002

003 004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

GRAPH DISTRIBUTIONAL ANALYTICS: ENHANCING GNN EXPLAINABILITY THROUGH SCALABLE EMBED-DING AND DISTRIBUTION ANALYSIS

Anonymous authors

Paper under double-blind review

Abstract

Graph Neural Networks (GNNs) have achieved significant success in processing graph-structured data but often lack interpretability, limiting their practical applicability. We introduce the Graph Distributional Analytics (GDA) framework, leveraging novel combinations of scalable techniques to enhance GNN explainability. The integration of Weisfeiler-Leman (WL) graph kernels with distributional distance analysis enables GDA to efficiently quantify graph data distributions, while capturing global structural complexities without significant computational costs. GDA creates high-dimensional embeddings employing WL kernels, measures the distribution of distances from measures of categorical central tendency, and assigns distribution scores to quantify each graph's deviation from this vector We evaluate GDA on the ENZYMES, ogbg-ppa, and MalNet-Tiny datasets. Our experiments demonstrate GDA not only accurately characterizes graph distributions but also outperforms baseline methods in identifying specific structural features responsible for misclassifications. This comprehensive analysis provides deeper insights into how training data distributions affect model performance, particularly with out-of-distribution (OOD) data. By revealing the underlying structural causes of GNN predictions through a novel synergy of established techniques, GDA enhances transparency and offers a practical tool for practitioners to build more interpretable and robust graph-based models. Our framework's scalability, efficiency, and ability to integrate with various embedding methods make it a valuable addition to the suite of tools available for GNN analysis.

034 1 INTRODUCTION

Graph Neural Networks (GNNs) have emerged as a powerful framework for learning from graphstructured data, with applications ranging from social network analysis to molecular biology (Wu
et al., 2020). While GNNs excel at learning complex, non-Euclidean relationships, their decisionmaking processes remain largely opaque (Bodria et al., 2023). This lack of interpretability poses significant challenges, especially in critical domains like bioinformatics and healthcare, where model
decisions must be both accurate and explainable (Sanchez-Lengeling et al., 2020; Adoni et al., 2020).
Explainability is crucial for understanding why models make specific predictions, detecting potential
biases, and improving model generalization (Ju et al., 2024).

In recent years, efforts have been made to develop explainability tools tailored for GNNs. However, most existing methods focus on explaining node- or edge-level predictions, relying on gradient-based or perturbation-based approaches (Rajabi & Kafaie, 2022). While useful in some contexts, these methods struggle with graph-level predictions, particularly in datasets where structural diversity within categories complicates classification (Georgousis et al., 2021). Furthermore, many of these tools are computationally expensive and require significant customization to adapt to different datasets, limiting their scalability and practical applicability (Agarwal et al., 2022).

This paper introduces Graph Distributional Analytics (GDA), a novel framework designed to ad dress the limitations of existing GNN explainability methods. GDA leverages Weisfeiler-Leman
 graph kernels to embed graph structures into high-dimensional vector spaces, allowing for the characterization of structural distributions within graph categories. By analyzing the distribution of

graph embeddings, GDA provides insights at both population and sample levels, offering a scalable
 and interpretable approach to understanding GNN model behavior.

Through a series of experiments on the ENZYMES, MalNet-Tiny, and ogbg-ppa datasets, we demonstrate how GDA can reveal structural anomalies and distributional shifts that contribute to model underperformance. Unlike previous explainability methods, GDA not only identifies potential sources of confusion within the dataset but also offers a preemptive tool for improving data splits and detecting out-of-distribution samples. The contributions of this paper are threefold: (1) the introduction of GDA for graph-level explainability, (2) validation of GDA's utility through case studies on real-world datasets, and (3) a discussion of how GDA can inform future research into scalable, interpretable GNN models.

- 065 066
- 067 068

2 RELATED WORK

- 069 070
- 071
- 072

Demystifying neural network decision-making is a Herculean task, particularly for non-Euclidean data, where complex interconnections stymie explainability (Adoni et al., 2020). Current explainability methods demand significant computational resources (Sanchez-Lengeling et al., 2020), and the explanations offered by these tools tend to be approximate and difficult to interpret (Li et al., 2022). If explainability tools are difficult to explain, their usefulness is questionable. Nearly all methods require extensive tailoring to specific datasets and models which reduces practicability for diverse, real-world environments (Agarwal et al., 2022). Despite increased community interest, most libraries have prioritized accuracy benchmarks over explainability (Agarwal et al., 2023).

Gradient-based explainers, such as Grad (Simonyan et al., 2014), Grad-CAM (Pope et al., 2019),
 GuidedBP (Baldassarre & Azizpour, 2019), and Integrated Gradients (Sundararajan et al., 2017),
 highlight specific feature importances by backpropagating gradients. These methods work well for
 image or node classification but struggle with graph-level predictions (Kakkad et al., 2023), which
 are critical for datasets like MalNet-Tiny, where node and edge features are absent.

Perturbation-based explainers, such as GNNExplainer (Ying et al., 2019), PGExplainer (Luo et al., 2024), SubgraphX (Yuan et al., 2021), and GraphLIME (Huang et al., 2023), evaluate how model predictions change in response to input perturbations. D4Explainer (Chen et al., 2024) offers a more efficient method by using a discrete denoising diffusion model tailored towards in-distribution explanations. While these methods offer valuable insights, they are computationally expensive and typically explore only a fraction of possible perturbations, often missing the structural nuances that impact graph-level decisions (Munikoti et al., 2023).

Surrogate-based explainers, like PGMExplainer (Vu & Thai, 2020), use simpler models such as
 Bayesian networks to approximate GNN behavior. Although they provide interpretability, they face
 scalability challenges when dealing with high-dimensional probability distributions, and they require
 substantial tuning to capture relevant patterns in large datasets (Zhu et al., 2022).

Several gaps remain in GNN explainability. First, most approaches focus on node- and edge-level tasks, neglecting graph-level explainability crucial for structural prediction problems (Pope et al., 2019). Additionally, their computational inefficiency renders them impractical for dynamic, large-scale graph environments (Kazemi et al., 2020). Finally, due to the difficulty of characterizing graph data, the literature primarily focuses on explainability for Euclidean data, which leaves non-Euclidean datasets like MalNet-Tiny underexplored (Kakkad et al., 2023).

Graph Distributional Analytics (GDA) addresses these issues by offering scalable, dataset-agnostic
 explanations for both population- and sample-level predictions with linear time complexity, making
 it suitable for real-world applications. GDA also tackles the challenge of out-of-distribution (OOD)
 detection, enabling better interpretation of model behavior on unseen data by characterizing graph
 distributions. GDA fills this gap by analyzing graph structures without relying on node or edge
 features, providing a more versatile approach for complex, graph-structured data.

3 GRAPH DISTRIBUTIONAL ANALYTICS (GDA) FOR EXPLAINABILITY

3.1 GRAPH EMBEDDING AND STRUCTURAL DISTRIBUTION CHARACTERIZATION

112 Let G = (V, E) denote a graph where V represents the vertices and E the edges. To analyze the 113 structural properties of graphs, we first embed each graph as a high-dimensional vector using the 114 Weisfeiler-Leman (WL) graph kernel. For each graph, we define an embedding function $\tau : G \rightarrow$ 115 \mathbb{R}^d , where $\tau_h(G)$ represents the graph embedding after h iterations of the WL kernel, and d is equal 116 to the cardinality of unique labels in $\tau_h(G)$. Formally, the embedding is given by:

$$\tau_h(G) = \sum_{v \in V} l_h(v) \quad \text{where} \quad l_h(v) = f(l_{h-1}(v), \{l_{h-1}(u) : u \in \mathcal{N}(v)\}),$$

where $l_0(v)$ is the initial label of vertex v, $\mathcal{N}(v)$ represents the neighborhood of v, and f is a hash function aggregating neighborhood information (Shervashidze et al., 2011). This process captures the structural features of the graph, producing a vectorized representation that preserves graph topological properties (Morris et al., 2019). This is repeated for each G in a dataset H to create a set of embeddings:

$$\tau_h(H) = [\tau_h(G) \; \forall G \in H] \tag{1}$$

However, the embeddings generated are heterogeneous in both cardinality and the correspondence of mappings from labels to specific dimensions. Let $L = \{L_1, L_2, \dots, L_n \mid n = |H|, L_j \neq L_{i < j}\}$ be the set of unique vectors generated by $\tau_h(H)$ using the WL graph kernel. To standardize the embeddings, we define the Hamel dimension, a, as the cardinality of the set of unique labels L, and construct a vector space T, such that $\dim(T) = T^{|L|} = \mathbb{R}^a$. A mapping $g(L) \mapsto \mathbb{R}^a$ ensures that each unique label in L corresponds to a specific dimension in \mathbb{R}^a . For each embedding in H, we generate a uniform vector $\eta(G)$ as follows:

$$\forall G \in H, \eta(G) = g(\phi(G)) \in \mathbb{R}^a \tag{2}$$

The vectors within $\eta(H)$ have very high dimensionality, and many dimensions within T are sparse, containing non-zero values for only a few elements in $\eta(H)$. We filter these non-informative dimensions by constructing a linear subspace $U \subset \mathbb{R}^{b \leq a}$, defined by a function $\phi(G) \mapsto \mathbb{R}^{b}$:

$$\forall j \in \{1, \cdots, a\}, \eta(G)_j \in \phi(G) \iff \sum_{i=1}^{n=|H|} \eta(H_i)_j \le |H| \times \kappa$$
(3)

where $0 \le \kappa \le 1$ denotes a percentage threshold such that $0 \le |H| \times \kappa \le |H|$. Numerous methods 146 exist among the feature selection literature that can be used to define κ to optimize the removal of 147 irrelevant dimensions from the embedding vector (Zebari et al., 2020; Gui et al., 2016). However, 148 during the investigation of the effectiveness of GDA towards GNN explainability, we opted for a 149 conservative approach that set κ such that $|H| \times \kappa = \max(1, 0.002 \times |H|)$. This still significantly 150 reduces the dimensionality of the embedding since a majority of embeddings, regardless of dataset 151 size, are unique to one sample due to the combinatorial immensity of possible graph structures 152 (Bondy et al., 1976). 153

After embedding, GDA characterizes the distribution of graph embeddings within each class. Let \mathcal{G}_C represent the set of graphs belonging to class C. The mean embedding for class C is calculated as:

158

159

160

which serves as the central tendency for the graphs within the class. We then calculate the cosine similarity $S_c(G)$ between each graph embedding $\phi(G)$ and the class mean μ_C :

 $\mu_C = \frac{1}{|\mathcal{G}_C|} \sum_{G \in \mathcal{G}_C} \phi(G),$

143 144 145

136 137 138

139

140 141 142

108

110

111

126

127

- 162
- 163 164 165

167

168

 $S_c(G) = \frac{\phi(G) \cdot \mu_C}{\|\phi(G)\| \|\mu_C\|}.$

This similarity score quantifies how closely each graph aligns with the structural norm of its class.

To assess the degree of deviation from the class central tendency, we introduce a normalized distribution score for each graph G:

 $z(G) = \frac{S_c(G)}{\sqrt{\frac{1}{|\mathcal{G}_C|}\sum_{G' \in \mathcal{G}_C} S_c(G')^2}}.$

169 170

173 174

175

176

177 178

214

This score reflects how much a graph diverges from its class's prototypical structure. Graphs with lower z(G) values indicate significant deviations and are often associated with misclassifications or outliers. By analyzing these scores, GDA can reveal structural anomalies and patterns responsible for misclassification.

Algorithm 1 Graph Embedding and Structural Distribution Characterization using Weisfeiler-179 Leman Kernel 180 **Input:** Dataset H of graphs, iteration count h, threshold κ 181 **Output:** Set of filtered embeddings $\eta(H)$ and distribution scores z(G) for each $G \in H$ for each graph $G = (V, E) \in H$ do 183 Initialize node labels $l_0(v)$ for all $v \in V$ for i = 1 to h do 185 for each vertex $v \in V$ do 186 Update label: $l_i(v) = f(l_{i-1}(v), \{l_{i-1}(u) \mid u \in \mathcal{N}(v)\})$ 187 end for end for 188 Compute graph embedding: $\tau_h(G) = \sum_{v \in V} l_h(v)$ 189 end for 190 Let $L = \{L_1, L_2, \dots, L_n \mid n = |H|, L_j \neq L_{i < j}\}$ represent the set of unique labels across all 191 graphs 192 Define vector space $T \in \mathbb{R}^a$ where a = |L|193 for each graph $G \in H$ do 194 Map embedding: $\eta(G) = q(\tau_h(G)) \in \mathbb{R}^a$ end for 196 for each dimension $j \in \{1, \dots, a\}$ do 197 if $\sum_{i=1}^{|H|} \eta(H_i)_j \leq |H| \times \kappa$ then Remove dimension j from all embeddings $\eta(G) \in H$ 199 end if 200 end for Compute the filtered embeddings $\eta(H) \in \mathbb{R}^b$ where b < a201 for each class C in dataset H do 202 Compute class mean embedding: $\mu_C = \frac{1}{|\mathcal{G}_C|} \sum_{G \in \mathcal{G}_C} \eta(G)$ 203 for each graph $G \in \mathcal{G}_C$ do 204 Compute cosine similarity: $S_c(G) = \frac{\eta(G) \cdot \mu_C}{\|\eta(G)\| \|\mu_C\|}$ 205 Compute normalized distribution score: $z(G) = \frac{S_c(G)}{\sqrt{\frac{1}{|\mathcal{G}_C|}\sum_{G' \in \mathcal{G}_C} S_c(G')^2}}$ 206 207 end for 208 end for 209 **Output:** Filtered embeddings $\eta(H)$ and distribution scores z(G) for each $G \in H$ 210 211 212 213

3.2 OUTLIER DETECTION AND DISTRIBUTION ANALYSIS

In the context of Graph Distributional Analytics (GDA), we define outliers as graphs whose struc-215 tural properties significantly deviate from the central tendency of their class. Using the embeddings 216 generated by the Weisfeiler-Leman (WL) kernel, we quantify these deviations based on cosine sim-217 ilarity between the graph embeddings and the mean embedding of the class. Additionally, we assess 218 abnormal distributions within classes using kurtosis, a measure of how heavily tailed a distribution 219 is compared to a normal distribution.

221 3.2.1 OUTLIER DETECTION

220

225 226 227

228

229

230

231

232

240 241

242 243

222 We define outliers as graphs whose similarity score falls significantly below the average similarity 223 for the class. Specifically, a graph $G \in \mathcal{G}_C$ is considered an outlier if: 224

$$S_c(G) < \mu_{S_c} - \alpha \sigma_{S_c},$$

where μ_{S_c} and σ_{S_c} are the mean and standard deviation of cosine similarity scores across all graphs in class C, and α is a user-defined threshold (typically set to 2 or 3) that determines how many standard deviations away from the mean qualifies as an outlier. This method ensures that graphs with significantly lower similarity to the class mean are identified as outliers, which can help explain misclassifications or structural anomalies in the dataset.

233 3.2.2 ABNORMAL DISTRIBUTION DETECTION VIA KURTOSIS 234

235 To detect abnormal distributions within a class, we use kurtosis, $\xi(G_c)$ which measures the "tailedness" of the distribution of cosine similarity scores. High kurtosis indicates that the distribution has 236 heavy tails, meaning there are more extreme outliers than expected under a normal distribution. Low 237 kurtosis suggests that the distribution has lighter tails, indicating fewer extreme values. 238

239 Kurtosis for a set of cosine similarity scores $\{S_c(G) \mid G \in \mathcal{G}_C\}$ is defined as:

$$\xi(G_c) = \frac{1}{n} \sum_{i=1}^{n} \left(\frac{S_c(G_i) - \mu_{S_c}}{\sigma_{S_c}} \right)^4 - 3,$$

244 where μ_{S_c} and σ_{S_c} are the mean and standard deviation of the similarity scores, and $n = |G_C|$ is 245 the number of graphs in class C. A kurtosis value greater than 3 indicates a distribution with heavy 246 tails, suggesting that the class contains a significant number of structural outliers or subgroups with 247 distinct structural features. Conversely, a kurtosis value less than 3 indicates a distribution with light 248 tails, where most graphs are similar to the class mean. 249

By analyzing kurtosis, GDA can identify classes with abnormally high or low variability in graph 250 structures, which can provide insights into model performance. For instance, classes with high kur-251 tosis may indicate the presence of substructures that cause misclassifications or other performance 252 issues, while low kurtosis may suggest that the class is highly homogeneous and thus easier for the 253 model to learn. 254

255 3.3 POST-HOC STRUCTURAL ATTRIBUTION FOR MISCLASSIFIED GRAPHS 256

257 For misclassified graphs, GDA provides a post-hoc explanation by tracing structural deviations. 258 By rerunning the WL kernel with degree sequence tracking, we can identify specific substructures 259 responsible for classification errors. The process involves examining how node labels evolve through each iteration, enabling us to pinpoint graph substructures that deviate from the class norm. This 260 structural attribution mechanism allows for an interpretable analysis of misclassifications, offering 261 insights into how structural features influence GNN predictions. 262

264

263

3.4 Scalability and Integration with Existing Explainability Methods

265 GDA's key strength lies in its scalability. The embedding and distribution analysis have a compu-266 tational complexity of $O(n \cdot m)$, where n is the number of graphs and m is the average number of 267 nodes per graph. This ensures that GDA can handle large-scale graph datasets efficiently. 268

Moreover, GDA complements existing gradient-based and perturbation-based explainability tech-269 niques. While Grad-CAM and GNNExplainer focus on the node or edge-level importance, GDA

270	Dataset	Avg. # Nodes	Avg. # Edges	Node Features	Edge Features	# Categories
271	MalNet-Tiny	1,410.3	2,647.3	N/A	N/A	5
272	ogbg-ppa	243.4	2,266.1	N/A	6 features	37
273	ENZYMES	32.6	62.1	Atom type	N/A	6

Table 1: Dataset Statistics

Statistics highlighting graph order, graph size, number of classification categories, and the presence or absence of node features.

offers a higher-level, graph-wide perspective. This allows GDA to serve as a versatile framework, providing insights at both population and sample levels, including outlier detection and out-of-distribution (OOD) analysis.

4 EXPERIMENTS

284 285 286

287

288

289

290

295

300

304

305

306

307

308

309

310

311

312

313

314

315

316

321

274 275

276

277 278 279

280

281

282 283

The experiments conducted in this study validate the utility of Graph Distributional Analytics (GDA) for explaining GNN performance. GDA provides insights into model behavior and data distributions, but it does not prescribe specific actions to improve model performance. The results presented in this section demonstrate how GDA can inform further analysis and lead to improvements, though the implied action remains open and is subject to further research.

We present specific case studies that highlight the power of GDA at both population and sample levels. The body of this paper focuses upon how GDA is employed for GNN explainability. Detailed quantitative accountings of experiments are included in the appendices.

4.1 EXPERIMENTAL SETUP

All experiments were conducted by both models on all three datasets for ten runs with seeds (numpy, torch, random) set to $\{1, 2, \dots, 10\}$ corresponding to the run of the experiment. Initial experiments, also referred to as baseline experiments, used prescribed training, test, and validation splits as defined by the dataset.

301 4.1.1 DATASETS:

302 303 Three datasets were examined for this study:

- **MalNet-Tiny:** MalNet-Tiny (Freitas & Dong, 2021) is a curated subset of 5000 of the MalNet dataset graphs where each graph is selected to have fewer than 5000 nodes. Both datasets seek to classify function calls as indicative of malware or benign. It is designed for rapid prototyping of methods and seeks to strike a balance between the problem difficulty of the larger dataset and the computational demands of processing large data.
- **ogbg-ppa:** ogbg-ppa (Szklarczyk et al., 2019) is a graph-structural prediction dataset focused on identifying protein interactions between different species and taxa of organisms from the Stanford Open Graph Benchmark (Hu et al., 2020).
- ENZYMES: ENZYMES (Schomburg et al., 2004) is a part of the larger TUDataset benchmark (Morris et al., 2020). It seeks to identify protein tertiary structures categorized into six enzyme classes, with graphs representing the structural relationships between amino acids.
- 317 318 4.1.2 MODELS

To fully place the emphasis of the study upon the application of GDA for GNN explainability, we examined all datasets using the same model architectures to ensure relevance and comparability.

GraphSAGE: A scalable, inductive graph learning model that generates embeddings by sampling and aggregating features from a node's local neighborhood. This model is particularly useful for learning representations in large-scale graphs (Hamilton et al., 2017).



representations are shown in Figure ??. This finding is consistent with biochemical literature that
 highlights the variability within enzyme families (Giegé et al., 2012; Breton et al., 2006). The
 structural diversity in transferases led to underperformance in this category, with the GNN models
 struggled to capture the dual signals introduced by these two clusters.



Figure 2: Distribution Shifts in Malnet-Tiny

In Class 0 (left), we see a relatively normal distribution of similarities for a classification category with nearly all samples lying close to the central tendency vector. However, in Class 1 (right), we see an atypical distribution that has a significant portion of samples being nearly orthogonal to the central tendency.

To address this issue, we computed Pearson correlation coefficients on the embedding vectors, which allowed us to separate the two clusters and treat them as distinct categories during model training. After the training phase, we recombined the categories, leading to a 2.3% performance improvement for the transferase category, which also translated to a 0.4% increase in overall dataset performance across both models repeated with ten separate seeds. These results highlight the importance of identifying and handling structural heterogeneity within a class, which can significantly impact model accuracy. This observation aligns with prior work showing that structural complexity and variability within functional groups can lead to model confusion (Chen & Wu, 2022; Ji et al., 2021).

Moving to the MalNet-Tiny dataset, we found that three out of five categories exhibited relatively 402 normal distributions with cosine similarities clustered near 1 and low kurtosis. However, two cate-403 gories showed significant divergence, with some graph embeddings nearly orthogonal to their class 404 mean vectors, as seen in Figure 2. These outliers indicated the presence of structurally distinct 405 graphs that deviated from the majority of samples in their respective categories. Our hypothesis 406 was that these outliers were injecting noise into the model, leading to poorer performance in those 407 categories. To test this, we removed the outliers and retrained the models. While the improvements 408 were modest, the reduction in false positives was notable. We hypothesize that the model was al-409 ready learning to ignore these outliers during training in the baseline runs, as their removal only 410 reflected an equivalent reduction in false negative samples. This observation reinforces the idea 411 that models can effectively disregard noisy samples under certain conditions, but identifying and 412 removing them still contributes to cleaner classification metrics.

413 A related experiment involved analyzing distribution shifts between the training, validation, and test 414 sets in the MalNet-Tiny dataset. GDA revealed significant shifts in the distributions, with the train-415 ing set embeddings clustering more tightly around the mean compared to the test and validation sets. 416 These shifts, shown in Figure 3, likely caused overfitting, as the model was able to perform well on the more homogeneous training set but struggled to generalize to the more diverse test and validation 417 samples. To address this, we restructured the dataset splits to ensure more consistent distributions 418 across the sets. After retraining the models, we observed an average 4.3% improvement in perfor-419 mance over the baseline across 10 runs, using both GraphSAGE and GIN architectures. These re-420 sults emphasize the critical importance of ensuring distributional consistency between dataset splits 421 to prevent overfitting and improve generalization, as noted in previous studies on distributional shifts 422 in machine learning (Ding et al., 2021; Fan et al., 2024). 423

424 425

387

388

389

390

391

392 393 394

395

396

397

398

399

400

401

4.2.2 SAMPLE LEVEL ANALYSIS

Finally, we applied GDA to analyze individual misclassifications at the sample level, particularly
in the ogbg-ppa dataset. In this case, GDA identified a structural motif (Figure 4) that was present
in 12% of misclassified samples from Category 5 but was prevalent in 68% of Category 27 samples. This structural overlap likely led to misclassifications, as graphs containing this motif were
frequently classified as belonging to Category 27, even when they were from Category 5. Misclassification due to structural imbalances is a well-known issue in graph-based learning systems, and
similar findings have been reported in the literature (Ding et al., 2021). By isolating this motif and



Figure 3: Distribution Shifts

In classification category 2 of the MalNet-Tiny dataset, we see a significant shift between the training distribution, which is closer to the central tendency, and the testing and validation shifts which tend to be less similar. This may result in overfitting by the model due to the disparity between training and testing/validation distributions.



Figure 4: Embedding Analysis

This structure was present in a majority of misclassified examples in one class. We discovered it is a highly prevalent structure in another class which may have led to models becoming stuck in local minima when trying to differentiate these samples.

adjusting the training process to account for this structural overlap, we observed a significant reduction in misclassification rates between the two categories. This case study demonstrates the utility of GDA in identifying structural imbalances and mitigating their effects on model performance.

While benchmark datasets like those used in this study are useful for prototyping, their simplicity
precludes in-depth study of individual samples. The sample-level analysis provided by GDA is
most powerful when employed by those with domain expertise. The lack of detailed information in
protein, enzyme, and malware identification datasets, which would be required to tie structures to
real-world examples, hindered the application of this method in this study.

5 CONCLUSION

Graph Distributional Analytics (GDA) offers a novel and scalable framework for understanding the structural complexities that influence model performance in Graph Neural Networks (GNNs). By leveraging Weisfeiler-Leman graph kernels to embed graphs and analyze their distributions, GDA provides insights at both population and sample levels. The experiments conducted on the ENZYMES, MalNet-Tiny, and ogbg-ppa datasets demonstrate that GDA can effectively identify structural anomalies, distributional shifts, and misclassified samples that traditional methods may overlook. While GDA does not prescribe specific solutions, it offers a powerful tool for preemptive data analysis, improving dataset splits, and detecting out-of-distribution samples. Future work will focus on integrating GDA with more systematic model refinement strategies and further exploring its applications in real-world GNN tasks.

486 REFERENCES 487 Hamilton Wilfried Yves Adoni, Tarik Nahhal, Moez Krichen, Brahim Aghezzaf, and Abdeltif El-488 byed. A survey of current challenges in partitioning and processing of graph-structured data in 489 parallel and distributed systems. Distributed and Parallel Databases, 38:495-530, 2020. 490 491 Chirag Agarwal, Marinka Zitnik, and Himabindu Lakkaraju. Probing gnn explainers: A rigorous 492 theoretical and empirical analysis of gnn explanation methods. In International Conference on Artificial Intelligence and Statistics, pp. 8969–8996. PMLR, 2022. 493 494 Chirag Agarwal, Owen Queen, Himabindu Lakkaraju, and Marinka Zitnik. Evaluating explainability 495 for graph neural networks. Scientific Data, 10(1):144, 2023. 496 Federico Baldassarre and Hossein Azizpour. Explainability techniques for graph convolutional 497 networks. In International Conference on Machine Learning (ICML) Workshops, 2019 Work-498 shop on Learning and Reasoning with Graph-Structured Representations, 2019. URL https: 499 //urn.kb.se/resolve?urn=urn:nbn:se:kth:diva-260507. 500 501 Francesco Bodria, Fosca Giannotti, Riccardo Guidotti, Francesca Naretto, Dino Pedreschi, and Sal-502 vatore Rinzivillo. Benchmarking and survey of explanation methods for black box models. Data Mining and Knowledge Discovery, 37(5):1719–1778, 2023. 504 John Adrian Bondy, Uppaluri Siva Ramachandra Murty, et al. Graph theory with applications, 505 volume 290. Macmillan London, 1976. 506 507 Christelle Breton, Lenka Šnajdrová, Charlotte Jeanneau, Jaroslav Koča, and Anne Imberty. Structures and mechanisms of glycosyltransferases. Glycobiology, 16(2):29R-37R, 2006. 508 509 Jialin Chen, Shirley Wu, Abhijit Gupta, and Rex Ying. D4explainer: in-distribution gnn explanations 510 via discrete denoising diffusion. In Proceedings of the 37th International Conference on Neural 511 Information Processing Systems, NIPS '23, Red Hook, NY, USA, 2024. Curran Associates Inc. 512 Yu Chen and Lingfei Wu. Graph Neural Networks: Graph Structure Learning. Springer, 2022. 513 514 Mucong Ding, Kezhi Kong, Jiuhai Chen, John Kirchenbauer, Micah Goldblum, David Wipf, Furong 515 Huang, and Tom Goldstein. A closer look at distribution shifts and out-of-distribution general-516 ization on graphs. In NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and 517 Applications, 2021. URL https://openreview.net/forum?id=XvgPGWazqRH. 518 Shaohua Fan, Xiao Wang, Chuan Shi, Peng Cui, and Bai Wang. Generalizing graph neural networks 519 on out-of-distribution graphs. IEEE Transactions on Pattern Analysis and Machine Intelligence, 520 46(1):322-337, 2024. doi: 10.1109/TPAMI.2023.3321097. 521 522 Scott Freitas and Yuxiao Dong. A large-scale database for graph representation learning. Advances in neural information processing systems, 2021. 523 524 Stavros Georgousis, Michael P Kenning, and Xianghua Xie. Graph deep learning: State of the art 525 and challenges. IEEE Access, 9:22106-22140, 2021. 526 Richard Giegé, Frank Jühling, Joern Pütz, Peter Stadler, Claude Sauter, and Catherine Florentz. 527 Structure of transfer rnas: similarity and variability. Wiley Interdisciplinary Reviews: RNA, 3(1): 528 37-61, 2012. 529 530 Jie Gui, Zhenan Sun, Shuiwang Ji, Dacheng Tao, and Tieniu Tan. Feature selection based on struc-531 tured sparsity: A comprehensive study. *IEEE transactions on neural networks and learning* 532 systems, 28(7):1490–1507, 2016. Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. 534 In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 535 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/ 536 file/5dd9db5e033da9c6fb5ba83c7a7ebea9-Paper.pdf.

Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural Information Processing Systems*, 33:22118–22133, 2020.

- 540 Qiang Huang, Makoto Yamada, Yuan Tian, Dinesh Singh, and Yi Chang. Graphlime: Local inter-541 pretable model explanations for graph neural networks. IEEE Transactions on Knowledge and 542 Data Engineering, 35(7):6968-6972, 2023. 543
- Yang Ji, Samuel Deslauriers-Gauthier, and Rachid Deriche. Structure-function mapping via graph 544 neural networks. In International Workshop on Machine Learning in Clinical Neuroimaging, pp. 135–144. Springer, 2021. 546
- 547 Wei Ju, Siyu Yi, Yifan Wang, Zhiping Xiao, Zhengyang Mao, Hourun Li, Yiyang Gu, Yifang Qin, 548 Nan Yin, Senzhang Wang, et al. A survey of graph neural networks in real world: Imbalance, noise, privacy and ood challenges. CoRR, 2024. 549
- 550 Jaykumar Kakkad, Jaspal Jannu, Kartik Sharma, Charu Aggarwal, and Sourav Medya. A survey on 551 explainability of graph neural networks. Data Engineering, pp. 35, 2023. 552
- Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, and Peter Forsyth. 553 Representation learning for dynamic graphs: A survey. Journal of Machine Learning Research, 554 21(70):1-73, 2020. 555
- 556 Peibo Li, Yixing Yang, Maurice Pagnucco, and Yang Song. Explainability in graph neural networks: An experimental survey, 2022. URL https://arxiv.org/abs/2203.09258. 558
- Dongsheng Luo, Tianxiang Zhao, Wei Cheng, Dongkuan Xu, Feng Han, Wenchao Yu, Xiao Liu, 559 Haifeng Chen, and Xiang Zhang. Towards inductive and efficient explanations for graph neural 560 networks. IEEE Transactions on Pattern Analysis and Machine Intelligence, 46(8):5245–5259, 561 2024. 562
- 563 Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. 564 In Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, pp. 4602–4609, 565 2019. 566

569

576

577

578

- Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion 568 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020), 2020. URL www. 570 graphlearning.io.
- 571 Sai Munikoti, Deepesh Agarwal, Laya Das, Mahantesh Halappanavar, and Balasubramaniam 572 Natarajan. Challenges and opportunities in deep reinforcement learning with graph neural net-573 works: A comprehensive review of algorithms and applications. IEEE transactions on neural 574 networks and learning systems, 2023. 575
 - Phillip E Pope, Soheil Kolouri, Mohammad Rostami, Charles E Martin, and Heiko Hoffmann. Explainability methods for graph convolutional neural networks. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pp. 10772–10781, 2019.
- 579 Enayat Rajabi and Somayeh Kafaie. Knowledge graphs and explainable ai in healthcare. Informa-580 tion, 13(10):459, 2022.
- 581 Benjamin Sanchez-Lengeling, Jennifer Wei, Brian Lee, Emily Reif, Peter Wang, Wesley Qian, 582 Kevin McCloskey, Lucy Colwell, and Alexander Wiltschko. Evaluating attribution for graph 583 neural networks. Advances in neural information processing systems, 33:5898–5910, 2020. 584
- 585 Ida Schomburg, Antje Chang, Christian Ebeling, Marion Gremse, Christian Heldt, Gregor Huhn, 586 and Dietmar Schomburg. Brenda, the enzyme database: updates and major new developments. *Nucleic acids research*, 32(suppl_1):D431–D433, 2004.
- 588 Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M Borg-589 wardt. Weisfeiler-lehman graph kernels. Journal of Machine Learning Research, 12:2539–2561, 590 2011. 591
- K Simonyan, A Vedaldi, and A Zisserman. Deep inside convolutional networks: visualising im-592 age classification models and saliency maps. In Proceedings of the International Conference on Learning Representations (ICLR), pp. 1-8. ICLR, 2014.

- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In International conference on machine learning, pp. 3319–3328. PMLR, 2017.
- Damian Szklarczyk, Annika L Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta-Cepas, Milan Simonovic, Nadezhda T Doncheva, John H Morris, Peer Bork, et al. String v11:
 protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Research*, 47(D1):D607–D613, 2019.
- Minh Vu and My T Thai. Pgm-explainer: Probabilistic graphical model explanations for graph
 neural networks. *Advances in neural information processing systems*, 33:12225–12235, 2020.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A
 comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In International Conference on Learning Representations, 2019. URL https: //openreview.net/forum?id=ryGs6iA5Km.
- 610 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gn-611 nexplainer: Generating explanations for graph neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Ad-612 vances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 613 URL https://proceedings.neurips.cc/paper_files/paper/2019/ 2019. 614 file/d80b7040b773199015de6d3b4293c8ff-Paper.pdf. 615
- Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neural networks via subgraph explorations. In *International conference on machine learning*, pp. 12241–12252. PMLR, 2021.
- Rizgar Zebari, Adnan Abdulazeez, Diyar Zeebaree, Dilovan Zebari, and Jwan Saeed. A comprehensive review of dimensionality reduction techniques for feature selection and feature extraction. *Journal of Applied Science and Technology Trends*, 1(1):56–70, 2020.
 - Xiubin Zhu, Dan Wang, Witold Pedrycz, and Zhiwu Li. Fuzzy rule-based local surrogate models for black-box model explanation. *IEEE Transactions on Fuzzy Systems*, 31(6):2056–2064, 2022.
 - A APPENDIX

624 625 626

627 628

629

630

645 646 647 A.1 EXPERIMENTAL RESULTS

A.1.1 ENZYMES-BEFORE AND AFTER SPLITTING CLUSTERS

Run	1	2	3	4	5	6	7	8	9	10	Avg
GraphSAGE											
Before	83.3	80.6	83.0	81.8	81.9	80.1	83.5	82.3	82.8	81.7	82.1
After	82.8	82.8	83.7	82.5	83.5	81.3	83.9	83.4	82.9	80.2	82.7
GIN											
Before	83.3	81.4	83.8	80.4	84.9	83.2	84.9	82.7	84.2	82.5	83.1
After	82.0	81.8	83.1	84.4	82.1	84.7	84.5	81.7	84.6	84.4	83.3

A.1.2 MALNET-TINY DATASET-BEFORE AND AFTER REMOVING OUTLIERS

Run	1	2	3	4	5	6	7	8	9	10	Avg
GraphSAGE											
Before	84.0	82.7	81.4	85.2	81.9	79.1	79.0	79.2	82.4	82.8	81.8
After	81.1	81.0	81.1	82.0	82.7	82.8	83.1	81.6	83.4	81.1	82.0
GIN											
Before	83.5	87.1	83.2	81.0	86.1	87.5	84.6	86.3	86.7	86.5	85.2
After	86.0	84.8	83.8	85.0	86.3	86.9	85.7	85.4	85.9	84.3	85.4

Run	1	2	3	4	5	6	7	8	9	10	Avg
GraphSAGE											
Before	84.0	82.7	81.4	85.2	81.9	79.1	79.0	79.2	82.4	82.8	81.8
After	84.1	87.1	86.4	85.3	87.6	87.6	83.9	84.9	87.3	85.1	85.9
GIN											
Before	83.5	87.1	83.2	81.0	86.1	87.5	84.6	86.3	86.7	86.5	85.2
After	88.4	89.0	89.5	88.8	90.0	89.6	89.8	89.8	90.1	89.1	89.4

A.1.3 MALNET-TINY: BEFORE AND AFTER CORRECTING DISTRIBUTION SHIFTS



A.2 EXAMPLE VISUALIZATION OF FUNCTION CALL GRAPH FROM MALNET-TINY



A.3 VISUALIZATION OF EMBEDDING VECTOR IN MALNET-TINY







702	A.4.1 WEISFEILER-LEMAN EMBEDDING PROCESS	
703	For each graph $G_{i} = (V_{i}, E_{i})$ in the dataset H the Weisfeiler-Leman (WL) embedded	ling process is
705	performed. The embedding involves iterating over the graph's nodes, refining their	labels through
706	the WL kernel, and aggregating these labels to form the final graph embedding. The tin	ne complexity
707	for this embedding process is given by:	
708		
709	$T_{\mathrm{WI}}\left(G_{i} ight) = 5h \cdot V_{i} + 1$	(4)
710	where h is the number of iterations of the WI kernel $ V $ is the number of nodes	in graph C
711	where $-n$ is the number of iterations of the well kerner, $- v_i $ is the number of nodes	in graph G _i .
712	This accounts for the operations required to hash and sort node labels during each iter	ration.
714	A.4.2 MEAN VECTOR CALCULATION	
715	The calculation of mean vectors for each classification category involves summing the	e embeddings
716	of all graphs within the same category Let: $-m$ represent the number of classification	n categories -
717	H represent the total number of graphs in the dataset H.	n eurogonies,
718	The time complexity for computing the mean vectors is	
719	The time complexity for computing the mean vectors is:	
720		(-)
721	$T_{\text{mean}} = m + H $	(5)
722	This reflects the time required to iterate over all graphs and categories to compute the	neir respective
723	mean embeddings.	ion respective
724		
726	A.4.3 COSINE SIMILARITY CALCULATION	
727	Determining the cosine similarity for each graph with respect to its classification of	eategory mean
728	vector is a key step in GDA. The time complexity for this operation across all graphs	is:
729		101
730	$T_{\text{cosine}} = m + 2 H $	(6)
731		1 .
732	similarity between the graph embeddings and the category mean vectors.	ute the cosine
734 735	A.4.4 DISTRIBUTION SCORE CALCULATION	
736 737 738	Calculating the distribution score for each graph involves standardizing its cosine sime to the category's standard deviation. The time complexity for this step is:	ilarity relative
739	$T = -m \pm 3 H $	(7)
740	$L_{\text{score}} = m + \delta L $	(7)
741	This reflects the additional operations required for calculating and applying the standard	dard deviation
742	to each graph's cosine similarity.	
743		
744	A.4.5 OVERALL TIME COMPLEXITY	
745	The total time complexity for the entire GDA process is the sum of the complexities of	of each step:
747		
748		
749	$T_{\text{GDA}} = (5h+1)\sum V_i + 6 H + 3m$	(8)
750	$\overline{i=1}$	
751	Here $\sum H V = V_{rel} $ represents the total number of nodes screepe all graphs in the	dataset U In
752	nere, $\sum_{i=1}^{N} v_i = v_H $ represents the total number of nodes across all graphs in the the worst-case scenario, the time complexity is dominated by the term $ V ^{-1}$ leading	ualasel Π . In
753	complexity of V_H , neutron complexity is dominated by the term $ V_H $, leading	
754	complexity of.	

$$T_{\rm GDA} \in O(|V_H|) \tag{9}$$



Figure 5: Time and Space Empirical Analysis of GDA

774 A.4.6 STORAGE COMPLEXITY

The storage requirements for the graph embeddings are also an important consideration. Each embedding vector is normalized to a standard length based on the number of unique labels in the dataset. If every node in the dataset has a unique label, the worst-case storage requirement is:

$$S_{\text{GDA}} = |H| \times |V_H| \tag{10}$$

However, in practice, the storage requirement is often significantly smaller. For instance, in the ogbg-ppa dataset, which contains |H| = 158, 100 graphs with an average of $|V_i| = 243$ nodes each, the length of the embedding vectors was only 2,021. This indicates that the practical storage needs are much less than the theoretical worst case.

787 788

793 794

796

797 798

799 800

801

802

804

805

806

771

772 773

776

777

778

779

781

A.4.7 EMPIRICAL VALIDATION

The time and storage complexities were empirically validated using the datasets discussed in this paper. As shown in Figure 5, the observed time complexity closely matches the predicted linear growth with respect to the sum of the graph orders in the dataset, confirming the efficiency of the GDA approach.

A.5 HARDWARE AND SOFTWARE SPECIFICATIONS

All experiments were conducted using Python with several key libraries integral to the development and evaluation of our models:

- **PyTorch:** The primary deep learning library used for implementing neural network models.
- **PyTorch Geometric:** A specialized extension of PyTorch designed for working with graph-structured data, crucial for implementing the Graph Neural Networks (GNNs) evaluated in this study.
- **NumPy:** Used for efficient numerical computations, particularly for handling matrix operations and dataset preprocessing.
- **Matplotlib:** Utilized for generating the figures and plots presented in the main text, providing clear visualizations of the experimental results.
- 808
- All computations were performed on a cloud-based server with the following hardware specifications:

810	• 2 NUIDIA TA COLlas These COLla are entimized for deen learning tests, each equipped
811	• 2 NVIDIA 14 GFOS: These OFOS are optimized for deep rearining tasks, each equipped with 16 GB of GDDR6 memory, providing the necessary computational power for training
812	complex GNN models on large datasets
813	
814	• 30 GB KAIVI: This amount of system memory facilitated the handling of large datasets and the execution of memory intensive operations, portionlarly during the proprocessing and
815	training phases
816	
817	• 8 vCPUs: Virtual CPUs were used to manage concurrent processing tasks efficiently, en-
818	suring smooth operation and timely completion of the experiments.
819	• Ubuntu 20.04 LTS: The operating system provided a stable and secure environment for
820	running the software libraries and managing computational tasks.
821	
822	I nese nardware and software configurations ensured that the experiments were executed efficiently and that the regults were both reliable and reproducible. The combination of DyTerch, DyTerch Co
823	and that the results were bour reliable and reproducible. The combination of Py forch, Py forch Ge-
824	and evaluation of the GDA framework across the various datasets and models used in this study
825	and evaluation of the ODA framework across the various datasets and models used in this study.
826	
827	
828	
829	
830	
831	
832	
833	
834	
835	
836	
837	
838	
839	
840	
841	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
053	
855	
856	
857	
858	
859	
860	
861	
862	
863	