41 13.59	46.59 10.00	49.31 13.24	46.63 7.09	47.80 13.24	46.65 9.19	49.34 8.29	52.22	45.32 8.07	49.67 10.90	45.01 8.71	44.01 15.48	44.69 8.38	48.02 5.41	49.90 3.59
990																				
991																				
992	-	4	6		8	10	12	1	4				-	6	<b></b>	B	10	1	2	
993															-					
994												Galaxy (4	4.5)						(12) CO	)F
995	Galaxy (3.) GMM (6.)	6) 6)							(14) De -(12) CO -(12) DA	epSVDL F		CBLOF (	7.4)— 7.5)—						L(12) De 	epSVDD F
996	COPOD (7.) CBLOF (7.)	3) 4)							-(12) LO -(12) LO	F DA		PCA ( HBOS (	<b>7.5)</b> 7.7)						—(11) DA —(10) SC	AGMM DD
997	HBOS (7. KMeans (7.	6)—— B)——							-(10) OC -(9.8) SC	SVM D	KN	COPOD ( 1eans (1	7.8)—— 8.1)——						—(10) LC —(9.9) O	DA CSVM
998	KNN (8.) PCA (8.)	2) 2)							-(8.4) EC	:OD		KNN (	8.2) 8.3)						(8.6) E	LOD
999				(a) A	AUC	-RO	С								(b)	AUC-P	R			
1000				. ,																
1001			F	igur	e 5:	Crit	ical	diffe	erence	e dia	gram	s for	AUC	C-ROO	C and	AUC	-PR.			
1002																				
1003																				
1004	D.2.7	Ref	PRES	ENT	ATIC	)N-F	BASE	ED M	IODE	LS										
1005	· ·						-			~~~			<b>D</b> 00							
1007	• Deep :	Sup	port	Vect	tor I	Jata	Des	scrip	otion	(Dee	epSV	DD)	Ruff	et al.	. (201	18): M	1n1m	izes	the vo	olume
1008	of a fly	per	spne	re en	cios	mg	netw	OIK	uata	repre	esent	ations	5, 180	lating	anoi	nanes	outs	ide i	ms sp	mere.
1009	. Doon	Icol	otion	For	oct -	for	Ano	malı	7 Dot	ootid		TE) V	7.11.01	al (	0022)	• T ]+;1;	700	laan	loorn	inato
1010	enhan	ce fr	aditi	onal	isol	atio	n fo	mary rest t	echn	ique	n (D s off	ering	imn	roved	anoi	nalv d	leteci	tion	in con	nng to nnlex
1011	datase	ts w	ith n	ninin	nal n	arar	nete	r tun	ing.	Ique	3, 011	ering	, imp	10,000	ano	inary c	ietee	lion	in coi	прих
1012	Guide				ran p															
1013	<b>F</b> 1	.1	1,			1			1	1.				• 1				c	. 1	
1014	Each me	etho	d's u	niqu	e me	echa	inisr	n an	d app	olicat	100 (	contex	xt pr	ovide	a ric	h land	lscap	e of	techn	iques
1015	approact	perv	ised	anor	nalý n an	uet	teculo Iv d	лі, 11. etect	ion c	umg hall⊴	ule I	ieid s	uive	ise m	euro	lologie	es an	u the	e orea	uui of
1016	approact	105 1	U lac	AIIII	5 all	oma	iy u	cicci	1011 0	114110	nges	•								
1017																				
1018	D 3 S	IIPPI	LEMI	ENT	ARY	Exu	DEBI	MEN	ΤΔΙ	RES	штя									
1019	2.5 0			., 17		-//1			1111	1100	0.010	,								
1020	In the ar	nen	div	we d	etail	l the	stat	istic	al an	alvei	s cor	ducte	ot he	comr	nare f	he ner	form	ance	e of v	arious
1021	anomaly	det	ecto	rs. W	Ve ol	btair	ned f	his c	liagra	am h		nduct	ing a	Fried	dman	test (	p-val	ue:	4.657	e-19).
1022	indicatir	ng si	ignif	icant	dif	fere	nces	amo	ong c	liffe	rent of	detect	tors.	We	utiliz	ed ave	erage	e rar	nks ar	nd the

Nemenyi test to generate the critical difference diagram, as shown in Figure 5. It is noteworthy that
 the vector version exhibits significantly superior performance compared to the scalar version across
 more methods. The detailed outcomes for the AUCROC and AUCPR metrics, spanning 30 datasets
 and against 17 baseline approaches, are showcased in Table 5 and Table 6.

Table 6: AUCPR of 17 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCPR is marked in red, the second highest in blue, and the third highest in green.

1029																				
1030	Dataset	SOD	HBOS	PCA	OC SVM	LOF	KNN	LODA	IForest	DA GMM	ECOD	COPOD	DB SCAN	CBLOF	DCOD	KMeans	Deep SVDD	DIF	UniCAD (Scalar)	UniCAD (Vector)
1001	annthyroid	18.84	16.99	16.12	10.37	15.71	16.74	7.06	30.47	9.64	25.35	16.58	7.60	13.74	10.01	15.41	21.75	18.93	26.37	25.03
1031	backdoor	37.07	4.96	31.29	8.79	26.14	44.37	13.84	4.75	5.47	10.72	7.69	21.04	7.03	6.77	15.47	55.70	41.46	37.77	36.36
1000	breastw	84.88	97.71	95.11	82.70	28.55	92.19	97.04	96.04	N/A	98.54	99.40	78.42	91.94	96.83	92.25	48.60	50.65	94.47	95.90
1032	campaign	19.14	38.01	27.90	29.25	14.59	27.18	14.11	32.26	14.54	36.65	38.58	11.43	20.88	19.61	18.86	16.75	26.52	27.66	27.12
	celeba	2.36	13.82	15.89	10.73	1.73	3.14	4.04	8.96	1.95	13.96	13.69	2.32	11.22	17.48	3.19	2.73	5.44	15.12	14.66
1033	census	8.54	8.68	10.02	6.82	5.48	9.04	5.03	7.78	9.03	9.46	9.92	7.52	7.52	10.92	8.13	8.42	7.42	9.70	9.75
	glass	18.73	11.82	10.05	8.02	20.11	20.26	13.37	10.99	24.58	15.35	9.78	6.88	11.57	9.66	14.66	8.46	18.86	13.29	15.33
103/	Hepatitis	24.73	37.73	36.65	29.44	13.67	21.95	30.90	26.25	22.93	32.80	41.50	22.31	36.54	19.53	25.14	30.04	34.93	36.08	43.37
1034	http	8.32	44.79	56.43	46.86	3.82	0.70	0.67	90.83	N/A	16.61	35.19	0.37	47.53	44.03	45.09	13.39	41.72	43.53	43.52
1025	Ionosphere	85.88	41.78	73.92	74.54	88.07	90.41	73.04	80.41	64.97	64.69	69.89	63.04	89.77	47.63	91.36	43.24	87.45	89.55	87.61
1035	landsat	26.38	22.03	16.18	16.21	24.69	24.65	18.86	19.81	24.48	16.24	17.48	20.80	31.05	15.57	22.40	36.92	24.35	20.84	23.27
	Lymphography	22.00	91.83	97.02	93.59	23.08	38.69	44.54	97.31	19.52	90.87	88.68	7.66	97.31	12.34	100.00	34.58	32.84	91.69	96.66
1036	mnist	19.15	12.51	39.93	33.20	20.90	35.55	25.86	27.71	23.75	17.45	21.35	9.21	30.60	23.59	37.12	20.18	44.55	41.19	41.94
	musk	1.59	100.00	99.89	10.61	2.82	9.65	47.60	99.61	32.76	50.13	34.79	3.16	100.00	2.87	37.55	8.78	/0./0	97.65	99.96
1037	pendigits	4.46	29.27	23.65	23.52	3.78	6.50	18.71	26.05	4.67	30.65	21.22	2.94	32.87	22.21	32.67	1.53	23.75	24.86	21.68
1007	Pima	48.24	50.01	54.03	50.00	47.18	55.14	44.09	55.82	41.55	50.45	55.19	36.65	52.99	50.24	53.50	35.02	46.34	54.66	54.23
1020	satellite	47.23	67.25	59.64	57.61	37.68	50.01	61.94	65.92	38.33	52.22	56.58	37.50	61.43	45.51	54.68	41.//	68.92	/1.68	75.13
1030	satimage-2	26.11	/8.04	85.69	82.71	4.50	39.14	80.52	93.45	22.07	64.49	/6.55	12.08	97.09	8.12	97.13	2.58	72.90	97.33	97.31
1000	snuttie	20.27	90.40	92.55	85.29	19.70	20.58	48.75	97.02	95.20 N/A	10.27	90.30	20.80	79.89	31.82	32.00	12.41	07.25	92.03	92.50
1039	SKIII	24.01	25.70	41.00	21.20	21.20	28.12	24.60	20.08	1N/A 42.72	22.21	17.99	20.89	28.34	20.29	25.58	12.00	23.50	42.30	28.72
	thyroid	23.56	50.08	44.34	21.23	20.81	3/ 08	14.68	63.11	16.06	51.06	19.64	9.44	29.90	10.56	31.69	2 70	50.36	60.00	60.06
1040	vertebral	11 70	0.23	10.49	10.04	14.24	10.57	0.68	10.46	15.24	11.84	8 80	13.11	11.43	11.58	10.54	10.62	14 31	0.78	12.96
	vowels	38.88	13.41	8.92	8 24	34.42	63.41	13.82	15.12	12.22	10.56	4 14	13.27	35.14	3 58	49.10	4 58	14.97	26.52	32.42
1041	Waveform	9.66	5.86	5 79	4 37	11 33	13.04	4 71	6.24	3.11	4 76	6.90	5 33	17.93	4 26	19.74	4 41	11.28	6.49	7.83
10-11	WBC	54.00	73 56	82.20	80.87	5 57	66.55	78.67	90.49	N/A	86.10	86.10	30.25	67.31	33.43	71.88	8 00	13 32	68 69	83.14
1010	Wilt	5 53	3.84	3.13	3.62	5.05	4 73	3 36	4 23	4 00	3.93	3.69	5 33	3 74	4.62	3.76	4.65	4.05	4.80	5.19
1042	wine	7.95	43.08	30.87	21.56	7.77	8.43	48.82	25.96	17.51	23.54	45.71	8.11	5.98	24.44	6.27	18.78	8.38	21.40	49.59
10.10	WPBC	25.62	23.04	23.01	22.93	20.29	21.49	25.39	22.42	22.49	21.24	22.81	23.86	21.08	22.86	20.58	25.00	20.73	22.71	24.90
1043	Avg Pank	10.83	8 10	8 31	11.14	13.24	9.36	11 79	7 29	11.96	9.36	9.53	14.91	8 53	11.07	9.03	13.41	9.10	631	4.74

#### D.3.1 DEGREES OF FREEDOM IN T-DISTRIBUTION

In fixed degrees of freedom scenarios, specifically when set to 1, the benefits of utilizing the t-distribution become less pronounced. Drawing from existing literature Xie et al. (2016); Van Der Maaten (2009), the flexibility to learn the degrees of freedom or to perform cross-validation on the validation set is particularly pertinent in unsupervised contexts. For the sake of simplicity and to minimize computational demands, we opted to maintain the degrees of freedom at 1, which provided robust performance while reducing complexity.

Table 7: Comparison of Performance: Learning vs. Fixed Degrees of Freedom

Metric	Learn v	Fix $v = 1$
AUC-ROC Avg. Rank	4.4	3.34
AUC-PR Avg. Rank	5.05	4.47

#### D.3.2 ABLATION STUDY ON HYPERPARAMETER SETTINGS

An ablation study was conducted to evaluate the impact of hyperparameters k and l. A grid search was performed over various values of these hyperparameters across 30 datasets, benchmarking against 17 baseline methods. The comprehensive results, showcasing average ranks based on AUC-ROC, are summarized in the following table:

Table 8: Results of Hyperparameter Grid Search

$l \backslash k$	10	20	30	40
0.01	3.34	4.31	4.69	4.71
0.05	4.44	4.23	4.65	4.88
0.10	4.27	4.46	4.48	4.88

> The findings indicate that the method exhibits robustness across specific parameter ranges. To ensure fair comparisons, a consistent parameter set (k = 10, l = 1%) was applied, demonstrating strong performance across the majority of datasets.

> Additionally, guidelines for selecting hyperparameters were examined. While techniques such as the elbow method and silhouette coefficient were considered for determining the optimal number of clusters, they proved to be time-consuming and exhibited weak correlation with anomaly detection

Dataset	# Nodes	# Edges	# Features	# Anomaly	Category
Disney	124	670	28	6	co-purchase network
Weibo	8,405	407,963	400	868	social media network
Reddit	10,984	168,016	64	366	user-subreddit network
T-Finance	39,357	42,445,086	10	1,803	trading network

Table 9: Statistics of graph benchmark datasets.

performance. An ensemble learning approach, which involved random searches of k values and aggregation of anomaly scores, showed promise in enhancing performance and model robustness for certain datasets. Future research will further explore this area.

## 1093 D.4 COMPLEXITY ANALYSIS

1080

1087 1088

1092

1095 The complexity of each iteration in UniCAD involves three parts: constructing the outlier set, updating the network parameters  $\Theta$ , and optimizing the mixture model using the EM algorithm. Constructing the outlier set requires a sorting operation, for which we use Numpy's built-in quantile calculation with a time complexity of  $\mathcal{O}(N \log N)$ . Considering the number of network parameters 1098 along with the computation of the loss function, the computational complexity for optimizing  $\Theta$  is 1099 approximately  $\mathcal{O}(TNDd + TNKd)$ . The EM algorithm for the Student's t mixture model includes 1100 two main steps: the E-step, where the complexity for computing the probability (or responsibility) 1101 of each data point belonging to each component is approximately  $\mathcal{O}(NKd)$ , and the M-step, where 1102 the full computational complexity of updating the parameters (mean, covariance matrix) of each 1103 component is  $\mathcal{O}(NKd^2)$ . In practice, we use diagonal covariance matrices, which reduces the 1104 update complexity to roughly  $\mathcal{O}(NKd)$ . If the EM algorithm requires T round to converge, its 1105 time complexity is approximately  $\mathcal{O}(TNKd)$ . Therefore, the time complexity for t-iterations is 1106  $\mathcal{O}(tN(\log N + Td(D+K))).$ 1107

# 1108 E ADDITIONAL EXPERIMENTS ON GRAPH

1110 E.1 BASELINES

Our proposed method was compared with 16 graph domain baseline methods grouped into three categories as follows:

Contrastive Learning-based Methods: This group includes CoLA Liu et al. (2021), SL-GAD Zheng et al. (2021), CONAD Xu et al. (2022), and ANEMONE Jin et al. (2021). These methods primarily assume that the contrastive loss between anomalous nodes and their neighborhoods is more significant.

Autoencoder-based Methods: This category consists of MLPAE Sakurada & Yairi (2014), GC-NAE Kipf & Welling (2016), DOMINANT Ding et al. (2019), GUIDE Yuan et al. (2021), ComGA Luo et al. (2022), AnomalyDAE Fan et al. (2020), ALARM Peng et al. (2020), DONE/AdONE Bandyopadhyay et al. (2020) and AAGNN Zhou et al. (2021). These methods focus on the reconstruction errors of anomalous nodes during the process of reconstructing the graph structure or features.

- Clustering-based Methods: This category of methods encompasses SCAN Xu et al. (2007), CBLOF He et al. (2003), and DCFOD Song et al. (2021). These methods generally identify anomalies by detecting if a sample deviates from the clustering.
- 1128

1129 E.2 DATASETS

We assess the performance of our model using four graph benchmark datasets containing organic anomalies. Table 9 presents the statistical summary for each dataset. These datasets contain naturally occurring real-world anomalies and are valuable for assessing the performance of anomaly detection algorithms in real-world scenarios. The sources and compositions of these datasets are as follows:

WeiboJiang (2016) is a labeled graph comprising user posts extracted from the social media platform Tencent Weibo. The user-user graph establishes connections between users who exhibit similar topic labels. A user is considered anomalous if they have engaged in a minimum of five suspicious events, whereas normal nodes represent users who have not.

RedditKumar et al. (2019) consists of a user-subreddit graph extracted from the popular social media platform Reddit. This publicly accessible dataset encompasses user posts within various subreddits over a month. Each user is assigned a binary label indicating whether they have been banned on the platform. Our assumption is that banned users exhibit anomalous behavior compared to regular Reddit users.

DisneyMüller et al. (2013) is a co-purchase network of movies that includes attributes such as price, rating, and the number of reviews. The ground truth labels, indicating whether a movie is considered anomalous or not, were assigned by high school students through majority voting.

T-FinanceTang et al. (2022) aims to identify anomalous accounts within a trading network. The nodes in this network represent unique anonymous accounts, each characterized by ten features related to registration duration, recorded activity, and interaction frequency. Graph edges denote transaction records between accounts. If a node is associated with activities such as fraud, money laundering, or online gambling, human experts will designate it as an anomaly.

#### 1152 E.3 EXPERIMENT SETTINGS 1153

Table 10: AUC-ROC and AUC-PR of 16 unsupervised algorithms on 4 graph benchmark datasets.

Group	Method	Weibo		Reddit		Disney		<b>T-Finance</b>	
		AUC-ROC	AUC-PR	AUC-ROC	AUC-PR	AUC-ROC	AUC-PR	AUC-ROC	AUC-PR
CL-Based	CoLA	0.382	0.087	0.527	0.036	0.455	0.060	0.243	0.031
	SL-GAD	0.421	0.109	0.594	0.040	0.494	0.061	0.442	0.041
	ANEMONE	0.320	0.082	0.536	0.036	0.454	0.068	0.226	0.030
	CONAD	0.806	0.432	0.551	0.037	0.600	0.138	N/A	N/A
AE-Based	MLPAE	0.880	0.629	0.501	0.035	0.563	0.064	0.299	0.030
	GCNAE	0.847	0.567	0.526	0.033	0.517	0.059	0.295	0.030
	GUIDE	0.897	0.692	0.566	0.040	0.521	0.060	N/A	N/A
	DOMINANT	0.927	0.797	0.561	0.037	0.590	0.077	N/A	N/A
	ComGA	0.925	0.809	0.568	0.037	0.494	0.058	N/A	N/A
	AnomalyDAE	0.892	0.694	0.560	0.037	0.520	0.070	N/A	N/A
	ALARM	0.952	0.843	0.559	0.037	0.595	0.123	N/A	N/A
	DONE	0.856	0.579	0.551	0.037	0.517	0.061	0.550	0.046
	AAGNN	0.804	0.530	0.564	0.045	0.479	0.059	N/A	N/A
Cluster-Based	SCAN	0.701	0.186	0.496	0.033	0.548	0.053	N/A	N/A
	CBLOF*	0.972	0.875	0.503	0.035	0.574	0.146	0.524	0.046
	DCFOD*	0.684	0.196	0.552	0.038	0.675	0.119	0.521	0.066
	UniCAD *	0.985	0.927	0.560	0.040	0.701	0.130	0.876	0.422

1169

1151

1154

1155

In this experiment, we compared graph-based methods on relational data. For methods originally
designed around feature vectors, including CBLOF, DCFOD, and our approach, we uniformly
employed the same graph representation learning technique as described in BGRL Thakoor et al.
(2021). Specifically, we used a two-layer Graph Convolutional Network (GCN) for encoding, which
produced output embeddings with a dimensionality of 128. The training epochs were set to 3000,
including a warm-up period of 300 epochs. The hidden size of the predictor was set to 512, and the
momentum was fixed at 0.99.

1177

#### 1178 E.4 PERFORMANCE ANALYSIS 1179

The performance of UniCAD compared to 16 baseline methods on the four datasets are summarized in
Table 10. From the results, we have the following observations: Our model consistently outperforms
the baseline methods on most datasets, underlining its effectiveness in anomaly detection even within
graph data contexts. This highlights the superiority of UniCAD in detecting anomalies in real-world
graph data.

1185 When comparing UniCAD with the four contrastive learning-based methods, it exhibits a distinct 1186 advantage, outperforming them by a substantial margin across all metrics. Unlike contrastive learning 1187 methods that rely on the local neighborhood for anomaly detection, UniCAD leverages the global clustering distribution. This key difference contributes to its consistently superior performance. Although CONAD incorporates human prior knowledge about anomalies, enabling it to outperform other similar methods on the Weibo and Disney datasets, it still falls short compared to our proposed UniCAD.

1191 Compared to the autoencoder-based methods, UniCAD offers the advantage of lower memory 1192 requirements along with better performance. Graph autoencoders typically reconstruct the entire 1193 adjacency matrix during full graph training, resulting in memory usage of at least  $\mathcal{O}(N^2)$ . In contrast, 1194 UniCAD, as a clustering-based method, only requires  $\mathcal{O}(N \times K)$ . Among the autoencoder-based 1195 methods, GCNAE, DONE, and AdONE can be extended to the T-Finance dataset as they only 1196 reconstruct the sampled subgraphs rather than the entire adjacency matrix. However, UniCAD still 1197 showcases superior performance while being more memory-efficient.

UniCAD also demonstrates superior performance compared to various other clustering-based methods, including traditional structural clustering (SCAN) methods that treat the embedding from BGRL as tabular data (CBLOF, DCFOD).