Adversarially Robust Graph Classification: A Pooling-Based Defense Framework

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

Paper under double-blind review

Abstract

Graph Neural Networks (GNNs) have shown great success across various domains but remain vulnerable to adversarial attacks. While most defense methodology focuses on node classification and enhancing robustness during training, this work shifts the focus to graph classification and inference-time defenses. We theoretically show that the final pooling operation, that is required for graph-level tasks, can have an impact on the graph classifier's underlying robustness. Based on this analysis, we propose a pre-pooling operation, called R-Pool (Robust-Pooling), which is based a novel filtering mechanism using Gaussian Mixture Models (GMMs) to detect and exclude nodes heavily impacted by attacks, thereby enhancing robustness at inference time. Our framework can be used with any pooling operation and any underlying model, and does not require re-training the model nor adapting its architecture. Our experiments demonstrate that this approach effectively mitigates adversarial effects while maintaining a balance between clean and attacked accuracy. Through extensive evaluations on state-of-the-art adversarial attacks, we show that the proposed framework significantly improves the robustness of the underlying GNNs in graph classification tasks compared to other available post-hoc defense methods.

028 1 INTRODUCTION

Graph Neural Networks (GNNs) (Kipf & Welling, 2017; Xu et al., 2019b; Veličković et al., 2018) 031 have emerged as a robust framework for learning representations of nodes and graphs, demonstrating notable success across a wide range of real-world applications. These models, which generalize neu-033 ral network architectures to handle graph-structured data, have been successfully applied in critical 034 domains such as protein function prediction (Kearnes et al., 2016), antibiotic resistance prediction (Qabel et al., 2022), session-based recommendation systems (Wu et al., 2019b) and lately tabular data (Alkhatib et al., 2023). Despite their success, recent studies (Günnemann, 2022) have high-036 lighted the vulnerability of GNNs to adversarial perturbations, small and deliberately introduced 037 changes in the adjacency matrix or node features that can lead to incorrect predictions. These attacks pose significant challenges for the reliable deployment of GNNs, particularly in critical sectors such as healthcare. In response, a growing body of research has focused on characterizing these vul-040 nerabilities through various adversarial attack strategies (Dai et al., 2018; Zügner et al., 2018), while 041 simultaneously advancing defense mechanisms (Wu et al., 2019a; Zhang & Zitnik, 2020) to mitigate 042 these risks and improve the robustness of GNN models. 043

Existing research predominantly focuses on node-classification tasks, with comparatively limited 044 attention given to robustness analysis in the context of graph-classification tasks (Jin et al., 2021). Although both node-level and graph-level tasks rely on the message-passing propagation mecha-046 nism (Gilmer et al., 2017), graph-level tasks introduce additional complexity, requiring the gener-047 ation of a global graph representation for input graphs of varying sizes and topologies. In such 048 tasks, the pooling mechanism (Cai et al., 2021; Duvenaud et al., 2015) plays a critical role, as it consolidates the node representations obtained from the propagation step into a smaller graph or a single vector. Pooling operations generally fall into two main categories (Liu et al., 2022). The first, 051 hierarchical pooling (Cai et al., 2021), incrementally reduces the graph size, ultimately producing the graph representation for downstream tasks. The second, flat pooling (Duvenaud et al., 2015), 052 directly constructs a graph-level representation in a single step by aggregating node representations, as seen in techniques like sum-pooling. This work specifically focuses on the latter category, examining commonly used flat pooling operations such as sum, average, and max pooling, which are favored in the graph literature for their simplicity and effectiveness.

While most existing approaches to adversarial defense concentrate on modifying the message-057 passing scheme-through mechanisms such as attention (Zhang & Zitnik, 2020), adjusting weights (ABBAHADDOU et al., 2024) or pre-processing the adjacency matrix (Wu et al., 2019a), our work takes a different approach by theoretically investigating the impact of the pooling operation 060 on the overall robustness of the model. Specifically, we aim to understand how the choice of pooling 061 method influences adversarial robustness. Additionally, the majority of methods focus on enhancing 062 robustness during the training phase, with relatively few addressing robustness at inference time. 063 This latter perspective is becoming increasingly important with the growing interest in foundation 064 models, where the model is pre-trained, and the goal is to adapt it to downstream tasks while ensuring its adversarial resilience. To the best of our knowledge, this work is the first to explore the effect 065 of pooling operations on the adversarial vulnerability of models in graph classification, specifically 066 at inference time. 067

068 In this work, we begin by introducing the concept of adversarial attacks in the context of graph-069 based models, followed by a formalization of a graph classifier's robustness. This formal framework, which derives an upper bound on the model's expected adversarial risk, enables us to quantify the 071 model's vulnerability within a defined neighborhood. We then apply this formalization to conduct a theoretical analysis of the robustness of various pooling operations, extracting their corresponding 072 upper bounds. Building on these theoretical insights, we propose Robust-Pooling ("R-Pool"), an 073 incremental component that can be integrated with any existing pooling operation to enhance the 074 model's robustness at inference time. Our theoretical findings suggest that certain nodes in the graph 075 are more prone to accumulating adversarial perturbations in their final representations. To address 076 this, the proposed framework aims to mitigate the impact of these perturbations by filtering out such 077 nodes. Specifically, the proposed method consists of fitting a Gaussian Mixture Models (GMMs) to node representations and using an out-of-distribution score (Morteza & Li, 2022) to detect nodes 079 heavily impacted by adversarial attacks, excluding them prior to the pooling operation. Although this filtering may lead to a loss of information for the downstream classification task, our experiments 081 demonstrate that with an appropriate threshold, a balance between clean and adversarial accuracy can be achieved. Finally, we validate our theoretical contributions through extensive experimental evaluations, using state-of-the-art graph-classification adversarial attacks. The results confirm the 083 validity and practical value of the proposed framework. The contributions of this work can be 084 summarized as follows: 085

- We theoretically analyze how various pooling operations impact a model's underlying robustness in the case of graph classification.
- We then introduce "R-Pool", a novel filtering mechanism, based on fitting Gaussian Mixture Models (GMMs) and using an out-of-distribution score to detect and exclude nodes heavily impacted by adversarial attacks, thereby enhancing robustness during inference.
- We validate experimentally the efficiency of the proposed framework, demonstrating improved robustness against adversarial attacks in graph-classification tasks while maintaining a balance between clean and adversarial accuracy.

2 RELATED WORK

090

092

093

094 095 096

097

098 There has been growing interest in adversarial attacks targeting Graph Neural Networks 099 (GNNs) Zügner et al. (2018); Dai et al. (2018). Various attack methods have been proposed de-100 pending on the attack setting, which may assume full access to the model's architecture and training 101 data (white-box) or limit the attacker to model queries or surrogate models (black-box). Most attacks 102 frame the problem as an optimization task, with solutions ranging from PGD-based methods (Xu 103 et al., 2019a) to meta-gradient approaches (Zügner & Günnemann, 2019) and reinforcement learning 104 strategies (Dai et al., 2018). While much of the literature focuses on node classification, with some 105 methods like PGD adaptable to graph classification, few approaches are specifically tailored to this task. For instance, Wan et al. (2021) employ Bayesian optimization, and Dai et al. (2018) explore 106 several strategies, including reinforcement learning, gradient-based attacks, and genetic algorithms, 107 where graph modifications evolve through a population-based approach guided by a fitness function.

108 On the defense side, recent efforts have emerged to protect GNNs from adversarial attacks. Some 109 approaches focus on pre-processing the adjacency matrix to improve robustness. For example, Wu 110 et al. (2019a) use Jaccard similarity to filter out adversarial edges, while Entezari et al. (2020) apply 111 SVD decomposition to denoise the adjacency matrix before feeding it into the model. Other defenses 112 target the structure of the GNN itself. Zhang & Zitnik (2020) introduce edge pruning through an attention mechanism to remove vulnerable connections, while ABBAHADDOU et al. (2024) and 113 Ennadir et al. (2024) propose modifications to the message-passing scheme to reduce the impact 114 of adversarial perturbations. Moreover, adversarial training methods, such as those presented by 115 Gosch et al. (2024), have been developed to increase the model's resilience by explicitly training the 116 GNN with adversarially perturbed graphs. Additionally, a growing interest in exploring robustness 117 certificates (Zügner & Günnemann, 2019; Bojchevski & Günnemann, 2019) have emerged such as 118 randomized smoothing (Bojchevski et al., 2020). However, most of these defense techniques are 119 tailored to node classification tasks and have limited direct application to graph classification. 120

Our work takes a different direction by focusing specifically on the graph classification task, which has received comparatively less attention in adversarial defense literature. In particular, we investigate the role of pooling operations in determining the robustness of GNNs. Pooling, a key operation in graph classification models, has not been thoroughly studied in terms of its influence on adversarial vulnerability, especially during inference time. This paper aims to fill that gap by offering a theoretical analysis of how pooling affects model robustness. To our knowledge, this is the first work to explore this aspect in the context of graph classification, contributing a novel perspective to the ongoing research in adversarial attacks on GNNs.

128 129

130 131

132

150 151 152

3 PRELIMINARIES

Before continuing with our contribution, we introduce some fundamental concepts and notations.

133 Notation and Setup. Let G = (V, E) be a graph where V is its set of vertices and E its set of 134 edges. We denote by n = |V| and m = |E| the number of vertices and number of edges, respectively 135 and $\mathcal{N}(v) = \{u : (v, u) \in E\}$ the set of neighbors of a node $v \in V$. The degree of a node is equal 136 to its number of neighbors, i.e., $|\mathcal{N}(v)|$ for a node $v \in V$. A graph is commonly represented by 137 its adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, where the (i, j)-th element of the adjacency matrix is equal to the 138 weight of the edge between the *i*-th and *j*-th node of the graph and a weight of 0 in case the edge does not exist. In some settings, the nodes of a graph might be annotated with feature vectors. We 139 use $\mathbf{X} \in \mathbb{R}^{n \times D}$ to denote the node features where D is the feature dimensionality. The feature of 140 the *i*-th node of the graph corresponds to the *i*-th row of **X**. 141

142 143 Message Passing GNNs. A GNN model consists of a series of neighborhood aggregation layers 144 which use the graph structure and the nodes' feature vectors from the previous layer to generate 145 new representations for the nodes. Specifically, GNNs update nodes' feature vectors by aggregating 146 local neighborhood information. Suppose we have a GNN model that contains T neighborhood 147 aggregation layers. Let also $\mathbf{h}_v^{(0)}$ denote the initial feature vector of node v, i. e., the row of matrix 148 X that corresponds to node v. At each iteration (t > 0), the hidden state $\mathbf{h}_v^{(t)}$ of a node v is updated 149 as follows:

$$\begin{split} \mathbf{a}_{v}^{(t)} &= \mathsf{AGGREGATE}^{(t)} \Big(\big\{ \mathbf{h}_{u}^{(t-1)} \colon u \in \mathcal{N}(v) \big\} \Big) \\ \mathbf{h}_{v}^{(t)} &= \mathsf{COMBINE}^{(t)} \Big(\mathbf{h}_{v}^{(t-1)}, \mathbf{a}_{v}^{(t)} \Big), \end{split}$$

where AGGREGATE is a permutation invariant function that maps the feature vectors of the neighbors of a node v to an aggregated vector. This aggregated vector is passed along with the previous representation of v (i. e., $\mathbf{h}_v^{(t-1)}$) to the COMBINE function which combines those two vectors and produces the new representation of v.

Pooling Operation. After T iterations of neighborhood aggregation, to produce a graph-level representation, GNNs apply a permutation invariant readout function to the feature vectors of all nodes of the graph as follows: $(f_{1}, f_{2}) = 0$

$$\mathbf{h}_G = \mathsf{READOUT}\Big(\big\{\mathbf{h}_v^{(T)} \colon v \in V\big\}\Big),$$

Our study focuses on the Flat Pooling Family which directly generated a graph-level representation in one step (e. g., sum operator, mean operator and max operator).

165 166

167

168

170

171

4 ON THE ROBUSTNESS OF POOLING OPERATIONS

This work focuses on the robustness of the graph classification task within the broader context of graph representation learning, with a particular emphasis on understanding the impact of pooling operations on model robustness. We begin by discussing the concept of robustness in graph classifiers, followed by a formal mathematical definition of model robustness. Finally, we provide theoretical insights into the robustness of widely-used flat pooling operations, offering a deeper understanding of their behavior under adversarial attacks.

181 182

4.1 ADVERSARIAL GRAPH ROBUSTNESS

Ţ

Given an input graph (with its corresponding node attributes), graph-level classification aims to learn a function that predicts a property of interest related to the graph. Let's therefore consider our set of input graphs $\{(G_1, X_1, y_1), \ldots, (G_N, X_N, y_N)\} \in (\mathcal{G}, \mathcal{X}, \mathcal{Y})^N$ considered to be sampled from an underlying distribution \mathcal{D} defined on $(\mathcal{G}, \mathcal{X}, \mathcal{Y})$. Graph classification aims to find the graph-classifier $f: (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ minimizing the classification risk w.r.t \mathcal{D} , which is defined as:

$$R[f] := \mathbb{E}_{(G,X,y)\sim\mathcal{D}}[\mathbf{1}\{f(G,X)\neq y\}].$$

In this work, we focus on the black-box evasion attack setting, where we consider that the user cannot 183 access/modify the trained and static victim model $f : (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ or the training dataset. For our 184 theoretical insights, we follow the same definition as the one provided in the work, by adapting it to 185 the case of graph-classification. Specifically, let's consider an input graph $(G, X) \in (\mathcal{G} \times \mathcal{X})$ with its corresponding label $y \in \mathcal{Y}$, an adversarial attack aims to degrade the performance of the considered 187 victim model by finding a graph \hat{G} and its corresponding features \hat{X} within the the input graph's 188 neighborhood for which the predicted classification is different from the original classification of the 189 consider graph input (G, X). In this perspective, assessing the adversarial risk of a classifier consists 190 of analyzing the input graph's neighborhood by quantifying the "chances" (in terms of expectancy) 191 of finding an adversarial graph that is similar to input graph and for which the model's output is 192 different than the original output. For a given input graph's neighborhood defined by a threshold ϵ , 193 this expected risk robustness quantification can be formulated as follows:

194 195 196

197

202 203

$$\mathcal{R}_{\epsilon}[f] = \underset{\substack{(A,X)\sim\mathcal{D}\\(\tilde{A},\tilde{X})\in\mathcal{N}_{\epsilon}(A,X)}}{\mathbb{E}} [d_{\mathcal{Y}}(f(\tilde{A},\tilde{X}),f(A,X))],$$
(1)

where we consider $d_{\mathcal{Y}} = \|.\|_2$ as a distance within our output space \mathcal{Y} . Additionally, $\mathcal{N}_{\epsilon}(A, X) = \{(\tilde{A}, \tilde{X}) : d_{\mathcal{A}, \mathcal{X}}([A, X], [\tilde{A}, \tilde{X}]) < \epsilon\}$ denotes the input graph's considered neighborhood defined by our attack/perturbation budget ϵ . We simply consider the following graph distance that reflects both the structure and the node features:

$$d_{\mathcal{A},\mathcal{X}}([A,X],[\tilde{A},\tilde{X}]) = \min_{P \in \Pi} \{ \|A - P\tilde{A}P^T\|_2 + \|X - P\tilde{X}\|_2 \},\$$

with Π being the set of permutation matrices. In the case of un-attributed graphs, the previous distance resolves to using only the first part related to the adjacency matrices.

206 Typically, from the previous formulation, from a defense perspective, we aim to have the smallest 207 possible value, meaning that within the considered perturbation budget, the distance in term of output 208 of the attacked graph and the clean graph is not very big. Hence, we would expect the two of them 209 to have the same classification, which reflects the failure of the attack. Deriving the precise expected 210 distance, as defined in Equation (1), is challenging. However, an effective and more manageable 211 approach is to establish an upper-bound on this risk. By deriving such upper-bound, users can 212 gain a comprehensive understanding of the GNN's susceptibility to adversarial attacks and make 213 informed assessments of its robustness based on the specific task at hand. For example, in certain scenarios like social networks, where a limited number of successful attacks may not have severe 214 consequences, a larger upper-bound on the adversarial risk might be tolerable. While in other more 215 sensitive areas, such as financial applications, we need to aim for a much tighter upper-bound to control the confidence level of the adversarial risk. From this perspective, Definition 4.1 introduces
 the notion of a GNN's robustness.

218 219 219 220 221 Definition 4.1. (Adversarial Robustness). The graph-based function $f : (\mathcal{A}, \mathcal{X}) \to \mathcal{Y}$ is said to be (ϵ, γ) – robust if its adversarial risk is upper-bounded i. e., $\mathcal{R}_{\epsilon}[f] \leq \gamma$ with respect to the chosen graph distances in the input and output measurable sets.

As previously precised by the work, the definition rather approaches the adversarial problem from an "average perspective", where we analyze the whole neighborhood rather than "worst-case", in which the focus is on a single adversary that results in the most harmful performance. Note that, *if f is* $(\epsilon, \gamma) -$ "*robust*", *then it is also* $(\epsilon, \gamma) -$ "*worst-case robust*". The complete difference between worst-case and average case have been thoroughly studied previously in the literature.

4.2 ON THE ROBUSTNESS OF FLAT POOLING OPERATIONS

Having established a formal definition and framework for robustness in the case of graph classification, we now apply this framework to examine the robustness of commonly used pooling operations. In this context, we focus on two well-established message-passing models, which are instances from the general perspective provided in Section 3 : Graph Convolutional Networks (GCNs) and Graph Isomorphism Networks (GINs).

To improve the quality of graph-level representations, various pooling mechanisms have been proposed in the literature, which can be broadly categorized into Flat Pooling and Hierarchical Pooling. In this study, we focus on three widely used flat pooling methods, typically the Sum, Average, and Max pooling (Duvenaud et al., 2015). These methods directly generate a graph-level representation in a single step by aggregating node embeddings. While our theoretical analysis focuses on these specific pooling operations, we consider that our theoretical analysis can be easily expanded to take into account other available pooling operations within the Flat pooling sub-family.

Our analysis covers both node feature-based adversarial attacks, which modify node attributes to change model predictions, and structural perturbations, where adversaries alter the graph's structure by adding or removing edges to achieve their goal.

Theorem 4.2. Let $f : (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ denote a graph-based function composed of L GCN layers, where the weight matrix of the *i*-th layer is denoted by $W^{(i)}$. Further, let $d^{0,1}$ be a graph distance. For adversarial attacks only targeting node features of the input graph, with a budget ϵ , we have:

• If f is Max-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

 $\gamma = \prod_{l=1}^{L} | W^{(l)} | \max_{u \in \mathcal{V}} \hat{w_u} \epsilon.$

• If f is Sum-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

$$\gamma = \prod_{l=1}^{L} \mid W^{(l)} \mid \sum_{u \in \mathcal{V}} \hat{w_u} \epsilon.$$

255 256

248

249

250 251

253 254

227

228

257

258 259

260

• If f is Average-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

 $\gamma = \frac{\epsilon}{|V|} \prod_{l=1}^{L} | W^{(l)} | \sum_{u \in \mathcal{V}} \hat{w_u},$

with \hat{w}_u denoting the sum of normalized walks of length (L-1) starting from node u.

261 Theorem 4.2 provides the respective upper bounds on the different considered pooling operations. 262 The derived bound depends on two main factors. The first factor relates to the model itself, where 263 we observe the norm of the weights. The second factor is associated with the underlying structure of 264 the input graph. When subject to targeted attacks, we find that the Average and Sum pooling oper-265 ations consistently demonstrate vulnerability to adversarial manipulation. In contrast, Max pooling 266 exhibits variable robustness. Specifically, it can be more resilient if the targeted attack happens to 267 affect a node that is not the most connected (in terms of walks). However, it can become significantly more vulnerable if the attack targets the most connected node, which is likely in scenarios 268 such as gradient-based attacks, as these focus on nodes with the highest impact. In the case of more 269 global, non-targeted attacks, the effect on Max pooling can be expected to be less pronounced than

270 on Average and Sum pooling. This is because Max pooling accumulates the effect through only 271 one node, while the others accumulate the effects of all nodes. Overall, our analysis suggests that 272 there isn't a universally more robust pooling operation. Rather, the relative robustness depends on 273 the specific attack setting and the different types of attacks employed. This underscores the im-274 portance of considering the anticipated threat model when selecting an appropriate pooling strategy for GNNs in adversarial environments. We note that while this latter study is focusing on GCN, a 275 similar approach can be applied for GINs. 276

277 **Theorem 4.3.** Let $f : (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ be composed of L GIN-layers (with its internal parameter 278 $\zeta = 0$) and let $W^{(i)}$ denote the weight matrix of the *i*-th MLP layer. We consider the input node feature space to be bounded i. e., $||X||_2 < B$ for some $B \in \mathbb{R}$. For node feature-based attacks, with 279 a budget ϵ , the function f is $(d^{0,1}, \epsilon)$ -robust with 280

• If f is Max-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

 $\gamma = \prod_{l=1}^{L} \|W^{(l)}\| \left[B \times L \times \max_{u \in V} \deg(u) + \epsilon\right].$

• If f is Sum-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

 $\gamma = \prod_{l=1}^{L} \left\| W^{(l)} \right\| \left[2B \times L \times |E| + |V|\epsilon \right].$

• If f is Average-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

$$\gamma = \prod_{l=1}^{L} \left\| W^{(l)} \right\| \left[\frac{2B \times L \times |E|}{|V|} + \epsilon \right],$$

with |E| being the number of edges and |V| the number of nodes. 293

Similar to the case of GCN, the computed upper-bounds reveal the impact of the pooling oper-295 ation on the model's underlying adversarial robustness depending on the input graph's structure. 296 Specifically, for the sum pooling, the perturbation scales with the total number of edges and nodes, 297 indicating that larger and denser graphs are more susceptible to adversarial feature perturbations. 298 In contrast, average pooling normalizes this effect by the number of nodes, potentially reducing the 299 overall sensitivity in larger graphs. In the case of max pooling, it seems that it depends on the max-300 imum node degree, suggesting that graphs with highly connected nodes may be more vulnerable to attacks targeting those nodes. Overall, both in the case of GCN and GIN, the theoretical results show that choosing appropriate pooling strategies can have an effect of the model's robustness. 302

303 304

305

301

281

282 283

284 285

287 288

5 **ROBUST POOLING THROUGH FILTERING**

306 From Section 4.2, it appears that for different pooling operations, the upper bound is dependent 307 on the underlying graph's structure. This observation aligns with the intuitive understanding of 308 message-passing mechanisms, which propagate information within node neighborhoods. In this 309 context, when an attack is injected into certain nodes, its effect propagates to others. Consequently, nodes that are more "influential" in terms of their interactions with other nodes have a higher prob-310 ability of being affected by the attack, which in turn impacts nodes within their neighborhood, 311 resulting in a compounded effect. 312

313 These insights suggest a direct and straightforward defense method: discarding the more highly 314 connected nodes within the graph to mitigate their cumulative attack effect. Similar approaches have 315 been proposed in the literature, employing pre-processing techniques such as Jaccard Similarity or SVD decomposition. However, these methods have primarily been applied to node classification 316 tasks, where graphs are typically large, and discarding a number of nodes does not necessarily result 317 in significant information loss. In contrast, graph classification often deals with smaller graphs, 318 where highly connected nodes are crucial for the downstream task, and their removal could lead to 319 substantial information loss, equivalent to reduced clean accuracy which is very important since a 320 priori we do not know if the graph is attacked or not. Therefore, finding the right trade-off between 321 clean and attacked accuracy is important. 322

Given these considerations, we propose an alternative solution based on post-message-passing filter-323 ing. Specifically, after propagating the messages (e.g., using convolutions in the case of GCN (Kipf & Welling, 2017)), the result is a matrix X ∈ ℝ^{n,e}, where n is the number of nodes and e is the embedding dimension. We suppose that the resulting node representations follow a certain distribution
Consequently, when subject to adversarial attacks, some of these node representations may exhibit anomalous behavior, which is likely to be out-of-distribution with respect to D. This approach is particularly practical for small graphs, where retaining all nodes ensures that information relevant to the downstream task is preserved through the message-passing scheme.

330 Given a trained model, which we consider static, we introduce a filtering mechanism prior to the 331 pooling operation. Specifically, for an input graph G with n nodes, we consider its node represen-332 tation X and aim to partition the set of nodes into two subsets: a set of non-affected nodes \mathcal{N} and a 333 set of affected nodes \mathcal{M} . Assuming the distribution \mathcal{D} , the first set can be viewed as in-distribution 334 points and the second set as out-of-distribution points. Based on this perspective, we approximate the distribution \mathcal{D} using a finite mixture of k components, in particular with a Gaussian Mixture 335 Model (GMM), which serves as a universal approximator for a wide range of density functions. The 336 process involves estimating the GMM parameters, i.e., the component weight η_i , the mean μ_i and 337 covariance matrix Σ_j , for each $j = 1, \ldots, k$ for each input graph using the node representation 338 X, implemented via the Expectation-Maximization (EM) algorithm. After computing the optimal 339 parameter values, we adopt an approach similar to that proposed in (Morteza & Li, 2022), an OOD 340 detection method. For each node u with representation x_u , we compute a score, denoted as the 341 Gaussian mixture based energy measurement (GEM), which can be expressed as 342

$$s(u) = \log \sum_{j=1}^{k} \eta_j \exp(-\frac{1}{2}(x_u - \mu_j)^\top \Sigma_j (x_u - \mu_j)),$$

345 where x_u represents the node representation of node u. Given a threshold value λ , we employ the 346 GEM score to determine whether a node u has been subjected to an attack, consequently supporting 347 our decision to retain or discard it. Our implementation computes the threshold λ based on the quan-348 tile values of the calculated scores. This approach has demonstrated superior adaptability compared 349 to a rigid threshold, particularly given the heterogeneity in graph sizes prevalent in our datasets. 350 Formally, we classify a node u as $u \in \mathcal{N}$ based on its score s(u) comapred to λ . It is noteworthy 351 that (Morteza & Li, 2022) provide theoretical insights into the performance characteristics of the 352 GEM methodology, which is out of this work's scope. We additionally note that while in this work, 353 we have chosen to work with the GEM score, any other OOD score could be used. After processing all the nodes and deciding on the set of nodes to be kept \mathcal{N} and to discard \mathcal{M} , we proceed with 354 the classical pooling function to produce the final graph representation. We refer to the proposed 355 filtering scheme as R-pool and as it can be seen from its construction, it is adaptable to any pooling 356 e.g., sum, mean or max. 357

358 The primary advantage of our method lies in its model-agnostic nature, operating independently of the underlying Graph Neural Network (GNN) architecture. R-Pool accepts any node representation 359 X as input, regardless of whether it is generated by GCN, (Kipf & Welling, 2017), GIN (Xu et al., 360 2019b), GAT (Veličković et al., 2018), or any other method. It then performs adaptive node filtering 361 to produce a refined set of node representations. This universality ensures broad applicability across 362 various GNN frameworks without necessitating modifications to the original model architecture or 363 retraining procedures. Our approach is particularly relevant in the context of Foundation Models, 364 where the pre-trained nature of the model makes it prohibitively costly to re-train for adapting the message-passing scheme or implementing adversarial training. By operating at inference time, R-366 Pool offers a practical solution for enhancing robustness without the computational and resource 367 overhead associated with model re-training.

368 Complexity of the Method The main computational complexity of our method is concentrated in 369 the Expectation-Maximization algorithm used for estimating the GMM's parameters. The EM algo-370 rithm is an iterative process comprising two main steps: the Expectation (E) step and the Maximiza-371 tion (M) step. In the E-step, for each data point, we compute the probability of it being generated 372 by each component of the model, given the current parameter estimates. The M-step then updates 373 these parameters to maximize the expected log-likelihood based on the probabilities calculated in 374 the E-step. The algorithm's convergence rate and the quality of the resulting fit are partially depen-375 dent on the number of iterations. In our implementation and experimental results, we observed that a relatively small number of iterations (between 100 and 200) yielded satisfactory results (both in 376 terms of clean and attacked accuracy). A comprehensive time and complexity analysis is provided 377 in Appendix D.

378 6 EXPERIMENTAL RESULTS

In this section, we aim to validate the validity of our proposed R-Pool when subject to adversarial attacks and using real world datasets. We start by give details about the experimental settings we will be following and then we will report and analyze the performance of our proposed approach in comparison to other available methods.

3856.1 EXPERIMENTAL SETUP386

387 **Datasets.** Consistent with our theoretical analysis, this section focuses on the graph classification 388 task. We base our evaluation on the standard graph dataset derived from bioinformatics (PROTEINS, 389 NCI1) and from social networks (IMDB-BINARY) (Morris et al., 2020). We note that the social 390 network graphs are unlabeled, while all other graph datasets come with either node labels or node attributes. We take those labels/attributes into account when available and we otherwise use the 391 node's degree as its node features. To mitigate the impact of randomness during training, each 392 experiment was repeated 10 times, using the public train/validation/test splits provided in the work 393 (Errica et al., 2020). 394

395 Attacks. To evaluate the effectiveness of our proposed R-Pool defense method, we consider four 396 adversarial attack strategies: (i) Random Attack, which performs random searches by randomly adding or deleting edges in the input graph; (ii) Genetic Attack (Dai et al., 2018), which modifies 397 graph structures using evolutionary computing principles by evolving a population of candidate 398 solutions through selection, crossover, and mutation to generate adversarial examples; (iii) Gradient-399 *Based Attack* (PGD) (Dai et al., 2018), which greedily adds or deletes edges based on the magnitude 400 of gradients computed with respect to the input graph. For all attacks, we set a perturbation budget 401 of $\epsilon = 0.3$, allowing the modification of up to 30% of the edges in the input graph. 402

Architecture. For all our experiments, a 2-layers GCN classifier with identical hyperparameters and
 activation functions was employed. The models were trained using the cross-entropy loss function,
 and consistent values for the number of epochs and learning rate were maintained across all analysis
 with the Adam optimizer (Kingma & Ba, 2015). Further implementation details can be found in Appendix E and the code implementation to replicate our experiments is provided in the supplementary
 material and would be available on GitHub upon publication.

Baselines. As discussed in Section 2, few methods exist for defending against adversarial attacks 409 on graph classification tasks, particularly at inference time (post-hoc methods). To evaluate our 410 proposed R-Pool defense, we compare it against two main baselines: (i) *Pre-processing techniques*, 411 where the adjacency matrix is pre-processed to remove nodes identified as malicious; specifically, we 412 employ the Jaccard similarity approach (Wu et al., 2019a). (ii) Randomized Smoothing (Bojchevski 413 et al., 2020), which involves adding random noise to the input graph and making predictions based 414 on a majority vote over multiple noisy samples. We evaluate the methods based on two key metrics. 415 First, the *clean accuracy*, which measures the method's performance on unperturbed graphs; this is 416 particularly important for post-hoc methods since, in practice, we may not know whether a graph 417 has been attacked. Second, we consider the *attacked accuracy*.

- 418 419 420
- 6.2 EXPERIMENTAL RESULTS

Table 6.2 reports the clean and attacked graph classification accuracy of our proposed R-Pool and
other considered benchmarks. A primary observation is that R-Pool maintains clean accuracy comparable to the baseline GCN, in contrast to some alternatives such as pre-processing approaches
which often compromise clean accuracy. This finding validates our hypothesis discussed in Section
5, where we argued for the filtering after the message-passing scheme rather than before. Based on
this approach, our proposed R-Pool preserves the flow of relevant information within node representations, thereby maintaining the relevant elements for the downstream graph classification task.

When subject to adversarial attacks, our proposed framework demonstrates robust performance,
 often surpassing or matching the considered benchmarks. In the case of the labeled PROTEINS
 dataset, R-Pool outperforms other baselines in two out of three attack scenarios, showcasing its
 robustness against various adversarial perturbations. For the NCI1 dataset, R-Pool exhibits per formance comparable to randomized smoothing. Interestingly, while the pre-processing technique

Table 1: Clean and Attacked classification accuracy (\pm standard deviation) of the considered baselines on different benchmark graph classification dataset when subject to adversarial attacks. The best accuracy in each setting and each dataset is typeset in **bold**.

Attack	Dataset	PROTEINS	NCI1	IMDB-BINARY
Clean	GCN	73.2 ± 0.4	64.6 ± 0.3	$\textbf{58.9} \pm \textbf{0.2}$
	+ Random Smoothing	74.7 ± 0.3	$\textbf{64.9} \pm \textbf{0.5}$	53.3 ± 0.5
	+ Pre-processing	71.9 ± 0.7	60.1 ± 0.4	-
	+ Ours	$\textbf{75.1} \pm \textbf{0.4}$	63.7 ± 0.9	57.4 ± 0.4
PGD	GCN	57.1 ± 0.8	37.2 ± 0.3	53.0 ± 0.4
	+ Random Smoothing	62.8 ± 1.3	40.3 ± 0.6	51.8 ± 0.7
	+ Pre-processing	59.3 ± 0.7	$\textbf{50.3} \pm \textbf{0.2}$	-
	+ Ours	$\textbf{64.6} \pm \textbf{1.1}$	43.6 ± 0.5	$\textbf{55.7} \pm \textbf{0.9}$
Random	GCN	68.8 ± 0.9	22.3 ± 0.6	54.2 ± 0.3
	+ Random Smoothing	70.3 ± 0.8	25.3 ± 0.7	54.8 ± 0.6
	+ Pre-processing	64.3 ± 0.9	$\textbf{37.3} \pm \textbf{0.4}$	-
	+ Ours	$\textbf{70.6} \pm \textbf{1.2}$	23.1 ± 0.8	$\textbf{56.4} \pm \textbf{0.8}$
Genetic	GCN	63.4 ± 0.8	18.4 ± 0.5	52.4 ± 0.6
	+ Random Smoothing	$\textbf{67.9} \pm \textbf{0.7}$	22.1 ± 0.4	52.9 ± 0.7
	+ Pre-processing	61.7 ± 0.5	$\textbf{35.8} \pm \textbf{0.5}$	-
	+ Ours	67.8 ± 0.9	22.3 ± 0.6	$\textbf{53.2} \pm \textbf{0.9}$

453 454 455

456

457

458

459

460

461

462

463

yields the highest attacked accuracy for NCI1, it does so at the cost of significantly reduced clean accuracy, which is a trade-off that may be undesirable in many real-world applications where maintaining performance on unperturbed data is crucial. This observation underscores the importance of considering both clean and attacked accuracy when evaluating defense mechanisms. Additionally, for the unlabeled IMDB-BINARY dataset, where the Jaccard-based method can't be used due to the absence of node features and randomized smoothing underperforms, R-Pool shows particular promise. This demonstrates the effectiveness of our method with different graph datasets, including those lacking node attributes. Overall, R-Pool demonstrates a good ability to balance clean accuracy and robustness against adversarial attacks, a characteristic not consistently observed in other techniques such as pre-processing.

464 465 466

467 468

7 CONCLUSION AND LIMITATIONS

In this work, we have demonstrated that the choice of pooling operation in graph classification tasks 469 can significantly influence a model's adversarial robustness. Our comprehensive study of three main 470 flat pooling methods revealed that their effectiveness varies across different settings, highlighting the 471 importance of context-specific selection. Building on these theoretical insights, we concluded that 472 filtering nodes after the message-passing scheme but before the pooling operation could enhance 473 robustness. This led to the development of R-Pool, a novel approach that employs Gaussian Mixture 474 Models (GMMs) and an out-of-distribution score to rank nodes and filter out those deemed vulner-475 able. The proposed method can be adapted to different architectures and doesn't require re-training 476 the model and can directly be employed in the inference time. Through extensive experimental val-477 idation across various graph classification datasets, we have demonstrated the efficacy of R-Pool in comparison to existing baselines. 478

The limitations of this work can be categorized into two main areas. Firstly, our study focused exclusively on flat pooling methods, and therefore extending the analysis to hierarchical pooling represents an important step to demonstrate the universality of our approach. Secondly, the current proposed filtering method, R-Pool, is inherently "post-hoc" in nature. This approach offers a new perspective on the inference time defense, which is important with the current surge of Foundation models. A natural progression of this work would be the development of end-to-end trainable methods that integrate robustness considerations directly into the model's training phase, potentially leading to more adversarial robust graph classifiers.

486 REFERENCES 487

517

519

523

524

525

526

527

528

488	Yassine ABBAHADDOU, Sofiane ENNADIR, Johannes F. Lutzeyer, Michalis Vazirgiannis, and
489	Henrik Boström. Bounding the expected robustness of graph neural networks subject to node
490	feature attacks. In The Twelfth International Conference on Learning Representations, 2024.
491	URL https://openreview.net/forum?id=DfPtC8uSot.
492	Amr Alkhatib, Sofiane Ennadir, Henrik Boström, and Michalis Vazirgiannis. Interpretable graph
493	neural networks for tabular data. arXiv preprint arXiv:2308.08945, 2023.
494	
495	Aleksandar Bojchevski and Stephan Günnemann. Certifiable robustness to graph perturbations,
496	2019. URL https://arxiv.org/abs/1910.14356.
497	Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness certifi-
498	cates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more, 2020.
499	URL https://arxiv.org/abs/2008.12952.
500	
501	Chen Cai, Dingkang Wang, and Yusu Wang. Graph coarsening with neural networks. arXiv preprint
502	<i>arxiv:2102.01550</i> , 2021.
503	Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial Attack
504	on Graph Structured Data. In Proceedings of the 35th International Conference on Machine
505	<i>Learning</i> , pp. 1115–1124, 2018.
506	

- David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, 507 Alan Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning 508 molecular fingerprints. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett 509 (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 510 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/ 511 file/f9be311e65d81a9ad8150a60844bb94c-Paper.pdf. 512
- Sofiane Ennadir, Yassine Abbahaddou, Johannes F Lutzeyer, Michalis Vazirgiannis, and Henrik 513 Boström. A simple and yet fairly effective defense for graph neural networks. In Proceedings of 514 the AAAI Conference on Artificial Intelligence, volume 38, pp. 21063–21071, 2024. 515
- 516 Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. All you need is low (rank) defending against adversarial attacks on graphs. In *Proceedings of the 13th* 518 international conference on web search and data mining, pp. 169–177, 2020.
- Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A Fair Comparison of Graph 520 Neural Networks for Graph Classification. In 8th International Conference on Learning Repre-521 sentations, 2020. 522
 - Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In ICLR Workshop on Representation Learning on Graphs and Manifolds, 2019.
 - Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In International Conference on Machine Learning, pp. 1263-1272. PMLR, 2017.
- 529 Lukas Gosch, Simon Geisler, Daniel Sturm, Bertrand Charpentier, Daniel Zügner, and Stephan 530 Günnemann. Adversarial training for graph neural networks: Pitfalls, solutions, and new direc-531 tions. Advances in Neural Information Processing Systems, 36, 2024.
- Stephan Günnemann. Graph neural networks: Adversarial robustness. In Graph Neural Networks: 533 Foundations, Frontiers, and Applications, pp. 149–176. Springer, 2022. 534
- 535 Wei Jin, Yaxing Li, Han Xu, Yiqi Wang, Shuiwang Ji, Charu Aggarwal, and Jiliang Tang. Adver-536 sarial attacks and defenses on graphs. SIGKDD Explor. Newsl., pp. 19-34, 2021.
- Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. Molecular graph 538 convolutions: moving beyond fingerprints. Journal of Computer-Aided Molecular Design, 30(8): 595-608, 2016.

- 540
 541
 542
 542
 543
 544
 544
 544
 545
 545
 546
 546
 547
 547
 548
 548
 549
 549
 549
 549
 540
 541
 541
 542
 542
 544
 544
 544
 545
 546
 546
 547
 548
 548
 549
 549
 549
 549
 549
 540
 541
 541
 542
 542
 544
 544
 544
 545
 546
 547
 548
 548
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
 549
- Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Net works. In *International Conference on Learning Representations (ICLR)*, 2017.
- Chuang Liu, Yibing Zhan, Jia Wu, Chang Li, Bo Du, Wenbin Hu, Tongliang Liu, and Dacheng Tao.
 Graph pooling for graph neural networks: Progress, challenges, and opportunities. *arXiv preprint arXiv:2204.07321*, 2022.
- Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *arXiv preprint arXiv:2007.08663*, 2020.
- Peyman Morteza and Yixuan Li. Provable guarantees for understanding out-of-distribution detec tion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7831–7840, 2022.
 - Aymen Qabel, Sofiane Ennadir, Giannis Nikolentzos, Johannes F Lutzeyer, Michail Chatzianastasis, Henrik Boström, and Michalis Vazirgiannis. Structure-aware antibiotic resistance classification using graph neural networks. In *NeurIPS 2022 AI for Science: Progress and Promises*, 2022.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
 Bengio. Graph Attention Networks. In *ICLR*, 2018.
- 561 Xingchen Wan, Henry Kenlay, Robin Ru, Arno Blaas, Michael A Osborne, and Xiaowen Dong. Adversarial attacks on graph classifiers via bayesian optimisation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 6983–6996. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/38811c5285e34e2e3319ab7d9f2cfa5b-Paper.pdf.
 - Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. Adversarial examples for graph data: Deep insights into attack and defense. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 4816–4823, 2019a.
- Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. Session-based Recommendation with Graph Neural Networks. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, pp. 346–353, 2019b.
 - Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin. Topology attack and defense for graph neural networks: An optimization perspective. *arXiv* preprint arXiv:1906.04214, 2019a.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural
 Networks? In 7th International Conference on Learning Representations, 2019b.
- 580 Xiang Zhang and Marinka Zitnik. Gnnguard: Defending graph neural networks against adversarial attacks. *Advances in Neural Information Processing Systems*, 33:9263–9275, 2020.
- Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta
 learning. In 7th International Conference on Learning Representations, 2019.
- Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, jul 2019. doi: 10.1145/3292500.3330905. URL https://doi.org/10.1145%2F3292500.3330905.
- Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural networks
 for graph data. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2018.

589

555

556

558

567

568

569

570

574

575

576

A APPENDIX

B PROOF OF THEOREM 4.2

Theorem 4.2 Let $f : (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ denote a graph-based function composed of L GCN layers, where the weight matrix of the *i*-th layer is denoted by $W^{(i)}$. Further, let $d^{0,1}$ be a graph distance. For adversarial attacks only targeting node features of the input graph, with a budget ϵ , we have:

• If f is Max-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

 $\gamma = \prod_{l=1}^{L} \|W^{(l)}\| \max_{u \in \mathcal{V}} \hat{w_u} \epsilon$

• If f is Sum-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

$$\gamma = \prod_{l=1}^{L} \|W^{(l)}\| \sum_{u \in \mathcal{V}} \hat{w_u} \epsilon$$

• If f is Average-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:

$$\gamma = \frac{\epsilon}{|V|} \prod_{l=1}^{L} \|W^{(l)}\| \sum_{u \in \mathcal{V}} \hat{w_u}$$

with $\hat{w_u}$ denoting the sum of normalized walks of length (L-1) starting from node u.

Proof. In this proof, we consider that f is a graph-function that is based on L layers of GCN. We recall taht the GCN message-passing propagation is formulated for a node u as

$$h_u^{(\ell)} = \sigma^{(\ell)} \left(\sum_{v \in \mathcal{N}(u) \bigcup \{u\}} \frac{W^{(\ell)} h_v^{(\ell-1)}}{\sqrt{(1+d_u)(1+d_v)}} \right)$$
(2)

620 where $W^{(\ell)} \in \mathbb{R}^{d_{\ell-1} \times d_{\ell}}$ is the learnable weight matrix with d_{ℓ} being the embedding dimension of 621 layer ℓ and $\sigma^{(\ell)}$ is the activation function of ℓ -th layer. We recall that $h^{(0)} = X \in \mathbb{R}^{n \times d}$ is set to the 622 initial node features.

Similar to the work (ABBAHADDOU et al., 2024), we denote X as the original node features and denote by X' the perturbed adversarial features. We consider a node $u \in V$, we denote by h_u its representation in the clean graph and h'_u its representation in the attacked graph. We consider that the activation functions $(\sigma^{(\ell)})_{1 \le \ell \le L}$ are *nonexpensive* (1-Lipschitz continuous). From the work, we have the following result:

$$\|h_{u}^{(L)} - {h'}_{u'}^{(L)}\| \le \prod_{l=1}^{L} \|W^{(l)}\|_{2} \|\sum_{v \in \mathcal{N}(u)} \sum_{\substack{v \in \mathcal{N}(v) \cup \{u\} \\ j \in \mathcal{N}(v) \cup \{v\}}} \dots$$

$$\sum_{z \in \mathcal{N}(y) \cup \{y\}} \frac{X_u - X'_u}{\sqrt{(1+d_u)}(1+d_w)(1+d_j)\dots(1+d_y)\sqrt{(1+d_z)}} \|$$

$$\leq \prod_{l=1}^L \|W^{(l)}\|\hat{w}_u \epsilon$$

with \hat{w}_u being the sum of normalized walks of length (L-1) starting from node u.

The previous results gives us an idea about the behavior of each node's representation when attacked. In the case of graph classification, an additional pooling operation is added, Specifically:

 $\mathbf{h}_{\text{graph}}^{(L)} = \text{Pool}\left(\{h_u^{(L)}\}_{u \in V}\right)$

Hence this proof's goal is to analyze the following quantity:

$$\|\mathbf{h}_{ ext{graph}}^{(L)}-\mathbf{h'}_{ ext{graph}}^{(L)}|$$

648 Let's consider the **Sum-pooling** operation, that can be written as:

$$\mathbf{h}_{ ext{graph}}^{(L)} = \sum_{u \in V} h_u^{(L)}$$

We have the following:

$$\begin{split} \|\mathbf{h}_{\text{graph}}^{(L)} - \mathbf{h}_{\text{graph}}^{\prime(L)}\| &= \left\|\sum_{u \in V} \left(h_u^{(L)} - h_u^{\prime(L)}\right)\right\| \\ &\leq \sum_{u \in V} \|h_u^{(L)} - h_u^{\prime(L)}\| \quad \text{(Triangle Inequality)} \\ &\leq \sum_{u \in V} \left(\prod_{l=1}^L \|W^{(l)}\|\right) \|\hat{w}_u\|\epsilon \\ &= \left(\prod_{l=1}^L \|W^{(l)}\|\right) \epsilon \sum_{u \in V} \hat{w}_u \end{split}$$

For the case of the **Average-pooling** operation, that can be written as:

$$\mathbf{h}_{\text{graph}}^{(L)} = \frac{1}{|V|} \sum_{u \in V} h_u^{(L)}$$

We have the following analysis:

$$\begin{aligned} \|\mathbf{h}_{\text{graph}}^{(L)} - \mathbf{h}_{\text{graph}}^{\prime(L)}\| &= \left\| \frac{1}{|V|} \sum_{u \in V} \left(h_u^{(L)} - {h'}_u^{(L)} \right) \right\| \\ &\leq \frac{1}{|V|} \sum_{u \in V} \|h_u^{(L)} - {h'}_u^{(L)}\| \quad \text{(Triangle Inequality)} \\ &\leq \frac{1}{|V|} \left(\prod_{l=1}^L \|W^{(l)}\| \right) \epsilon \sum_{u \in V} \hat{w}_u \end{aligned}$$

In the case of **Max-pooling**:

$$\begin{split} \left| [\mathbf{h}_{\text{graph}}^{(L)}]_{k} - [\mathbf{h}'_{\text{graph}}^{(L)}]_{k} \right| &\leq \max_{u \in V} \left| [h_{u}^{(L)}]_{k} - [h'_{u}^{(L)}]_{k} \right| \\ &\leq \max_{u \in V} \|h_{u}^{(L)} - h'_{u}^{(L)}\| \\ &\leq \left(\prod_{l=1}^{L} \|W^{(l)}\| \right) \max_{u \in V} \hat{w}_{u} \epsilon \end{split}$$

By taking into account the expectancy (as shown in Definition 1), we get the desired results.

C PROOF OF THEOREM 4.3

699 Theorem 4.3 Let $f : (\mathcal{G}, \mathcal{X}) \to \mathcal{Y}$ be composed of *L* GIN-layers (with its internal parameter $\zeta = 0$) **700** and let $W^{(i)}$ denote the weight matrix of the *i*-th MLP layer. We consider the input node feature **701** space to be bounded i. e., $||\mathcal{X}||_2 < B$ for some $B \in \mathbb{R}$. For node feature-based attacks, with a budget ϵ , the function f is $(d^{0,1}, \epsilon)$ -robust with

702	• If f is Max-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:
703	$\gamma = \prod_{l=1}^{L} \ W^{(l)}\ \left[B \times L \times \max_{u \in V} \deg(u) + \epsilon\right]$
705	• If f is Sum pooling based classifier then f is $d^{0,1}(c, \alpha)$ robust with:
706	• If f is sum-pooling based classifier, then f is $a = (e, f)$ footst with.
707	$\gamma = \prod_{l=1}^{L} \left\ W^{(l)} \right\ \left[2B \times L \times E + V \epsilon \right]$
708	• If f is Average-pooling based classifier, then f is $d^{0,1}$ - (ϵ, γ) robust with:
710	$\gamma = \Pi^L W^{(l)} \left[\frac{2B imes L imes E }{2B imes L imes E } + \epsilon ight]$
711	$V = \prod_{l=1}^{N} V_l + V_l \left[\frac{ V_l }{ V_l } + c\right]$
712 713	with $\mid E \mid$ being the number of edges and $\mid V \mid$ the number of nodes.
714 715 716	<i>Proof.</i> In this proof, we consider that f is based on L GIN-layers (with a parameter $\zeta = 0$, usually denoted as ϵ). The GIN message-passing propagation process can be written for a node u as:
717 718	$h_u^{(\ell+1)} = T^{(\ell+1)}((1+\zeta)h_u^{(\ell)} + \sum_{v \in \mathcal{N}(u)} h_v^{(\ell)})$
719 720	with T denoting a Neural Networks (a MLP) for example and ζ denotes the parameter of the GIN. We recall that $h^{(0)} = X \in \mathbb{R}^{n \times d}$ is set to the initial node features.
721 722 723 724 725	Similar to the previous proof, we base our proof on previous work (ABBAHADDOU et al., 2024), we denote X as the original node features and denote by X' the perturbed adversarial features. We consider a node $u \in V$, we denote by h_u its representation in the clean graph and h'_u its representation in the attacked graph. We consider that the activation functions $(\sigma^{(\ell)})_{1 \le \ell \le L}$ are <i>nonexpensive</i> (1-Lipschitz continuous).
726 727 728 729	We use the same assumptions as the one considered in (ABBAHADDOU et al., 2024). Specifically, we consider that the input feature space \mathcal{H}_0 is bounded, thus each hidden space \mathcal{H}_i of the iterative process of message passing is bounded and let $B = \max_{\ell \leq L} B_{\ell}$ be its global maximum bound. We
730 731 732	additionally consider that GIN-parameter $\zeta \approx 0$ (which is very frequent in the literature). We have therefore the following result:
733 734 735	$\ h_{u}^{(\ell+1)} - {h'}_{u'}^{(\ell+1)}\ \leq \prod_{l=1}^{L} \ W^{(l)}\ [B \times L \times deg(u) + \epsilon]$
736 737 738	From this perspective, let's consider the case of graph classification, where we start by the sumpooling operation:
739 740 741 742	$\left\ h_{G}^{(L)} - {h'}_{G}^{(L)} \right\ = \left\ \sum_{u \in V} \left(h_{u}^{(L)} - {h'}_{u}^{(L)} \right) \right\ $
743 744	$\leq \sum_{u \in V} \left\ h_u^{(L)} - {h'}_u^{(L)} \right\ $ (by the triangle inequality)
745 746 747	$\leq \prod_{l=1}^{L} \left\ W^{(l)} \right\ \sum_{u \in V} \left[B \times L \times \deg(u) + \epsilon \right]$
748 749 750	$= \prod_{l=1}^{L} \left\ W^{(l)} \right\ \left[B \times L \times \sum_{u \in V} \deg(u) + V \epsilon \right].$
751 752 753	In the case of undirected graph, since $\sum_{u \in V} \deg(u) = 2 E $, we have:
754 755	$\left\ h_G^{(L)} - {h'}_G^{(L)} \right\ \leq \prod_{l=1}^L \left\ W^{(l)} \right\ \left[2B \times L \times E + V \epsilon \right]$

⁷⁵⁶ Similar for the case of **Average-pooling**, we have:

$$\begin{split} \left\| h_{G}^{(L)} - {h'}_{G}^{(L)} \right\| &= \frac{1}{|V|} \left\| \sum_{u \in V} \left(h_{u}^{(L)} - {h'}_{u}^{(L)} \right) \right\| \\ &\leq \frac{1}{|V|} \sum_{u \in V} \left\| h_{u}^{(L)} - {h'}_{u}^{(L)} \right\| \\ &\leq \prod_{l=1}^{L} \left\| W^{(l)} \right\| \left[\frac{2B \times L \times |E|}{|V|} + \epsilon \right]. \end{split}$$

In the case of **Max-pooling** operation, we have the following:

$$\left\| h_G^{(L)} - h'_G^{(L)} \right\| \le \prod_{l=1}^L \left\| W^{(l)} \right\| \left[B \times L \times \max_{u \in V} \deg(u) + \epsilon \right].$$

By taking into account the expectancy (as shown in Definition 1), we get the desired results.

D TIME COMPLEXITY ANALYSIS

As explained in Section 5, the main computational complexity of our method is concentrated in the EM algorithm used for estimating the GMM's parameters. The EM algorithm is an iterative process and hence the complexity mainly depends on the number of iterations that have been chosen. For our experimentation, we have seen that 100 iterations was a good number to reach a satisfactory accuracy. In table 2, we provide a time complexity comparison of our R-Pool to other considered baselines on the used graph classification datasets.

Table 2: Mean training time analysis (in s) of a our R-Pool in comparison to the other considered benchmarks on the graph classification datasets.

DATASET	GCN	RANDOMZIED SMOOTHING	PRE-PROCESSING	R-POOL
PROTEINS	0.001	0.014	0.037	0.013
NCI1	0.008	0.019	0.01	0.015
IMDB-BINARY	0.0007	0.013	-	0.011

E EXPERIMENTAL DETAILS

E.1 DATASETS

For our experimentation, we mainly used the classical graph datasets derived from bioinformatics
and chemoinformatics (PROTEINS, NCI1) and social networks (IMDB-BINARY) (Morris et al.,
2020). We used the public folds and the experimental setting that was provided by the work Errica
et al. (2020). Details about the dataset are provided in Table 3.

- 802 E.2
 - E.2 IMPLEMENTATION DETAILS

Our implementation is available in the supplementary materials (and will be publicly available afterwards). It is built using the open-source library *PyTorch Geometric* (PyG) under the MIT license (Fey & Lenssen, 2019) and DGL in the case of the genetic and random attacks. We leveraged the publicly available implementation of both the attacks and the benchmarks. The experiments have been run on both a NVIDIA A100 GPU.

For all the attacks, we set the number of attack epochs to 100 and in a "re-wiring" mode, meaning that the attack can either add/delete an edge.

811

812	-	DATAGET	#CD + DUG	#Norra	#Ep cpc	#CL + 0070	
813	-		#GRAPHS	#NODES	#EDGES	#CLASSES	
814		NCI1	4110	19.77 29.87	32 30	2	
815		PROTEINS	1113	39.06	72.82	2	
816	-						
817							
818	T		1 . 2 1		1 1 .		
819	For all our exper	iments, we employ	ed a 2-layer	convolution	hal archited	cture (consistin	g of two itera-
820	using the Adam	Optimizer Kingma	$r R_2 (2014)$	Will a Mul 5) We trair	the mode	1 for 100 epocl	(as this was)
821	sufficient to reac	h the state-of-the ar	t accuracy fo	r these mod	els) and we	e use a learning	rate of $1e - 02$
822	sumerent to reach	in the state of the ar	t decuracy fo	r mese mou	cis) and we	use a learning	
823							
824							
825							
826							
827							
828							
829							
830							
831							
832							
833							
834							
835							
836							
837							
838							
839							
840							
841							
842							
843							
044							
040 946							
040 9/17							
8/8							
849							
850							
851							
852							
853							
854							
855							
856							
857							
858							
859							
860							
861							
862							
863							

Table 3: Statistics of the graph classification datasets used in our experiments.