
Uncertainty-Aware Robust Learning on Noisy Graphs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph neural networks have shown impressive capabilities in solving various graph
2 learning tasks, particularly excelling in node classification. However, their effective-
3 ness can be hindered by the challenges arising from the widespread existence of
4 noisy measurements associated with the topological or nodal information present
5 in real-world graphs. These inaccuracies in observations can corrupt the crucial
6 patterns within the graph data, ultimately resulting in undesirable performance
7 in practical applications. To address these issues, this paper proposes a novel
8 uncertainty-aware graph learning framework motivated by distributionally robust
9 optimization. The framework aims to alleviate the challenges by considering the
10 distributional uncertainty associated with the graph data. Specifically, we use
11 a graph neural network-based encoder to embed the node features and find the
12 optimal node embeddings by minimizing the worst-case risk through a minimax
13 formulation. Such an uncertainty-aware learning process leads to improved node
14 representations and a more robust graph predictive model that effectively mitigates
15 the impact of uncertainty arising from data noise. The learned Least Favorable
16 Distributions (LFDs) also provide a means to quantify the predictive uncertainty,
17 which is valuable in some uncertainty-sensitive scenarios where incorrect decisions
18 can have severe consequence. In addition, we adopt the idea of differentiable
19 optimization and develop an end-to-end learning algorithm that seamlessly inte-
20 grates graph learning and distributionally robust optimization. Our experimental
21 result shows that the proposed framework achieves superior predictive performance
22 compared to the state-of-the-art baselines under various noisy settings.

23 1 Introduction

24 The field of graph learning has witnessed significant advancements in recent years, fueled by the
25 remarkable performance of graph neural networks (GNNs) [1, 2]. In general, GNNs leverage message-
26 passing techniques to enable efficient information exchange among nodes in a graph, leading to
27 improved graph embeddings. Among the various graph-based learning tasks, node classification,
28 especially in a semi-supervised setting, stands out as a prominent and widely applicable problem that
29 has greatly benefited from GNNs. The objective of semi-supervised node classification is to learn
30 high-quality node embeddings and make predictions on unlabeled nodes on a graph that has only a
31 small subset of nodes labeled.

32 However, similar to other deep neural networks, GNNs are susceptible to the influence of noisy
33 input, including both noisy node features and topological structure [3, 4, 5] of the graph. This issue
34 becomes even more critical in the low-data setting, where the performance can be greatly impaired by
35 noisy observations when only a limited number of labeled nodes are available [3, 4, 6]. For example,
36 in a social network, new users may inconsistently engage with content that aligns with their actual
37 interests or express their preferences due to limited options available at that time. Such noises would
38 introduce uncertainty and incorrect information about the user’s preferences into the training data.

39 With only a few labeled samples available, the GNN heavily relies on the noisy graph to learn the
 40 user’s preferences and generalize them to make less accurate recommendations. Hence, the GNN’s
 41 capacity to capture the underlying patterns of different nodes can be significantly compromised when
 42 the input graph is noisy. Therefore, it is challenging but imperative to achieve robust learning on
 43 noisy graphs, especially for the problem of semi-supervised node classification.

44 In this paper, we propose a novel uncertainty-aware graph learning framework based on the idea
 45 of distributionally robust optimization (DRO) [7, 8, 9, 10]. The proposed *Distributionally Robust*
 46 *Graph Learning* (DRGL) framework significantly improves node embedding quality and enhances
 47 model performance, particularly in the presence of noisy node features and topological structure that
 48 hinder the direct observation of key patterns. To tackle the challenge, we first allow the underlying
 49 node feature patterns to vary within a pre-defined family of distributions and seek the least favorable
 50 distribution (LFD) that minimizes the risk in the worst-case scenario, as illustrated by Figure 1. It
 51 is important to note that the uncertainty arising from data noise is represented by an uncertainty
 52 set [9, 11, 12], characterized by a Wasserstein ball [9, 11, 13, 14]. Then, by minimizing the risk
 53 associated with this worst-case distribution, we can uncover the most robust node embeddings that
 54 effectively mitigate the impact of uncertainty stemming from the noisy observation. To seamlessly
 55 integrate this minimax formulation into the graph learning process, we also leverage differentiable
 56 optimization techniques [15, 16] and develop a tractable end-to-end learning algorithm. Hence, the
 57 resulting minimax solution also provides a means to estimate the potential uncertainty associated
 58 with predictions, allowing decision-makers to assess the uncertainty of different outcomes. This
 59 information is particularly valuable in some high-stakes scenarios where incorrect decisions can have
 60 severe consequences. In summary, the contributions of this paper are threefold:

- 61 • We study the problem of robust graph learning under a challenging context, where both
 62 node attributes and topological structure are noisy.
- 63 • We present a model-agnostic graph learning framework – *Distributionally Robust Graph*
 64 *Learning* (DRGL), which improves the reliability of GNNs against noisy graph inputs, espe-
 65 cially when labeled data is limited.
- 66 • We conduct extensive experiments that demonstrate the proposed framework significantly
 67 enhances the quality of node embeddings and improves the predictive performance of GNNs
 68 on semi-supervised node classification.

69 2 Related Work

70 Graph neural networks (GNNs) have emerged
 71 as powerful models aiming to learn expressive
 72 representations for graph-structure data [1, 2,
 73 17, 18]. While these approaches proved to be
 74 successful [1, 2, 19], they are highly sensitive to
 75 the quality of the given graphs in terms of either
 76 node feature or graph structure [3, 5].

77 To address the robustness limitations of GNNs
 78 in the presence of noises in node features and
 79 edges, researchers have explored different meth-
 80 ods to purify the noisy graph [3, 4, 20]. For ex-
 81 ample, graph structure learning (GSL) focuses
 82 on the joint learning of an optimized graph struc-
 83 ture and its corresponding representations. GSL
 84 methods typically leverage GNNs as graph en-
 85 coders and employs specialized graph structure
 86 learners to refine and recover accurate graph
 87 structures, thereby enhancing the robustness of
 88 GNNs. [21, 22, 23, 24]. Another line of research focuses on improving the robustness of GNNs
 89 against inaccurate or missing node features [25, 26, 27] through developing noise-resistant aggrega-
 90 tion and propagation techniques.

91 However, a significant limitation of these methods is that they require a task-specific design of
 92 structure learning or regularization term tailored to address specific types of noise. This inherent

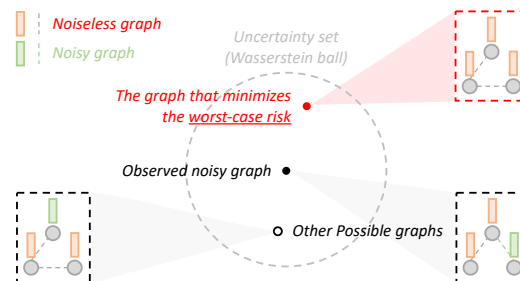


Figure 1: An illustration of the uncertainty set in our proposed framework. The goal is to search for the graph distribution that minimizes the worst-case risk.

93 task-specific nature limits their applicability in more general application scenarios. For instance,
 94 [5, 28] propose parametric learners that aim to denoise graph through pruning edges, but such methods
 95 cannot handle noise due to inaccuracies in node features. In contrast, our method handles both edge
 96 and feature noise directly in the GNN embedding space, offering a more unified and direct approach
 97 for robust graph representation learning.

98 Another related research area is defense against graph attacks, which often leverages graph properties
 99 observed in real-world graphs to regularize the learning process. Notable methods in this area include
 100 GCN-SVD [29], which preprocesses GCN with the low-rank approximation of the perturbed graph.
 101 This approach is inspired by the observation that only the high-rank (low-valued) singular components
 102 of the graph are susceptible to perturbation. Additionally, ProGNN [24] is designed to guide the
 103 propagation process by intrinsic graph properties such as low-rank and sparsity in the graph structure,
 104 as well as the tendency for adjacent nodes to exhibit similar features. RGCN [30] models node
 105 representations as Gaussian distributions to mitigate adversarial attacks. However, these methods
 106 are primarily designed to handle various types of attacks, including random perturbations, in graphs
 107 with ample samples. To the best of our knowledge, they have not been extensively tested in low-data
 108 scenarios against random graph noise, which distinguishes our setting from this body of research.

109 In contrast to the aforementioned methods, our approach handles noisy graphs based on the DRO
 110 [7, 8], which aims to enhance the robustness of GNNs by minimizing the worst-case risk in the
 111 presence of data noise. Only a handful of recent research attempts, such as [11, 31], have embraced
 112 similar concepts and investigated the potential applications of DRO in the context of graph-structured
 113 data. For instance, Sadeghi et al. [31] use GNNs to encode the input graph and employ DRO for
 114 model training. However, their work is mainly focused on the noisy node features, while we consider
 115 a more general scenario involving both noisy features and edges. Furthermore, their method solves
 116 for worst-case distribution and updates model parameters separately at each iteration [31], while our
 117 approach offers an end-to-end approach updating parameters with gradient-based learning techniques.
 118 Additionally, our approach provides uncertainty quantification under Least Favorable Distributions
 119 (LFDs). This capability has meaningful applications in uncertainty-sensitive scenarios [32, 33], such
 120 as molecular classification tasks in graphs [34, 35].

121 3 Uncertainty-aware Graph Learning

122 3.1 Problem Setup

123 Consider an attributed graph $(\mathcal{I}, \mathcal{E})$, where $\mathcal{I} = \{i = 1, \dots, n\}$ represents the set of n nodes, and
 124 $\mathcal{E} = \{(i, j), i, j \in \mathcal{I}\}$ represents the set of edges connecting the nodes. Each node i is associated
 125 with a d -dimensional feature vector $x_i \in \mathbb{R}^d$. The collection of all node features is denoted as
 126 $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{n \times d}$. To describe the graph more generally, we can represent it as
 127 $\mathcal{G} = (X, A)$, where $A = (a_{ij}) \in \{0, 1\}^{n \times n}$ is the adjacency matrix, which provides information
 128 about the connectivity between nodes. If $(i, j) \in \mathcal{E}$, then $a_{ij} = 1$; otherwise, $a_{ij} = 0$. Each
 129 node can be assigned one of the discrete labels $y \in \{m = 1, \dots, M\}$. In the context of semi-
 130 supervised node classification, we have access to labels for only a subset of nodes, which we denote
 131 as $y_o = [y_1, \dots, y_{n'}]^\top$, where n' represents the number of observed nodes. The remaining nodes
 132 have no assigned labels and are denoted as $y_u = [y_{n'+1}, \dots, y_n]$. In our setting, the assigned labels
 133 are assumed to be correct, but there might be errors or inaccuracies in the observed edges or node
 134 features due to noisy measurements. Our objective is to accurately predict the labels y_u for these
 135 unobserved nodes in the graph.

136 3.2 Node Embedding

137 We use a graph encoder to embed the node features, including both the nodal information and the
 138 corresponding graph structure. We emphasize that our framework is not tied to any specific graph
 139 model and offers flexibility in selecting graph encoders. In this study, we use Graph Convolutional
 140 Networks (GCNs) [1, 36] and Graph Attention Networks (GATs) [2] to encode both the node feature
 141 and the graph topology. For each node i , the graph encoder functions as a nonlinear transformation,
 142 denoted as ϕ_θ , taking the nodal features and the corresponding graph topology as input and returning
 143 their node embeddings, denoted by $\xi \in \Xi$. Formally,

$$\phi_\theta(\cdot, \mathcal{G}) : \mathcal{I} \rightarrow \Xi.$$

144 Here $\theta \in \Theta$ denotes the parameters specific to the model used in our framework.

145 It's important to note that the node embeddings obtained through the minimization of standard loss
 146 functions (e.g., cross-entropy loss) may not accurately capture the key feature patterns of the graph in
 147 our problem setting [3, 4, 5]. The presence of noise in the observed graph \mathcal{G} , compounded by the
 148 limited number of labeled nodes, can result in the model fitting too closely to the noise, being misled
 149 by incorrect gradients, and generalizing poorly. Such models can be vulnerable to outliers, which
 150 might compromise the quality of the learned representations.

151 3.3 Distributionally Robust Graph Learning

152 In light of these challenges, one of the goals is to enhance the robustness of graph encoders directly
 153 within their embedding space. This robustification process ensures that the embeddings better align
 154 with the underlying graph structure, even in the presence of noise and perturbations.

155 We develop a novel graph learning framework
 156 that improves the node embeddings, resulting in
 157 more robust predictive performance, particularly
 158 when confronted with noisy data. Figure 2 sum-
 159 marizes the architecture of the proposed model.

160 Our approach assumes that node embeddings
 161 share the same label m adhere to an underly-
 162 ing distribution ($\xi_i \sim P_m \in \mathcal{P}_m, \forall i : y_i = m$)
 163 within an uncertainty set \mathcal{P}_m encompassing all
 164 potential distributions P_m . However, obtaining
 165 a direct observation of this distribution is chal-
 166 lenging due to inaccuracies in the nodal or topo-
 167 logical information, and any changes in node
 168 embeddings will influence their corresponding
 169 distributions. The key idea of our proposed
 170 framework is to find the most robust node em-
 171 beddings, parameterized by θ , that minimize
 172 the worst-case risk over the uncertainty set \mathcal{P}_m
 173 in the probability simplex $\Delta_M = \{\pi \in \mathbb{R}_+^M : \sum_{m=1}^M \pi_m = 1\}$. This results in the definition
 174 of our distributionally robust minimax problem:
 175

$$176 \min_{\pi \in \Delta_M} \max_{\substack{P_m \in \mathcal{P}_m \\ 1 \leq m \leq M}} \Psi(\pi; P_1, \dots, P_M), \quad (1)$$

177 where Ψ is the risk function of a classifier π . We define the risk function as the summation of error
 178 probabilities under each class, i.e., $\Psi(\pi; P_1, \dots, P_M) := \sum_{m=1}^M \mathbb{E}_{\xi \sim P_m} [1 - \pi_m(\xi)]$. We note that
 179 the optimal solution P_1^*, \dots, P_M^* to the inner maximization problem is known as the *least favorable*
 180 *distributions* (LFDs) in statistics literature [9, 37]. The risk associated with these distributions is
 181 considered the worst-case risk [9].

182 As shown in Figure 3, we choose the uncertainty set \mathcal{P}_m to be a Wasserstein ball of radius ϑ_m
 183 centered at the empirical distribution \hat{P}_m :

$$184 \mathcal{P}_m := \{P_m \in \mathcal{P}(\Xi) : \mathcal{W}_1(P_m, \hat{P}_m) \leq \vartheta_m\}, \quad (2)$$

184 where $\mathcal{P}(\Xi)$ denotes the set of all probability distributions on Ξ . The Wasserstein distance of order
 185 one, \mathcal{W}_1 , is defined as $\mathcal{W}_1(P, P') := \min_{\gamma} \mathbb{E}_{(\xi, \xi') \sim \gamma} [c(\xi, \xi')]$, where $c(u, v)$ is some cost function
 186 transferring from u to v , $c(u, v) \geq 0$. The empirical distribution \hat{P}_m is represented by the Dirac point
 187 mass, denoted as:

$$188 \hat{P}_m := \frac{1}{|\{i : y_i = m\}|} \sum_{i=1}^n \delta_{\xi_i} \mathbb{1}\{y_i = m\}, m = 1, \dots, M, \quad (3)$$

188 Here, δ refers to the Dirac delta function, $|\cdot|$ represents the cardinality of a set, and $\mathbb{1}$ denotes the
 189 indicator function.

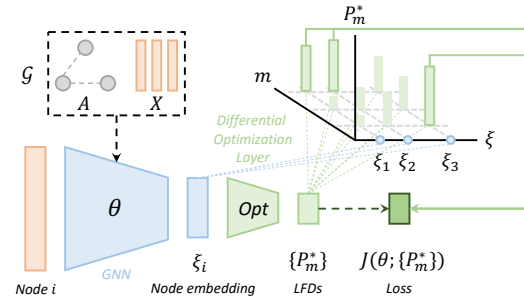


Figure 2: The architecture of the proposed framework consists of two cohesive modules: (1) a graph encoder parameterized by θ , which produces the node embedding ξ given the graph information \mathcal{G} ; (2) a differential optimization layer, which generates the corresponding least favorable distributions (LFDs) $\{P_m^*\}$ for ξ by solving the convex optimization defined in (4).

Algorithm 1 Learning algorithm of DRGL

```

1: Input:  $\mathcal{G} = (X, A)$ ;  $y_o = [y_1, \dots, y_{n'}]^\top$ ;  $\theta_0$ ;
2: Output:  $\theta_T$ 
3: for  $t \leftarrow [0 \dots T]$  do
4:    $\mathcal{L}(\theta_t) = 0$ ;
5:   for each mini-set do
6:     Compute the node embeddings  $\{\xi\}$  given  $\theta_t$  for all labeled nodes in the mini-set;
7:     Calculate  $\hat{P}_1, \dots, \hat{P}_M$  given  $\{\xi\}$  and  $y_o$  using (3);
8:     Obtain LFDs  $P_1^*, \dots, P_M^*$  by solving (4) with DO given  $\{\xi\}$  and  $\hat{P}$ ;
9:      $\mathcal{L}(\theta_t) \leftarrow \mathcal{L}(\theta_t) + J(\theta_t; P^*)$  using (5);
10:  end for
11:   $\theta_{t+1} \leftarrow \theta_t - \alpha \nabla \mathcal{L}(\theta_t)$  ( $\alpha$  is the learning rate);
12: end for
13: return  $\theta_T$ 

```

190 As the original problem (1) entails an intractable infinite dimensional functional optimization, we
 191 follow [9, 10] and introduce the proposition below. This proposition reformulates the problem into
 192 a computationally tractable convex optimization problem due to our careful selection of the risk
 193 function and uncertainty sets. We note that this selection exploits the structure of the least favorable
 194 distributions yielding from Wasserstein uncertainty sets [9].

195 **Proposition 1** *For the uncertainty sets defined in (2), the least favorable distribution of problem (1)*
 196 *can be obtained by solving the following problem:*

$$\begin{aligned}
 & \min_{p_1, \dots, p_M \in \mathbb{R}_+^n} \sum_{i=1}^n \max_{1 \leq m \leq M} p_m^i \\
 & \text{subject to } \sum_{i=1}^n \sum_{j=1}^n \gamma_m^{i,j} c(\xi_i, \xi_j) \leq \vartheta_m \\
 & \sum_{i=1}^n \gamma_m^{i,j} = \hat{P}_m(\xi_j), \quad \sum_{j=1}^n \gamma_m^{i,j} = p_m^i, \\
 & \forall 1 \leq i, j \leq n, 1 \leq m \leq M.
 \end{aligned} \tag{4}$$

197 The decision variable $\gamma_m \in \mathbb{R}_+^{n \times n}$ can be viewed as a joint distribution on n empirical points with
 198 marginal distributions \hat{P}_m and P_m , represented by a vector $p_m \in \mathbb{R}_+^n$. The inequality constraint
 199 controls the Wasserstein distance between P_m and \hat{P}_m .

200 **Remark 1** *The maximization in (4) measures the margin between the maximum likelihood of ξ_i*
 201 *among all classes and the likelihood of the m -th class. Thus, the objective can be equivalently*
 202 *rewritten as the minimization of the total margin. When $M = 2$, the total margin reduces to the total*
 203 *variation distance.*

204 3.4 Model Estimation

205 The proposed learning method for robust node embeddings can be carried out in an end-to-end
 206 fashion. To propagate the error backward through the convex optimization problem described in (4)
 207 to the graph encoder, we adopt the idea of differentiable optimization (DO) [15, 16].

208 This approach enables us to differentiate through certain subclasses of convex optimization problems.
 209 In other words, the convex solver can be seen as a function that maps the data of the problem to its
 210 corresponding solutions, making it amenable to gradient-based learning techniques. Therefore, the
 211 learning objective can be jointly written as:

$$J(\theta; P_1^*, \dots, P_M^*) := \sum_{i=1}^{n'} \max_{1 \leq m \leq M} P_m^*(\phi_\theta(i, \mathcal{G})), \tag{5}$$

212 where $P_m^*(\cdot)$ can be regarded as the output layer of our model, which takes the node embeddings ξ_i
 213 as input and returns their LFDs by solving (4) with DO.

214 We apply the mini-batch stochastic gradient descent as summarized in Algorithm 1. It is necessary
 215 that each mini-batch, which is provided to the convex solver, includes at least one training sample
 216 from every class to maintain the integrity of the optimization and its ability to generate valid solutions.

217 4 Experiments

218 4.1 Experimental Setup

219 In our experiments, we use three widely-used
 220 data sets, including Cora [38], Citeseer [39], and
 221 Pubmed [40]. These data sets consist of citation
 222 networks among 2, 708, 3, 327, and 19, 717 sci-
 223 entific publications, respectively [41]. In these
 224 networks, each node represents a text document,
 225 and its feature vector corresponds to a bag-of-
 226 words representation. We primarily focus on
 227 few-shot learning tasks where each data set con-
 228 tains M -class and K training samples per class.
 229 For each class, we randomly select K labeled
 230 nodes for a K -shot low-data setting.

231 To assess the robustness of our method, we in-
 232 troduce random noise into these data sets in
 233 the following two ways: (1) We add Gaussian
 234 noise $\epsilon \sim \mathcal{N}(0, \sigma)$ to the node feature matrix
 235 X , where σ is set proportionally to the standard
 236 deviation of the bag-of-words representation of
 237 all nodes in each graph; (2) We randomly re-
 238 move or add a certain percentage r of edges in
 239 the graph. We also test the performance of DRGL in K -shot low-data setting without adding noise
 240 to the original graphs. We repeat each test three times with three different seeds and calculate the
 241 average accuracy.

242 As our framework is designed to be GNN-agnostic, we chose to conduct experiments using Graph
 243 Convolutional Networks (GCNs) [1] and Graph Attention Networks (GATs) [2]. The purpose of
 244 adopting two basic GNNs is to test the generalizability of our DRGL framework. The performance of
 245 these basic GNN models aslo serves as a reference point for evaluating the robustness enhancements
 246 achieved by DRGL in our experiments. Specifically, for GCN, we employed two graph convolutional
 247 layers with hidden dimensions of 16 to learn 16-dimensional node representations. In the case of
 248 GAT, we utilized two graph attentional operator layers, with the first layer producing 8 attention
 249 heads, and ultimately obtaining 16-dimensional node embeddings.

250 The implementation of DRGL was carried out in PyTorch. Implementing DRGL for node classification
 251 tasks involved two stages. In the first stage, DRGL is used to enhance a pre-trained GNN encoder
 252 for robustness, following the algorithm detailed in Algorithm 1. In the second stage, we train a
 253 classification model based on either the learned embeddings ξ or the corresponding Least Favorable
 254 Distributions (LFDs) P_m^* of ξ . Kernel density estimation methods or k -nearest neighbors (k -NN) could
 255 also be employed to estimate the predicted probability based on the LFDs of the node embeddings.
 256 In our experiments, we evaluated our framework using two classifiers: a shallow (2-layer) neural
 257 network with Softmax output based on the learned embedding ξ and a weighted k -NN based on
 258 LFDs.

259 4.2 Robust GNN Baselines

260 To assess the effectiveness of DRGL, we conducted a comparative evaluation against baseline GNNs
 261 and state-of-the-art defence methods. An overview of these defense methods is provided below:

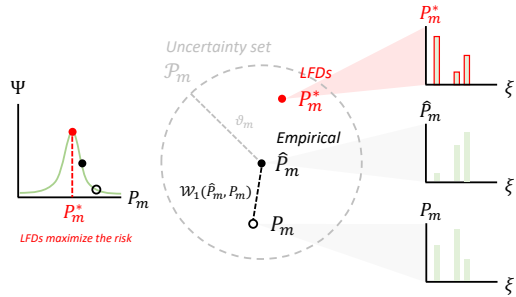


Figure 3: The minimax problem (1) aims to find the least favorable distributions (LFDs) by searching the optimal P_m in the uncertainty set \mathcal{P}_m that maximizes the risk Ψ . The uncertainty set is defined by Wasserstein distance.

Table 1: Model performances with Gaussian noise in node features ($K = 5$).

Models	Cora (M = 7)		Citeseer (M = 6)		Pubmed (M = 3)	
	σ					
LP	47.70	47.70	21.73	21.73	28.90	28.90
GCN+ Softmax	66.17	52.43	37.83	30.35	63.43	57.83
GAT+ Softmax	<u>66.48</u>	<u>59.03</u>	65.90	<u>60.43</u>	63.73	58.70
ProGNN	63.53	53.60	43.57	43.90	OOM	OOM
RGCN	61.43	52.87	45.93	36.13	<u>64.10</u>	59.03
GCN-SVD	55.27	55.27	33.57	31.53	49.80	49.80
GCN _{DRGL} + Softmax	66.20	54.40	39.80	34.90	63.20	<u>59.70</u>
GAT _{DRGL} + Softmax	67.83	59.40	<u>65.60</u>	61.27	64.57	59.77

Table 2: Model performances with random edge removal ($K = 5$).

Models	Cora (M = 7)		Citeseer (M = 6)		Pubmed (M = 3)	
	r					
LP	41.67	32.10	17.83	11.60	26.53	23.53
GCN+ Softmax	<u>67.30</u>	50.83	39.25	35.70	<u>66.07</u>	<u>62.07</u>
GAT+ Softmax	67.28	<u>57.08</u>	<u>64.07</u>	<u>60.13</u>	64.80	61.96
ProGNN	63.70	56.60	46.73	46.73	OOM	OOM
RGCN	62.76	56.33	47.00	43.77	63.93	58.20
GCN-SVD	32.53	27.80	32.53	27.80	47.40	44.37
GCN _{DRGL} + Softmax	66.70	52.15	40.15	41.80	66.63	64.30
GAT _{DRGL} + Softmax	69.02	58.88	65.93	61.20	65.93	60.53

262 **GCN-SVD** [29]: This approach is a preprocessing method that suggests robustifying GCN with
 263 the low-rank approximation of the perturbed graph. It was originally proposed to defend against the
 264 *netattack* attack [42], which includes perturbations in both node features and edges;

265 **RGCN** [30]: RGCN adopts a strategy of modeling node representations as Gaussian distributions to
 266 mitigate the impact of adversarial structural attacks. Additionally, it employs an attention mechanism
 267 to penalize nodes with high variance, enhancing robustness;

268 **ProGNN** [24]: ProGNN jointly learns a structural graph and a robust GNN model guided by
 269 intrinsic graph properties such as low-rank and sparsity in the graph structure, as well as the tendency
 270 for adjacent nodes to exhibit similar features. Since ProGNN [30] was originally introduced as a
 271 defense mechanism against graph structural attacks and noises, it can be readily extended and tested
 272 under our low-data graph noise scenarios.

273 4.3 Main Results

274 We evaluated the node classification accuracy of various methods in different noisy scenarios,
 275 including random noise, random edge removal or addition, and the standard low-data case. The
 276 best-performing method is highlighted in bold in the result tables, while the second-best is underlined.
 277 We note that the ProGNN method encountered out-of-memory (OOM) issues when trained with
 278 Pubmed using the official implementation, despite running on a 24GB RTX 4090 GPU. In such
 279 cases, its results are marked as OOM in the tables. Based on the experimental results, we make the
 280 following observations regarding the performance of DRGL in different settings:

281 **Random Gaussian noise in node features** Table 1 presents the average classification accuracy,
 282 revealing that GCN and GAT models trained with DRGL consistently outperforms the standard models
 283 using the same classifiers and baseline robust graph methods across almost all scenarios. These find-
 284 ings underscore the effectiveness of DRGL in improving the embedding’s robustness when confronted
 285 with noise in node features within low-data settings.

286 **Random edge perturbations** In the edge removal setting, Table 2 illustrates that both GCN and GAT
 287 exhibit significant improvements when enhanced with DRGL, particularly when a higher percentage of
 288 edges are removed. Their performance also consistently surpasses that of benchmark robust baselines,
 289 highlighting the effectiveness of DRGL in strengthening GNNs against missing edges. In the case of

Table 3: Model performances with random edge addition ($K = 5$).

Models	Cora (M = 7)		Citeseer (M = 6)		Pubmed (M = 3)	
	20%	50%	20%	50%	20%	50%
LP	47.70	47.70	21.73	21.73	28.90	28.90
GCN+ Softmax	61.73	55.03	47.67	39.83	59.27	60.07
GAT+ Softmax	68.93	63.00	47.20	44.27	63.70	<u>59.53</u>
ProGNN	66.63	56.60	43.00	38.00	OOM	OOM
RGCN	57.23	50.97	42.63	39.20	61.76	56.00
GCN-SVD	47.80	44.37	32.06	31.30	51.90	50.07
GCN _{DRGL} + Softmax	62.63	55.90	47.67	40.26	60.56	59.80
GAT _{DRGL} + Softmax	<u>67.60</u>	<u>59.13</u>	47.30	<u>43.97</u>	<u>62.30</u>	59.03

Table 4: Model performances without graph noise.

Models	Cora (M = 7)		Citeseer (M = 6)		Pubmed (M = 3)	
	K = 5	K = 10	K = 5	K = 10	K = 5	K = 10
LP	47.70	52.80	21.73	28.33	28.90	38.63
GCN+ k -NN	36.20	66.50	28.50	52.47	45.57	68.27
GCN+ Softmax	63.47	70.67	45.37	<u>57.83</u>	67.10	71.27
ProGNN	71.20	75.80	45.10	54.53	<u>67.70</u>	OOM
RGCN	<u>66.13</u>	71.23	50.00	53.00	66.43	70.53
GCN-SVD	52.67	60.43	33.57	34.34	52.17	57.43
GCN _{DRGL} + k -NN	44.60	68.23	34.07	53.30	51.10	68.67
GCN _{DRGL} + Softmax	<u>66.13</u>	<u>72.60</u>	<u>49.83</u>	59.33	67.83	72.20

290 edge addition, while Table 3 demonstrates a more moderate impact of DRGL on GATs, it still improves
 291 performance of GCN in most settings and outperforms robust methods.

292 **Standard few-shot setting** In this setting, each node class is provided with a limited number
 293 of labeled samples, along with the original features and edges. The results presented in Table 4
 294 demonstrate DRGL improves GCN in across all settings. This suggests that while DRGL is primarily
 295 designed to enhance robustness in noisy graph scenarios, it also enhances model performance in
 296 standard few-shot scenarios. We complemented the Softmax classifier with a k -NN classifier to
 297 further support this observation, where DRGL demonstrates similar improvement.

298 4.4 Learned Embeddings and Uncertainty

299 To gain a more intuitive understanding of the learned embedding space produced by DRGL, we
 300 conducted an additional ablation study. In this study, we set the output of the learned node embeddings
 301 to be two-dimensional and visualize them as scatter plots. Figure 4 and Figure 5 give real examples
 302 using three classes of nodes from Cora data set [38]. In these figures, the training data points are
 303 represented by large dots, while the testing data points are represented by small dots. The color of the
 304 dots corresponds to their true categories, providing a visual reference. The color depth of the regions
 305 suggests the likelihood of a sample being classified into the predicted category.

306 Figure 4 presents a comparison between the learned node embeddings and those generated by the
 307 baseline GCN used for DRGL in various noisy scenarios. We can observe that the between-class distance
 308 demonstrates a slight increase in contrast to the vanilla GCN. Conversely, the intra-class distance shows
 309 a minor reduction, as denoted by the denser distribution of dots with matching colors. It is worth
 310 mentioning that such differences in distribution can become more prominent in higher dimensional
 311 embedding spaces. Although seemingly subtle, this shift in the distribution of embeddings due to
 312 DRGL plays a pivotal role in significantly improving the accuracy of the classification outcomes,
 313 indicated by the accuracy displayed in the lower right corner and results in Section 4.3.

314 Figure 5 showcases the visual representation of the LFDs generated by DRGL. The shades of grey
 315 in the visualization are obtained using kernel density estimation, where darker shades indicate a
 316 higher level of uncertainty between the different categories. These LFDs serve as a valuable tool for
 317 uncertainty quantification (UQ), enabling us to pinpoint the data points that are susceptible to the
 318 greatest impact in worst-case scenarios. Let $\tilde{p}^* = [\tilde{p}_1^*, \dots, \tilde{p}_M^*]$ denote the predicted probabilities
 319 of a classifier built on learned LFDs of a graph. The predictive uncertainty can be expressed using

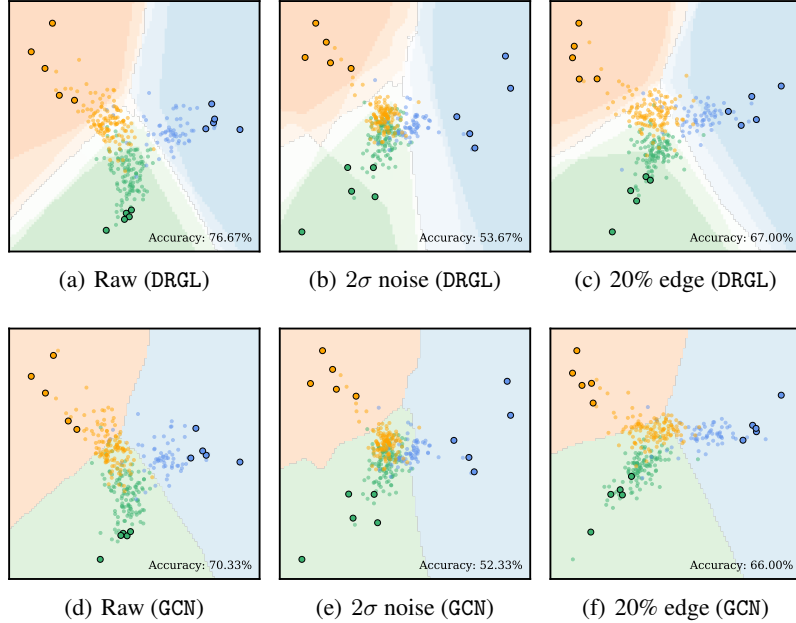


Figure 4: The impact of noise on the learned feature spaces. (a) and (b) show the embeddings from graphs without noise; (b) and (e) show the embeddings when the graphs have nodal features with 2σ noise; and (c) and (f) present the representations from graphs where 20% of the edges are removed.

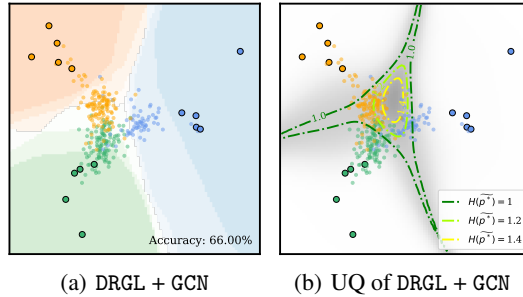


Figure 5: The learned embeddings and the uncertainty.

320 entropy [43, 44, 45]:

$$H(\tilde{p}^*) = - \sum_{m=1}^M \tilde{p}_m^* \log \tilde{p}_m^* \quad (6)$$

321 The result highlights the improved capability of our approach in capturing and quantifying uncertainty
 322 compared with the vanilla GCN method.

323 5 Conclusion

324 To address the challenges posed by noisy graphs, we proposed a novel GNN-agnostic framework,
 325 DRGL, that enhances the robustness of node embeddings and predictive performance. Our proposed
 326 framework improves model robustness by accounting for the uncertainties arising from data noise
 327 within the graph, leading to substantial improvements over state-of-the-art baselines on different
 328 benchmark datasets in various few-shot noisy graph settings. Extending this framework to other
 329 graph representation models beyond GCN and GATs and citation network data sets will be left to our
 330 future work.

331 **References**

- 332 [1] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
333 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 334 [2] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
335 Bengio. Graph attention networks, 2018.
- 336 [3] Jintang Li, Bingzhe Wu, Chengbin Hou, Guoji Fu, Yatao Bian, Liang Chen, Junzhou Huang,
337 and Zibin Zheng. Recent advances in reliable deep graph learning: Inherent noise, distribution
338 shift, and adversarial attack. *arXiv e-prints*, pages arXiv–2202, 2022.
- 339 [4] Yanqiao Zhu, Weizhi Xu, Jinghao Zhang, Yuanqi Du, Jieyu Zhang, Qiang Liu, Carl Yang, and
340 Shu Wu. A survey on graph structure learning: Progress and opportunities. *arXiv e-prints*,
341 pages arXiv–2103, 2021.
- 342 [5] Dongsheng Luo, Wei Cheng, Wenchao Yu, Bo Zong, Jingchao Ni, Haifeng Chen, and Xiang
343 Zhang. Learning to drop: Robust graph neural network via topological denoising. In *Proceedings*
344 *of the 14th ACM International Conference on Web Search and Data Mining, WSDM '21*, page
345 779–787, New York, NY, USA, 2021. Association for Computing Machinery.
- 346 [6] Kaize Ding, Elnaz Nouri, Guoqing Zheng, Huan Liu, and Ryen White. Toward robust graph
347 semi-supervised learning against extreme data scarcity. *arXiv preprint arXiv:2208.12422*, 2022.
- 348 [7] John C Duchi and Hongseok Namkoong. Learning models with uniform performance via
349 distributionally robust optimization. *The Annals of Statistics*, 49(3):1378–1406, 2021.
- 350 [8] Shai Shalev-Shwartz and Yonatan Wexler. Minimizing the maximal loss: How and why. In
351 *International Conference on Machine Learning*, pages 793–801. PMLR, 2016.
- 352 [9] Rui Gao, Liyan Xie, Yao Xie, and Huan Xu. Robust hypothesis testing using wasserstein
353 uncertainty sets. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
354 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran
355 Associates, Inc., 2018.
- 356 [10] Shixiang Zhu, Liyan Xie, Minghe Zhang, Rui Gao, and Yao Xie. Distributionally robust
357 weighted k-nearest neighbors. *Advances in Neural Information Processing Systems*, 35:29088–
358 29100, 2022.
- 359 [11] Xiang Zhang, Yinfei Xu, Qinghe Liu, Zhicheng Liu, Jian Lu, and Qiao Wang. Robust graph
360 learning under wasserstein uncertainty, 2021.
- 361 [12] Rizal Fathony, Ashkan Rezaei, Mohammad Ali Bashiri, Xinhua Zhang, and Brian Ziebart.
362 Distributionally robust graphical models. *Advances in Neural Information Processing Systems*,
363 31, 2018.
- 364 [13] Matthew Staib and Stefanie Jegelka. Distributionally robust deep learning as a generalization of
365 adversarial training. In *NIPS workshop on Machine Learning and Computer Security*, volume 3,
366 page 4, 2017.
- 367 [14] Soroosh Shafieezadeh Abadeh, Peyman M Mohajerin Esfahani, and Daniel Kuhn. Distribu-
368 tionally robust logistic regression. *Advances in Neural Information Processing Systems*, 28,
369 2015.
- 370 [15] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and Z. Kolter. Differentiable convex
371 optimization layers. In *Advances in Neural Information Processing Systems*, 2019.
- 372 [16] Brandon Amos and J. Zico Kolter. OptNet: Differentiable optimization as a layer in neural
373 networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70
374 of *Proceedings of Machine Learning Research*, pages 136–145. PMLR, 2017.
- 375 [17] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally
376 connected networks on graphs, 2014.

- 377 [18] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
378 on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing*
379 *Systems*, 2016.
- 380 [19] Petar Velickovic, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon
381 Hjelm. Deep graph infomax. *ICLR (Poster)*, 2(3):4, 2019.
- 382 [20] Kaize Ding, Zhe Xu, Hanghang Tong, and Huan Liu. Data augmentation for deep graph
383 learning: A survey. *ACM SIGKDD Explorations Newsletter*, 2022.
- 384 [21] Yu Chen, Lingfei Wu, and Mohammed Zaki. Iterative deep graph learning for graph neural
385 networks: Better and robust node embeddings. In H. Larochelle, M. Ranzato, R. Hadsell, M.F.
386 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
387 pages 19314–19326. Curran Associates, Inc., 2020.
- 388 [22] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph
389 structure learning for robust graph neural networks. In *Proceedings of the 26th ACM SIGKDD*
390 *International Conference on Knowledge Discovery & Data Mining*, KDD '20, page 66–74, New
391 York, NY, USA, 2020. Association for Computing Machinery.
- 392 [23] Wei Jin, Yaxing Li, Han Xu, Yiqi Wang, Shuiwang Ji, Charu Aggarwal, and Jiliang Tang.
393 Adversarial attacks and defenses on graphs. *SIGKDD Explor. Newsl.*, 22(2):19–34, 1 2021.
- 394 [24] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph
395 structure learning for robust graph neural networks. In *Proceedings of the 26th ACM SIGKDD*
396 *international conference on knowledge discovery & data mining*, pages 66–74, 2020.
- 397 [25] L. Wang, W. Yu, W. Wang, W. Cheng, W. Zhang, H. Zha, X. He, and H. Chen. Learning robust
398 representations with graph denoising policy network. In *2019 IEEE International Conference*
399 *on Data Mining (ICDM)*, pages 1378–1383, Los Alamitos, CA, USA, 12 2019. IEEE Computer
400 Society.
- 401 [26] Emanuele Rossi, Henry Kenlay, Maria I. Gorinova, Benjamin Paul Chamberlain, Xiaowen
402 Dong, and Michael Bronstein. On the unreasonable effectiveness of feature propagation in
403 learning on graphs with missing node features, 2022.
- 404 [27] Bo Jiang and Ziyang Zhang. Incomplete graph representation and learning via partial graph
405 neural networks. *arXiv preprint arXiv:2003.10130*, 2020.
- 406 [28] Cheng Zheng, Bo Zong, Wei Cheng, Dongjin Song, Jingchao Ni, Wenchao Yu, Haifeng Chen,
407 and Wei Wang. Robust graph representation learning via neural sparsification. In Hal Daumé
408 III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine*
409 *Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 11458–11468.
410 PMLR, 7 2020.
- 411 [29] Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. All
412 you need is low (rank) defending against adversarial attacks on graphs. In *Proceedings of the*
413 *13th International Conference on Web Search and Data Mining*, pages 169–177, 2020.
- 414 [30] Dingyuan Zhu, Ziwei Zhang, Peng Cui, and Wenwu Zhu. Robust graph convolutional networks
415 against adversarial attacks. In *Proceedings of the 25th ACM SIGKDD international conference*
416 *on knowledge discovery & data mining*, pages 1399–1407, 2019.
- 417 [31] Alireza Sadeghi, Meng Ma, Bingcong Li, and Georgios B. Giannakis. Distributionally robust
418 semi-supervised learning over graphs, 2021.
- 419 [32] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad
420 Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U. Rajendra Acharya, Vladimir
421 Makarenkov, and Saeid Nahavandi. A review of uncertainty quantification in deep learning:
422 Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- 423 [33] Edmon Begoli, Tanmoy Bhattacharya, and Dimitri Kusnezov. The need for uncertainty quanti-
424 fication in machine-assisted medical decision making. *Nature Machine Intelligence*, 1(1):20–23,
425 2019.

- 426 [34] Seongok Ryu, Yongchan Kwon, and Woo Youn Kim. A bayesian graph convolutional network
427 for reliable prediction of molecular properties with uncertainty quantification. *Chemical science*,
428 10(36):8438–8446, 2019.
- 429 [35] Yao Zhang et al. Bayesian semi-supervised learning for uncertainty-calibrated prediction of
430 molecular properties and active learning. *Chemical science*, 10(35):8154–8163, 2019.
- 431 [36] Qimai Li, Zhichao Han, and Xiao-ming Wu. Deeper insights into graph convolutional networks
432 for semi-supervised learning. *Proceedings of the AAAI Conference on Artificial Intelligence*,
433 32(1), 4 2018.
- 434 [37] Erich Leo Lehmann, Joseph P Romano, and George Casella. *Testing statistical hypotheses*,
435 volume 3. Springer, 1986.
- 436 [38] Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating
437 the construction of internet portals with machine learning. *Information Retrieval*, 3:127–163,
438 2000.
- 439 [39] C Lee Giles, Kurt D Bollacker, and Steve Lawrence. Citeseer: An automatic citation indexing
440 system. In *Proceedings of the third ACM conference on Digital libraries*, pages 89–98, 1998.
- 441 [40] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
442 Rad. Collective classification in network data. *AI Magazine*, 29(3):93, 9 2008.
- 443 [41] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning
444 with graph embeddings. In *International conference on machine learning*, pages 40–48. PMLR,
445 2016.
- 446 [42] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural
447 networks for graph data. In *Proceedings of the 24th ACM SIGKDD International Conference*
448 *on Knowledge Discovery & Data Mining*, KDD '18, page 2847–2856, New York, NY, USA,
449 2018. Association for Computing Machinery.
- 450 [43] C. E. Shannon. Communication theory of secrecy systems. *The Bell System Technical Journal*,
451 28(4):656–715, 1949.
- 452 [44] Alireza Namdari and Zhaojun Li. A review of entropy measures for uncertainty quantification
453 of stochastic processes. *Advances in Mechanical Engineering*, 11(6):1687814019857350, 2019.
- 454 [45] Stephan Schwill. Entropy analysis of financial time series. *arXiv preprint arXiv:1807.09423*,
455 2018.