
If You Want to Be Robust, Be Wary of Initialization

Sofiane Ennadir*
KTH
Stockholm, Sweden

Johannes F. Lutzeyer
LIX, Ecole Polytechnique
Paris, France

Michalis Vazirgiannis
KTH & Ecole Polytechnique
Stockholm, Sweden

El Houcine Bergou
UM6P
Benguerir, Morocco

Abstract

1 Graph Neural Networks (GNNs) have demonstrated remarkable performance across
2 a spectrum of graph-related tasks, however concerns persist regarding their vul-
3 nerability to adversarial perturbations. While prevailing defense strategies focus
4 primarily on pre-processing techniques and adaptive message-passing schemes, this
5 study delves into an under-explored dimension: the impact of weight initialization
6 and associated hyper-parameters, such as training epochs, on a model’s robustness.
7 We introduce a theoretical framework bridging the connection between initializa-
8 tion strategies and a network’s resilience to adversarial perturbations. Our analysis
9 reveals a direct relationship between initial weights, number of training epochs and
10 the model’s vulnerability, offering new insights into adversarial robustness beyond
11 conventional defense mechanisms. While our primary focus is on GNNs, we extend
12 our theoretical framework, providing a general upper-bound applicable to Deep
13 Neural Networks. Extensive experiments, spanning diverse models and real-world
14 datasets subjected to various adversarial attacks, validate our findings. We illustrate
15 that selecting appropriate initialization not only ensures performance on clean
16 datasets but also enhances model robustness against adversarial perturbations, with
17 observed gaps of up to 50% compared to alternative initialization approaches.

18 1 Introduction

19 Neural networks have demonstrated remarkable prowess across various domains, ranging from
20 computer vision [7] to natural language processing [28], proving their ability to model and extract
21 complex insights from real-world datasets. Recently, Graph Neural Networks (GNNs) [20, 35, 29]
22 have emerged as a powerful extension of neural networks specifically tailored to tackle graph-
23 structured data. These models have led to rapid progress in solving tasks such as node and graph
24 classification where their application have spanned from drug design [19], protein resistance anal-
25 ysis [23] to session-based recommendations [32]. Concurrently with their success, deep learning
26 architectures have been shown to be unstable when subject to adversarial perturbations [14], resulting
27 in unreliable predictions, consequently questioning these models’ applicability in critical domains.
28 While most adversarial robustness studies focus on the domain of computer vision, recent work [15]
29 studying the robustness of GNNs has emerged. Given their rich nature, graphs allow different attack
30 schemes, where the attacker can either choose to edit the graph structure (by adding/deleting edges)
31 or edit the node/edge features. In parallel, recent studies have been devoted to studying approaches to
32 defend against these attacks and enhance GNN’s robustness, such as input pre-processing techniques
33 [31], low-rank approximation [10], edge-pruning [37] or adapting the message-passing schemes [1].

*Corresponding Author: ennadir@kth.se

34 The majority of available defense studies focus on understanding the inner dynamics of GNNs to
35 pinpoint and mitigate adversarial vulnerabilities. While analyzing the message-passing mechanism
36 and implementing input pre-processing techniques remains a viable direction, comprehensive under-
37 standing necessitates exploration beyond traditional avenues. In this sense, investigating factors such
38 as weight initialization strategies and the impact of other hyperparameters, notably those associated
39 with optimization mechanisms, can offer new insights and perspectives on achieving GNN’s global
40 robustness. Hyperparameter choices and tuning play a critical role in striking a balance between learn-
41 ing the underlying signals in the data and preventing overfitting to ensure the model’s generalization.
42 Hence, existing studies on initialization mainly evolves around understanding its effect on the model’s
43 convergence, stability and performance [33, 22]. In contrast, the current work primarily focuses on
44 examining the effect of initialization on a model’s underlying adversarial robustness, representing
45 to the best of our knowledge the first exploration of its kind. Our main objective is to provide a
46 theoretical understanding of the link between weight initialization and other dynamics such as the
47 number of training steps and the resulting model’s robustness. With this perspective in mind, we start
48 by formalizing robustness in the context of GNNs when subjected to structural and node feature-based
49 adversarial attacks. Subsequently, we derive an upper bound that connects the model’s robustness to
50 the weight initialization strategies. Specifically, we illustrate that this bound depends on the initial
51 weight norms and the number of training epochs. Finally, we validate our theoretical findings by
52 demonstrating the effects of employing various initialization strategies on the model’s robustness
53 using benchmark adversarial attacks on real-world datasets. Note that while our analysis primarily
54 focuses on the widely used Graph Convolutional Networks (GCNs) [20] and Graph Isomorphism
55 Networks (GIN) [35], we highlight the versatility of our approach by providing a general upper bound
56 applicable to any Deep Neural Networks in Section 5. This underlines the potential for extending our
57 analysis to a wide range of architectures, showcasing its broad applicability in understanding and
58 enhancing adversarial robustness in neural networks. We summarize our contributions as follows:

- 59 • We provide a theoretical analysis that links weight initialization strategies with adversarial
60 robustness in GNNs. We specifically derive an upper bound connecting a model’s robust-
61 ness to weight initialization and the number of training epochs, demonstrating that the
62 initialization strategy can significantly influence the network’s adversarial robustness.
- 63 • We validate our theoretical findings by conducting extensive experiments across various
64 models using different benchmark adversarial attacks on real-world datasets. These exper-
65 iments demonstrate that certain weight initialization strategies can enhance the model’s
66 defense against adversarial attacks, without degrading its performance on clean datasets.
- 67 • While our primary focus is on GNNs, we extend our analysis to Deep Neural Networks,
68 illustrating the broader applicability of our theoretical analysis and its corresponding insights.

69 2 Related Work

70 **Graph Adversarial Attacks.** Multiple studies focus on designing adversarial attacks capable of
71 fooling a graph-based classifier [15, 34, 9]. The majority of these methods [41, 36] approach the
72 adversarial aim as an optimization problem and employs different methods to solving it such as meta-
73 learning [40]. Furthermore, Nettack [39] constrained the problem by preserving degree distribution
74 and imposing constraints on feature co-occurrence to generate unnoticeable perturbations. Finally,
75 reinforcement learning was proposed recently as a mean to generate graph adversarial attacks [6].

76 **Graph Adversarial Defenses.** Recent efforts have emerged to defend against the aforementioned
77 adversarial attacks. In particular, methods such as low-rank matrix approximation coupled with graph
78 anomaly detection [21] have been used. For example, GNN-Jaccard [31] proposed to pre-process
79 the graph’s adjacency matrix to detect potential manipulation of edges. Other methods such as edge
80 pruning [37] and transfer learning [27] have been leveraged to limit the effect of poisoning attacks.
81 Additionally, adaptations of the message-passing scheme, such as employing orthogonal weights
82 [1] or introducing noise during training [8], have been shown to perform well in term of defense.
83 Furthermore, there is a growing interest in exploring robustness certificates [41, 3] as a means of
84 ensuring model robustness. For instance, [4] used randomized smoothing to provide a highly scalable
85 model-agnostic certificate for graphs. Additionally, other robustness certificates for GCN-based graph
86 classification under topological perturbations have been proposed [18].

87 **Weight Initialization.** The impact of weight initialization has been extensively studied both theoreti-
 88 cally and empirically where the main line of study consists of understanding the interplay between
 89 initialization techniques and the implicit regularization they induce, thereby elucidating their influence
 90 on a model’s generalization capabilities [33, 22]. For instance, it has been showcased that sampling
 91 initial weights from the orthogonal group can speed up convergence [17]. Similarly, alternative
 92 initialization approaches such as Glorot Initialization [12] and Kaiming Initialization [16] have been
 93 proposed in efforts to improve the model’s performance.

94 Our work stands apart from existing research on adversarial robustness as it represents, to the best of
 95 our knowledge, the first attempt to theoretically investigate the impact of initialization on a model’s
 96 underlying robustness. Moreover, our approach diverges fundamentally from existing literature on
 97 weight initialization as our focus lies in theoretically understanding the effect of initialization on a
 98 model’s robustness rather than its implications for generalization or convergence.

99 3 Graph Adversarial Robustness

100 In this section, we start by introducing the notation and some fundamental concepts related to GNNs.
 101 We afterwards establish the problem setup together with the set of considered assumptions. We finally
 102 lay out a GNN’s robustness formalization on which we will build our theoretical analysis.

103 3.1 Preliminaries

104 Let $G = (V, E)$ be a graph where V ($|V| = n$) is its set of vertices and E its set of edges. We
 105 denote $A \in \mathcal{A} \triangleq \{0, 1\}^{n \times n}$ its adjacency matrix. The graph nodes are annotated with feature vectors
 106 $X \in \mathcal{X} \subseteq \mathbb{R}^{n \times d}$ (the i -th row of X corresponds to the feature of node i). We denote by $\mathcal{N}(i)$ the
 107 neighbors of node $i \in V$ and $\|\cdot\|_2$ the Euclidean (resp., spectral) norm for vectors (resp., matrices).

108 In this work, we consider the task of node classification. In this task, every node is assigned exactly
 109 one class from $\mathcal{C} = \{1, 2, \dots, C\} \subset \mathcal{Y}$ and we consider $d_{\mathcal{Y}}$ as a distance within the output space \mathcal{Y} .
 110 The learning objective is to find a function f_W , parameterized by W , that assigns each node $i \in V$ a
 111 class $c \in \mathcal{C}$ while minimizing some classification loss (e. g., cross-entropy loss), denoted as \mathcal{L} .

112 **GNNs.** A GNN model consists of a series of neighborhood aggregation layers that use the graph
 113 structure and the node features from the previous layers to generate new nodes representations.
 114 Specifically, GNNs update node feature vectors by aggregating local neighborhood information. In
 115 the particular case of GCNs, this process is described by the following iterative propagation:

$$h^{(\ell)} = \phi^{(\ell)}(\widehat{A}h^{(\ell-1)}W^{(\ell)}), \quad (1)$$

116 with $W^{(\ell)} \in \mathbb{R}^{p \times q}$ being the weight matrix in the ℓ -th layer, q is the embedding dimension and $\phi^{(\ell)}$
 117 is a non-linear activation function. Moreover, $\widehat{A} \in \mathbb{R}^{n \times n}$ denotes the normalized adjacency matrix
 118 $\widehat{A} = D^{-1/2}AD^{-1/2}$ where $D = \text{diag}(|\mathcal{N}(1)|, |\mathcal{N}(2)|, \dots, |\mathcal{N}(n)|)$ denotes the degree matrix.

119 **Problem Setup.** For our theoretical analysis, we assume that the model is based on 1-Lipschitz
 120 activation functions (which is a characteristic of commonly used activation functions such as TanH).
 121 Additionally, we consider the training loss function \mathcal{L} to be L -smooth and that it is minimized using
 122 gradient descent. We denote by W_* the local optimum towards which gradient descent iterates
 123 converge. Specifically, for a learning rate $\eta < \frac{1}{L}$, the update at time step t for a layer i is:

$$W_{t+1}^{(i)} = W_t^{(i)} - \eta \nabla \mathcal{L}(W_t^{(i)}).$$

124 It is worth emphasizing that although we focus on the node classification task, which is prevalent and
 125 well-studied in the literature of adversarial robustness, our analysis is equally applicable to other tasks
 126 such as graph classification. Moreover, while our theoretical analysis predominantly centers around
 127 using the gradient descent as the optimizer, this choice doesn’t limit the generality of our findings.
 128 One can employ a different optimizer and still yield the same insights and results by following a
 129 similar approach as the one outlined in this paper. Consequently, this specific setup should not be
 130 perceived as a limitation but rather as an analytical choice.

131 **3.2 Adversarial Robustness for Graph Neural Networks**

132 Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ be a GNN-classifier following the framework outlined in Section 3.1. An
 133 adversarial attacks consists of generating an alternative graph (\tilde{A}, \tilde{X}) that perturbs the original
 134 prediction $f(A, X)$ while not being far (semantically) from the original graph. Typically, this
 135 generated graph must adhere to a number of constraints related to its similarity to the original graph,
 136 defined by a perturbation budget ϵ controlling the number of edited edges or features. The set of
 137 these graphs is written as $B([A, X]; \epsilon) = \{(\tilde{A}, \tilde{X}) : \min_{P \in \Pi} (\|A - P\tilde{A}P^T\|_2 + \|X - P\tilde{X}\|_2) \leq \epsilon\}$,
 138 where Π represents the set of permutations of the adjacency matrix. While the previous formulation
 139 relies on the ℓ_2 norm, other norms may be used depending on the domain of application and the
 140 specific use case. Building on previous work [8], the adversarial risk of a GNN can be defined as the
 141 expected error of adjacent graphs within the considered graph’s neighborhood defined by ϵ written as:

$$\mathcal{R}_\epsilon[f] = \mathbb{E}_{(A, X) \sim \mathcal{D}} \left[\sup_{(\tilde{A}, \tilde{X}) \in B([A, X]; \epsilon)} d_{\mathcal{Y}}(f(\tilde{A}, \tilde{X}), f(A, X)) \right]. \quad (2)$$

142 In the current analysis, we focus on the ℓ_2 norm as our output distance $d_{\mathcal{Y}}$ (which can be substituted
 143 by any norm – giving the existence of norm’s equivalence). We theoretically approach the introduced
 144 adversarial risk by deriving an upper-bound, which reflects the model’s expected error under input
 145 perturbation. Intuitively, a smaller upper bound reflects a smaller adversarial risk which in turn
 146 suggests a robust behavior locally. In this perspective, Definition 1 draws the link between the
 147 considered risk quantity and a model’s robustness.

148 **Definition 1.** (Adversarial Robustness). The graph-based function $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ is said to be
 149 (ϵ, γ) – robust if its adversarial risk is upper-bounded by γ , i. e., $\mathcal{R}_\epsilon[f] \leq \gamma$.

150 The current definition addresses adversarial risk from a worst-case scenario perspective, which is
 151 the most prevalent approach in the literature. This means we aim to identify the neighbor graph
 152 that maximizes the harm (i. e., causes the greatest deviation from the original prediction). By upper-
 153 bounding the risk associated with this "worst-case" graph, we inherently account for all other potential
 154 adversaries within the same neighborhood, as their risk will be less than or equal to that of the worst-
 155 case scenario. We note that the nuances between the “average” and “worst-case” approaches have
 156 been thoroughly examined and justified in previous research [24].

157 **4 On the Effect of Initialization**

158 We start by considering the Graph Convolutional Networks (GCNs) within the broader context of
 159 Message Passing Neural Networks for node classification. This study investigates how initialization
 160 and other hyper parameters impact the final model’s robustness. In this context, we aim to establish a
 161 connection between the introduced adversarial risk (Equation 2) and the initial weight distribution
 162 and its evolution during training. Specifically, we seek to demonstrate that different choices in the
 163 initialization distribution and other relevant parameters lead to varying levels of model robustness,
 164 offering new insights into the potential trade-offs between initialization strategies and robustness. In
 165 this sense, we derive an upper-bound (denoted as γ in Definition 1) on the stability of a GCN-based
 166 classifier when the input graph’s node features are subject to adversarial attacks.

167 **Theorem 2.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers, where
 168 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. For adversarial attacks only targeting
 169 node features of the input graph, with a budget ϵ , we have (in respect to Definition 1):

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

170 with t being the number of training epochs and \hat{w}_u denoting the sum of normalized walks of length
 171 $(T - 1)$ starting from node u .

172 Theorem’s proof is provided in Section A of the Appendix. Theorem 2 provides a formal connection
 173 between the robustness of a GCN-based classifier and its initial weights, offering valuable insights
 174 into their effects. From a first perspective, the derived upper-bound depends on the initial weight’s

175 norm. Specifically, a lower norm corresponds to a smaller upper-bound, indicative of a more robust
 176 model. However, while setting all initial weights to zero theoretically yields the smallest upper-
 177 bound and consequently the optimum robustness, this direction can detrimentally affect the model’s
 178 performance on the learning task. Empirical evidence suggests that initializing weights to zero (or a
 179 constant) often leads to poor learning outcomes, as it constrains weight behavior during propagation,
 180 limiting subsequent back-propagation operations and resulting in convergence to unsatisfactory local
 181 minima (e. g., see page 301 in [13]). From a second perspective, it appears that a higher number
 182 of training epochs leads to the looseness of the upper-bound, resulting in increased adversarial
 183 vulnerability. This latter observation provides proofs and highlights on the existence of the usually
 184 discussed trade-off between clean and attacked accuracy. Achieving a balance between increasing the
 185 number of epochs to achieve satisfactory clean accuracy and limiting them to attain a robust model is
 186 hence essential. While theoretically challenging to identify this equilibrium point, our experimental
 187 results demonstrate its existence. We note that the dependence of γ on t can be sharpened by having
 188 $(1 + \eta L)^t$ instead of 2^t . With small η (which is the case usually in practice), $(1 + \eta L)^t \approx 1 + t\eta L$
 189 resulting in a bound which depends linearly in t . The same remark applies for the remaining bounds
 190 derived in the paper. These insights, in the case of node-feature-based adversarial attacks, also extends
 191 to structural perturbations where Theorem 3 provides the exact bound for this case.

192 **Theorem 3.** *Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers, where
 193 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. Let f be the number of used training
 194 epochs. When f is subject to structural attacks, with a budget ϵ , we have (in respect to Definition 1):*

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \left\| W_0^{(i)} \right\| + 2^{t+1} \left\| W_*^{(i)} \right\| \right) \|X\| \left(1 + T \prod_{i=1}^T \left(2^t \left\| W_0^{(i)} \right\| + 2^{t+1} \left\| W_*^{(i)} \right\| \right) \right)$$

195 The computed upper-bound suggests that the effect of initialization is more important in the case of
 196 structural perturbations. This emphasis is resulting from the distinct dynamics within the message
 197 passing mechanism, where the influence of the adjacency matrix and node features varies during each
 198 propagation step. Precisely, for structural perturbations, the effect of the attack is considered at each
 199 propagation step through the perturbed adjacency matrix (in the aggregation step). Moreover, the
 200 impact is also amplified by the affected residual layers from previous iterations, resulting in a more
 201 significant attack result. This is different in the case of node-feature based adversarial attacks, since
 202 the node features are only taken into account in the first propagation. Overall, the main takeaway
 203 of the provided analysis in Theorem 2 and 3 is that “approximately-free” robustness enhancements
 204 can be derived from choosing the right initial weight’s distribution and the right number of training
 205 epochs. We illustrate this specific point by analyzing the effect of the initial distributions choices on
 206 the model’s robustness. Specifically, we consider the case of the Gaussian distribution, where Lemma
 207 4 studies how the parameters of this distribution – namely, the mean and variance – exert an influence
 208 on the expected (in respect to the initial distribution) value of the adversarial risk.

209 **Lemma 4.** *Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers for
 210 which the initial weight are drawn from the Gaussian distribution $\mathcal{N}(\mu, \Sigma)$. When subject to node
 211 features based adversarial attacks, we have the following:*

$$\mathbb{E}_{W_0 \sim \mathcal{N}(\mu, \Sigma)} [\mathcal{R}_\epsilon[f]] \leq \epsilon \prod_{i=1}^T \left(2^t \sqrt{\mu^2 + \text{tr}(\Sigma)} + 2^{t+1} \left\| W_*^{(i)} \right\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

212 Proof is provided in Section C. Given that a tighter upper bound inherently results in a higher level of
 213 robustness, the results derived in Lemma 4 illustrate the clear effect of initialization in the case of the
 214 Gaussian distribution. The derived bound shows that increasing the distribution parameters, both the
 215 mean and variance values, leads to a decrease in the victim model’s underlying robustness. While one
 216 might intuitively aim to set these parameters as low as possible to achieve optimal robustness, doing
 217 so could potentially compromise the model’s performance on clean datasets. Therefore, as previously
 218 mentioned, striking the right balance between clean accuracy and adversarial robustness is crucial.

219 **Extending the results to the GIN.** The same previously applied analysis for the GCN-based models
 220 can be extended to take into account GIN-based classifiers. We consider the same set of assumptions
 221 and the same problem setup considered during the previously studied GCN case. We additionally
 222 assume that the input node feature space to be bounded, i. e., $\|X\| \leq B$. We note that this bound is a
 223 realistic assumption and that the value B can be easily computed for any real-world dataset.

224 **Theorem 5.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GIN layers, where
 225 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. For adversarial attacks only targeting
 226 node features of the input graph, with a budget ϵ , we have:

$$\gamma = \prod_{l=1}^T \left(2^l \|W_0^{(i)}\| + 2^{l+1} \|W_*^{(i)}\| \right) \left[BT \max_{u \in \mathcal{V}} \text{deg}(u) + \epsilon \right]$$

227 with t being the number of training epochs and $\text{deg}(u)$ is the degree of node u .

228 Proof of the theorem is provided in the appendix (Section D). Theorem 5 establishes an upper bound
 229 on the robustness of a GIN-based classifier against adversarial attacks targeting node features. We
 230 observe analogous insights, to the ones derived for a GCN-based classifier, regarding the influence of
 231 the initialization distribution and number of training epoch on the model’s underlying robustness.

232 5 Generalization to Other Models

233 While our primary research focus lies within the domain of graph representation learning, a sub-field
 234 of the broader landscape of Deep Learning models, the fundamental principles of our theoretical
 235 analysis hold applicability across various model architectures. Notably, and to our knowledge, the
 236 absence of a comparable study in current adversarial literature motivates our endeavor to bridge this
 237 gap. In this section, we aim to fill this void by presenting a comprehensive analytical framework that
 238 provide the connection between weight initialization and the robustness of neural networks.

239 Let $x \in \mathbb{R}^{n_0}$ denote an input vector where n_0 is the input dimension. Let $W^{(l)} \in \mathbb{R}^{n_{l-1}, n_l}$ be the
 240 weight matrix and $b_l \in \mathbb{R}^{n_l}$ the bias of the l^{th} layer with n_l being its dimensionality. We focus on
 241 the general family of neural networks for which the computation during layer l , using an activation
 242 function $\phi^{(l)}$, can be written as :

$$h^{(l)} = \phi^{(l)}(W^{(l)}h^{(l-1)} + b^{(l)}).$$

243 We consider the same set of assumptions (stated in Section 3.1) as the one from previous section. We
 244 consider the ℓ_2 norm as our input and output distances within the metric space \mathbb{R}^{n_0} and we consider
 245 an input attack budget ϵ . The introduced adversarial risk in Equation 2 can be easily extended
 246 and tailored to the family of considered neural networks discussed in this section. Further clarification
 247 on this extension is provided in the Appendix (Section G.1). From this standpoint, by adapting the
 248 Definition 1, analogous effects of the weight initialization, provided in Theorem 6, can be observed.

249 **Theorem 6.** Let $f : \mathcal{X} \subseteq \mathbf{R}^{in} \rightarrow \mathcal{Y} \subseteq \mathbf{R}^{out}$ be a T -layers neural network with $W_0^{(i)}$ denoting the
 250 initial weight matrix of the i -th layer. When subject to adversarial attacks, f is (ϵ, γ) – robust with:

$$\gamma = \epsilon \prod_{i=1}^T \left(2^i \|W_0^{(i)}\| + 2^{i+1} \|W_*^{(i)}\| \right)$$

251 The proof of Theorem 6 can be found in Section E of the Appendix. Similar to previous findings,
 252 the upper bound relies on key elements of the initialization process, specifically the initial weight
 253 norm and the number of training epochs. These results validate and extend the established link
 254 between initialization and a model’s robustness in neural networks, highlighting the importance
 255 of selecting appropriate parameters. From the derived upper bound, which is also applicable to
 256 GCN and GIN cases, we observe that the number of training epochs exerts an effect on the bound.
 257 Specifically, while increasing the number of epochs can improve the model’s performance on a clean
 258 dataset, it simultaneously leads to a deterioration in the model’s adversarial robustness. Ideally,
 259 adversarial defense strategies aim to avoid this trade-off between clean and attacked accuracy, striving
 260 for robust models that do not compromise the initial performance. In this context, considering the
 261 strong-convexity of the loss function \mathcal{L} , in addition to the previously made assumptions, we observe
 262 that the effect of the number of training epochs becomes less pronounced. Lemma 7 specifically
 263 provides the computed bound under these assumptions.

264 **Lemma 7.** Let $f : \mathcal{X} \subseteq \mathbf{R}^{in} \rightarrow \mathcal{Y} \subseteq \mathbf{R}^{out}$ be a T -layers neural network trained with a μ -strongly
 265 convex and L -smooth loss function. Let $W_0^{(i)}$ denote the initial weight matrix of the i -th layer. When

266 *subject to adversarial attacks, with a budget ϵ , we have that f is (ϵ, γ) – robust with:*

$$\gamma = \epsilon \prod_{i=1}^T \left((1 - \mu/L)^t \|W_0^{(i)}\| + 2\|W_*^{(i)}\| \right)$$

267 The proof of the Lemma is provided in Section F of the Appendix. Since $\mu \leq L$, increasing the
268 number of training epochs results in the diminishing influence of the initialization weights. In this
269 scenario, the bound depends solely on the final weights, a phenomenon previously explored in works
270 such as Parseval networks [5] for neural networks and GCORN [1] for GNNs. This observation
271 highlights the necessity of convexity in the loss function when training a neural network, as it plays a
272 crucial role in enhancing the model’s robustness, beyond the traditional considerations of classical
273 training optimization perspectives.

274 6 Experimental Results

275 This section aims to empirically validate our theoretical findings using real-world benchmark datasets.
276 We start by laying out the used experimental, then we study the impact of various initialization
277 strategies on a GCN’s robustness. Next, we analyze the influence of training epochs on adversarial
278 robustness. Finally, we extend our experimentation to considered family of DNNs in Section 5.

279 6.1 Experimental Setting

280 **Experimental Setup.** Consistent with our theoretical analysis, this section focuses on the node
281 classification task. We leverage the citation networks Cora and CiteSeer [26], with additional results
282 on other datasets provided in the Appendix G. To mitigate the impact of randomness during training,
283 each experiment was repeated 10 times, using the train/validation/test splits provided with the datasets.
284 A 2-layers GCN classifier with identical hyperparameters and activation functions was employed
285 across all the experiments. The models were trained using the cross-entropy loss function, and
286 consistent values for the number of epochs and learning rate were maintained across all analysis.
287 Further implementation details can be found in Appendix H and the code implementation to replicate
288 our experiments is provided in the supplementary material.

289 **Adversarial Attacks.** We consider two main gradient-based structural adversarial attacks: (i)
290 ‘Mettack’ (with the ‘Meta-Self’ training strategy) [40] that formulates the problem as a bi-level
291 problem solved using meta-gradients (ii) and the Proximal Gradient Descent (PGD) [34] which
292 consists of iteratively adding small crafted perturbations using the gradient of the classifier’s loss.
293 We additionally provide results for the ‘Dice’ attack [40] in Appendix G. For our experiments, we
294 considered perturbation rates ranging from 10% (i. e., $0.1E$) to 40% (i. e., $0.4E$).

295 **Evaluation Metrics.** We report the experimental findings in terms of the ‘Attacked Accuracy’, which
296 is the model’s test accuracy when subject to the attacks. Additionally, given that initialization have an
297 impact on the model’s generalization and performance, solely reporting the attacked accuracy fails in
298 some specific cases to provide a comprehensive perspective. Thus, we adopt for some experiments
299 the “Success Rate” metric, also commonly employed in adversarial literature, which encompasses the
300 number of successfully attacked nodes while taking into account the model’s initial clean accuracy.

301 6.2 Effect Of Training Epochs

302 The theoretical analysis presented in Section 4 established a connection between the number of
303 training epochs and the model’s resultant robustness. The derived bound suggests that increasing
304 the number of epochs results in the model becoming more vulnerable to adversarial attacks. The
305 objective of this experimental section is to empirically validate this assertion using real-world datasets.
306 To this end, at each training epoch, we assess the model’s performance on the test set, considering
307 both its clean accuracy and its accuracy under adversarial attacks.

308 Figure 1 illustrates the results of this analysis. The initial two subplots (a,b) displays the findings on
309 the Cora dataset, while the subsequent (c,d) subplots presents results from the CiteSeer dataset. For
310 each dataset, the first plot showcases the clean and attacked accuracy, while the second plot shows
311 the Success Rate (the discrepancy between the clean and attacked accuracy for each budget). The

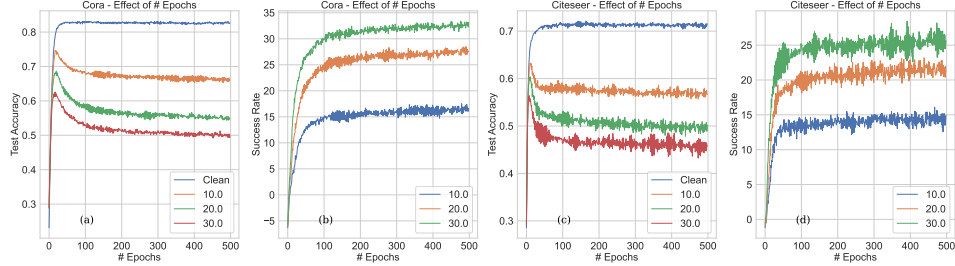


Figure 1: Effect of training epochs on the model’s robustness on Cora (a,b) and CiteSeer (c,d).

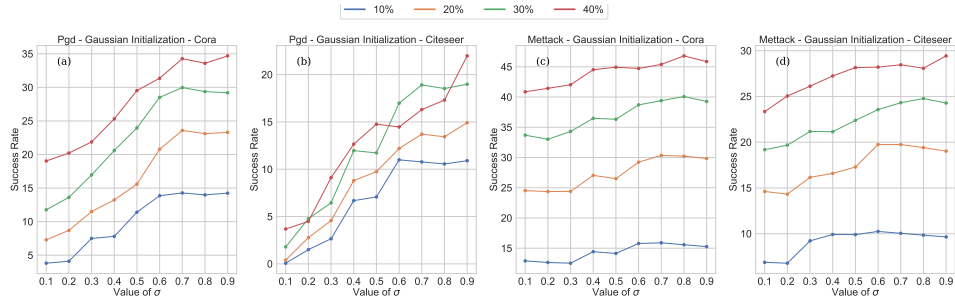


Figure 2: Effect of the variance parameter on the model’s robustness in the case of Gaussian Initialization using on PGD [on Cora (a) and CiteSeer (b)] and Mettack [on Cora (a) and CiteSeer (b)].

312 experimental results demonstrate the existence of the previously discussed trade-off between clean and
 313 robust accuracies. Specifically, as anticipated, the clean accuracy exhibits a continual increase until
 314 reaching a plateau, corresponding to the convergence of the loss function to a minimum. Conversely,
 315 the attacked accuracy demonstrates a rising trend until reaching an inflection point, beyond which it
 316 begins to decline. These findings confirms the observations from the derived upper-bound, indicating
 317 that a higher number of epochs leads to increased vulnerability in the model. Ideally, users would
 318 aim to stop training at the inflection point, where the attacked accuracy is maximized while the clean
 319 accuracy remains proximal to its convergence point.

320 6.3 Effect Of Initial Weight Distribution

321 We aim to validate the impact of the initial weight norms on the model’s adversarial robustness. As
 322 previously discussed in Section 4, a larger weight norm leads to the relaxation of the upper-bound,
 323 potentially resulting in the model being more susceptible to adversarial attacks.

324 In this perspective, we start by investigating the effect of sampling from a Gaussian distribution, as
 325 outlined in Lemma 4. We hence consider this latter by setting the mean value μ to a constant, and
 326 analyzing the impact of the variance parameter σ . Intuitively, based on the upper-bound analysis, a
 327 higher variance value is anticipated to result in reduced model robustness. Figure 2 illustrates the
 328 resultant Success Rate across various variance values for both the "PGD" and "Mettack" methods,
 329 applied to the Cora and CiteSeer datasets. The findings unequivocally validate the theoretical insights,
 330 demonstrating a direct correlation between increasing the variance (σ) and a higher Success Rates,
 331 indicating heightened vulnerability and reduced robustness of the model. Moreover, the impact of
 332 initialization becomes more pronounced when considering larger attack budgets, as outlined in the
 333 computed upper-bound. Notably, for certain budgets (e.g., 30% and 40%), the observed gap ranges
 334 between 5% and 15%, underscoring the initial weights significant implications on the robustness.

335 Within the same context, we explore alternative initialization strategies, focusing on two primary
 336 cases. First, we investigate sampling initial weights from a uniform distribution $\mathcal{U}(-\beta, \beta)$, where
 337 β can be seen as a scaling parameter for weight norms. Second, we consider employing a scaled
 338 orthogonal weight initialization strategy. While this our aim can be approached by sampling weights
 339 from a scaled random Gaussian distribution, we adopt the orthogonal initialization strategy proposed
 340 in prior work [25], which we further rescale by a factor β to examine the impact on weight norms.

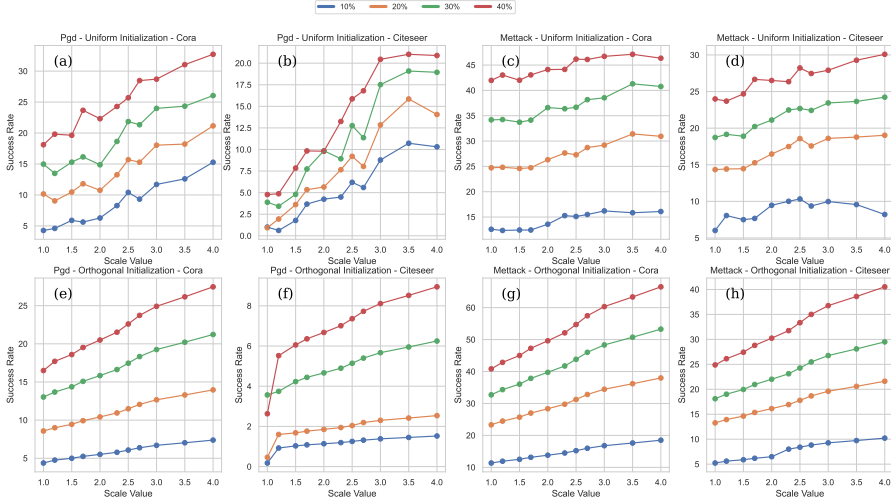


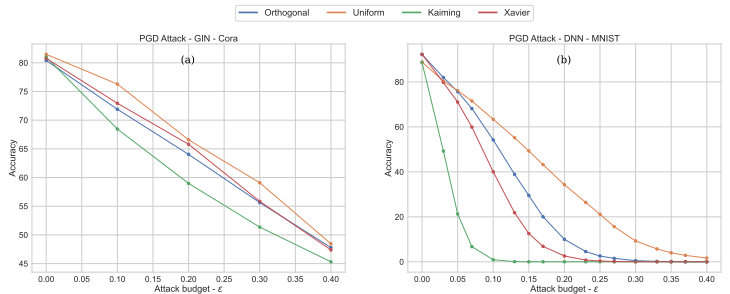
Figure 3: Effect of the scaling parameter β on the model’s robustness in the case of Uniform (a-d) and Orthogonal (e-h) Initialization when subject to PGD and Mettack using Cora and CiteSeer.

341 In both cases, higher scaling parameter values of β are anticipated to theoretically yield higher
 342 upper-bounds and consequently render the model more vulnerable, as indicated by our computed
 343 bounds. We conduct numerical computations on both the Cora and CiteSeer datasets to assess the
 344 resulting adversarial robustness of a GCN across various β values, as provided in Figure 3. The
 345 experimental results are exactly aligned with our theoretical findings showcasing the effect of the
 346 weight norm in the adversarial robustness. To summarize, while traditionally overlooked in prior
 347 studies on adversarial robustness, our experimentation underscores the critical importance of selecting
 348 appropriate initialization distributions and strategies for enhancing model robustness.

349 6.4 Experimental Generalization

350 We extend our experimenta-
 351 tion to empirically validate
 352 the theoretical generalizations
 353 provided in both Section 4 for
 354 the GINs and Section 5 for a
 355 DNNs. To this end, we con-
 356 sider these two models with
 357 various initialization schemes,
 358 including the previously used
 359 Orthogonal [25] and Uniform
 360 initialization in addition to the
 361 Kaiming [16] and Xavier Ini-
 362 tialization [12]. Our analysis
 363 primarily focuses on the PGD
 364 adversarial attack, using iden-
 365 tical attack budgets as in the
 366 previous sections. Figure 4
 367 presents the results on the
 368 GIN (a) using the Cora data-
 369 set and (b) on the DNN using
 370 the MNIST dataset. Notably,
 we observe that the different
 initialization methods yield
 similar clean accuracy ($\epsilon = 0$),
 yet as the attack budget in-
 creases, the discrepancy in
 attacked accuracy between the
 best and worst initialization
 methods for $\epsilon = 0.1$ ranges
 around 60%, proving our
 main assumption related to
 the impact of initialization on
 the model’s robustness.

Figure 4: Effect of initialization on the GIN (a) and DNN (b) for different attack budgets.



7 Conclusion & Limitations

The current study shows that the dynamics of learning in GNNs and DNNs have an important effect on the model’s final robustness. Specifically, we theoretically showed that the model’s robustness is connected to the weight initialization and the number of training epochs. We empirically validate our findings, where we can see that choosing the right initialization can yield huge “almost-free” robustness improvement. We additionally showed the existence of a trade-off between choosing the right number of epochs to have the best clean accuracy and the most robust model. While the current work didn’t propose an alternative or a solution, it has introduced a new perspective, which in our knowledge, was absent from the adversarial literature, opening the door to new research direction either by proposing new initialization scheme to improve robustness while guaranteeing a good generalization or new gradient-based weight updates to enforce the robustness of the model.

References

- [1] Yassine Abbahaddou, Sofiane Ennadir, Johannes F. Lutzeyer, Michalis Vazirgiannis, and Henrik Boström. Bounding the expected robustness of graph neural networks subject to node feature attacks. In *The Twelfth International Conference on Learning Representations*, 2024.
- [2] Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In *International Conference on Machine Learning*, pages 291–301. PMLR, 2019.
- [3] Aleksandar Bojchevski and Stephan Günnemann. Certifiable robustness to graph perturbations, 2019.
- [4] Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more, 2020.
- [5] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. In *International conference on machine learning*, pages 854–863. PMLR, 2017.
- [6] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial Attack on Graph Structured Data. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1115–1124, 2018.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [8] Sofiane Ennadir, Yassine Abbahaddou, Johannes F. Lutzeyer, Michalis Vazirgiannis, and Henrik Boström. A simple and yet fairly effective defense for graph neural networks, 2024.
- [9] Sofiane Ennadir, Amr Alkhatib, Giannis Nikolentzos, Michalis Vazirgiannis, and Henrik Boström. Unboundattack: Generating unbounded adversarial attacks to graph neural networks. In *International Conference on Complex Networks and Their Applications*, pages 100–111. Springer, 2023.
- [10] 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 international conference on web search and data mining*, pages 169–177, 2020.
- [11] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.

- 419 [13] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 420 [14] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adver-
421 sarial examples. In *International Conference on Learning Representations (ICLR)*, 2015.
- 422 [15] Stephan Günnemann. Graph neural networks: Adversarial robustness. In *Graph Neural*
423 *Networks: Foundations, Frontiers, and Applications*, pages 149–176. Springer, 2022.
- 424 [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers:
425 Surpassing human-level performance on imagenet classification. In *2015 IEEE International*
426 *Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015.
- 427 [17] Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in
428 optimizing deep linear networks. In *International Conference on Learning Representations*,
429 2020.
- 430 [18] Hongwei Jin, Zhan Shi, Venkata Jaya Shankar Ashish Peruri, and Xinhua Zhang. Certified
431 robustness of graph convolution networks for graph classification under topological attacks. In
432 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural*
433 *Information Processing Systems*, volume 33, pages 8463–8474. Curran Associates, Inc., 2020.
- 434 [19] Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. Molecular
435 graph convolutions: moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*,
436 30(8):595–608, 2016.
- 437 [20] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional
438 Networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- 439 [21] Xiaoxiao Ma, Jia Wu, Shan Xue, Jian Yang, Chuan Zhou, Quan Z. Sheng, Hui Xiong, and
440 Leman Akoglu. A comprehensive survey on graph anomaly detection with deep learning. *IEEE*
441 *Transactions on Knowledge and Data Engineering*, pages 1–1, 2021.
- 442 [22] Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. The emergence of spectral
443 universality in deep networks. In *International Conference on Artificial Intelligence and*
444 *Statistics*, pages 1924–1932. PMLR, 2018.
- 445 [23] Aymen Qabel, Sofiane Ennadir, Giannis Nikolentzos, Johannes F. Lutzeyer, Michail Chatzianas-
446 tasis, Henrik Boström, and Michalis Vazirgiannis. Advancing antibiotic resistance classification
447 with deep learning using protein sequence and structure. *bioRxiv*, 2023.
- 448 [24] Leslie Rice, Anna Bair, Huan Zhang, and J Zico Kolter. Robustness between the worst and
449 average case. *Advances in Neural Information Processing Systems*, 34:27840–27851, 2021.
- 450 [25] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear
451 dynamics of learning in deep linear neural networks, 2014.
- 452 [26] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
453 Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- 454 [27] Xianfeng Tang, Yandong Li, Yiwei Sun, Huaxiu Yao, Prasenjit Mitra, and Suhang Wang.
455 Transferring robustness for graph neural network against poisoning attacks. In *Proceedings of*
456 *the 13th International Conference on Web Search and Data Mining*. ACM, jan 2020.
- 457 [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
458 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
459 *processing systems*, 30, 2017.
- 460 [29] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
461 Bengio. Graph Attention Networks. In *ICLR*, 2018.
- 462 [30] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. Het-
463 erogeneous graph attention network. In *The world wide web conference*, pages 2022–2032,
464 2019.

- 465 [31] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. Ad-
466 versarial examples for graph data: Deep insights into attack and defense. In *Proceedings of*
467 *the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages
468 4816–4823. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- 469 [32] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. Session-based
470 Recommendation with Graph Neural Networks. In *Proceedings of the 33rd AAAI Conference*
471 *on Artificial Intelligence*, pages 346–353, 2019.
- 472 [33] Lechao Xiao, Yasaman Bahri, Jascha Sohl-Dickstein, Samuel Schoenholz, and Jeffrey Pen-
473 nington. Dynamical isometry and a mean field theory of cnns: How to train 10,000-layer
474 vanilla convolutional neural networks. In *International Conference on Machine Learning*, pages
475 5393–5402. PMLR, 2018.
- 476 [34] Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin.
477 Topology attack and defense for graph neural networks: An optimization perspective. 2019.
- 478 [35] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural
479 Networks? In *7th International Conference on Learning Representations*, 2019.
- 480 [36] Haoxi Zhan and Xiaobing Pei. Black-box Gradient Attack on Graph Neural Networks: Deeper
481 Insights in Graph-based Attack and Defense. *arXiv preprint arXiv:2104.15061*, 2021.
- 482 [37] Xiang Zhang and Marinka Zitnik. Gnn-guard: Defending graph neural networks against
483 adversarial attacks, 2020.
- 484 [38] Dingyuan Zhu, Ziwei Zhang, Peng Cui, and Wenwu Zhu. Robust graph convolutional networks
485 against adversarial attacks. In *Proceedings of the 25th ACM SIGKDD international conference*
486 *on knowledge discovery & data mining*, pages 1399–1407, 2019.
- 487 [39] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial Attacks on Neural
488 Networks for Graph Data. In *Proceedings of the 24th ACM SIGKDD International Conference*
489 *on Knowledge Discovery & Data Mining*, pages 2847–2856, 2018.
- 490 [40] Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta
491 learning. In *7th International Conference on Learning Representations*, 2019.
- 492 [41] Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph
493 convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on*
494 *Knowledge Discovery & Data Mining*. ACM, jul 2019.

Supplementary Material: If You Want to Be Robust, Be Wary of Initialization

495 A Proof Of Theorem 2

496 **Theorem.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers, where
 497 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. For adversarial attacks only targeting
 498 node features of the input graph, with a budget ϵ , we have (in respect to Definition 1):

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

499 with t being the number of training epochs and \hat{w}_u denoting the sum of normalized walks of length
 500 $(T - 1)$ starting from node u .

501 *Proof.* Let's consider a graph-function f that is based on T GCN-layers. The gradient descent update
 502 at epoch t for a layer i is written as:

$$W_{t+1}^{(i)} = W_t^{(i)} - \eta \nabla \mathcal{L}(W_t^{(i)})$$

503 Since we consider that our loss function \mathcal{L} to be L -smooth, we have the following result:

$$\|\nabla \mathcal{L}(W_t^{(i)})\| \leq L \|W_t^{(i)} - W_*^{(i)}\|$$

504 Consequently, after t training epochs, we can write:

$$\begin{aligned} \|W_t^{(i)}\| &= \|W_{t-1}^{(i)} - \eta \nabla \mathcal{L}(W_{t-1}^{(i)})\| \\ &\leq \|W_{t-1}^{(i)}\| + \eta L \|W_{t-1}^{(i)} - W_*^{(i)}\| \\ &\leq (1 + \eta L) \|W_{t-1}^{(i)}\| + \eta L \|W_*^{(i)}\| \end{aligned}$$

505 In addition, we have that $\eta \leq \frac{1}{L}$. Hence, by recursion, we find that:

$$\|W_t^{(i)}\| \leq (1 + \eta L)^t \|W_0^{(i)}\| + \sum_{h=0}^{t-1} 2^h \|W_*^{(i)}\| \quad (3)$$

$$\leq (1 + \eta L)^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \quad (4)$$

506 Giving that we are considering feature-based adversarial attacks, let X denote the original node
 507 features and X' denote the perturbed adversarial features. With an attack budget ϵ , from the work [1],
 508 we have the following result:

$$\forall [A, X'] \in B([A, X], \epsilon), \|f(A, X) - f(A, X')\| \leq \prod_{i=1}^T \|W_t^{(i)}\| \epsilon \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right). \quad (5)$$

509 with \hat{w}_u denoting the sum of normalized walks of length $(T - 1)$ starting from node u . Consequently:

$$\sup_{[A, X'] \in B([A, X], \epsilon)} \|f(A, X) - f(A, X')\| \leq \prod_{i=1}^T \|W_t^{(i)}\| \epsilon \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right). \quad (6)$$

510 From Result 3 and 6, we conclude that:

$$\sup_{[A, X'] \in B([A, X], \epsilon)} \|f(A, X) - f(A, X')\| \leq \epsilon \prod_{i=1}^T [2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\|] \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

511 We conclude that f is $(\epsilon; \gamma)$ -robust with:

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

512

□

513 B Proof Of Theorem 3

514 **Theorem.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers, where
 515 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. Let f be the number of used training
 516 epochs. When f is subject to structural attacks, with a budget ϵ , we have (in respect to Definition 1):

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \|X\| \left(1 + T \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \right)$$

517 *Proof.* Similar to the previous proof, let's consider a graph-function f that is based on T GCN-layers
 518 and trained using gradient descent for t epochs. We have the following result from Equation 3:

$$\|W_t^{(i)}\| \leq 2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \quad (7)$$

519 For this proof, we are considering the model f to be subject to structural perturbations. In this
 520 perspective, let \tilde{A} denote the input non-attacked adjacency and \tilde{A}' denote the attacked/perturbed
 521 adjacency, with h' denoting its corresponding hidden representation. From the work [1], we have:

$$\forall [A', X] \in B([A, X], \epsilon), \|f(\tilde{A}, X) - f(\tilde{A}', X)\| \leq \prod_{i=1}^T \|W^{(i)}\| \|X\| \epsilon \left(1 + T \prod_{i=1}^T \|W^{(i)}\| \right)$$

522 By combining the two previous results, we get that following inequality and hence the desired result:

$$\sup_{[A', X] \in B([A, X], \epsilon)} \|f(\tilde{A}, X) - f(\tilde{A}', X)\| \leq \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \|X\| \left(1 + T \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \right).$$

523

□

524 C Proof Of Lemma 4

525 **Lemma.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GCN layers for which
 526 the initial weight are drawn from the Gaussian distribution $\mathcal{N}(\mu, \Sigma)$. When subject to node features
 527 based adversarial attacks, we have the following:

$$\mathbb{E}_{W_0 \sim \mathcal{N}(\mu, \Sigma)} [\mathcal{R}_\epsilon[f]] \leq \epsilon \prod_{i=1}^T \left(2^t \sqrt{\mu^2 + \text{tr}(\Sigma)} + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

528 *Proof.* Let's consider f to be a graph classifier based on T -GCN layers for which the initial weight
 529 are drawn from the Gaussian distribution. Specifically, $\forall i \leq L, W_0^{(i)} \sim \mathcal{N}(\mu, \Sigma)$. We have that:

$$\mathbb{E}[\|W_0^{(i)}\|] \leq \sqrt{\|\mu\|^2 + \text{tr}(\Sigma)}$$

530 From Theorem 2, we have the following:

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

531 Hence, combining the two elements results in the following:

$$\mathbb{E}_{W_0 \sim \mathcal{N}(\mu, \Sigma)} [\mathcal{R}_\epsilon[f]] \leq \epsilon \prod_{i=1}^T \left(2^t \sqrt{\mu^2 + \text{tr}(\Sigma)} + 2^{t+1} \|W_*^{(i)}\| \right) \left(\sum_{u \in \mathcal{V}} \hat{w}_u \right)$$

532

□

533 D Proof Of Theorem 5

534 **Theorem.** Let $f : (\mathcal{A}, \mathcal{X}) \rightarrow \mathcal{Y}$ denote a graph-based function composed of T GIN layers, where
 535 the initial weight matrix of the i -th layer is denoted by $W_0^{(i)}$. For adversarial attacks only targeting
 536 node features of the input graph, with a budget ϵ , we have:

$$\gamma = \prod_{l=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left[BT \max_{u \in \mathcal{V}} \text{deg}(u) + \epsilon \right]$$

537 with t being the number of training epochs and $\text{deg}(u)$ is the degree of node u .

538 *Proof.* Let's consider a graph-function f that is based on T GIN-layers and trained using gradient
 539 descent for t epochs. We have the following result from Equation 3:

$$\|W_t^{(i)}\| \leq (1 + \eta L)^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \leq 2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \quad (8)$$

540 Let X denote the original node features and X' the perturbed adversarial features. For an attack
 541 budget ϵ , from the work [1], we have the following:

$$\forall [A', X] \in B([A, X], \epsilon), \|f(A, X) - f(A, X')\| \leq \prod_{l=1}^T \|W^{(l)}\| [B \times T \times \max_{u \in \mathcal{V}} \text{deg}(u) + \epsilon] \quad (9)$$

542 Consequently, we can merge the two inequalities resulting in the following:

$$\gamma = \prod_{l=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right) \left[B \times T \times \max_{u \in \mathcal{V}} \text{deg}(u) + \epsilon \right]$$

543

□

544 E Proof Of Theorem 6

545 **Theorem.** Let $f : \mathcal{X} \subseteq \mathbf{R}^{in} \rightarrow \mathcal{Y} \subseteq \mathbf{R}^{out}$ be a T -layers neural network with $W_0^{(i)}$ denoting the
 546 initial weight matrix of the i -th layer. When subject to adversarial attacks, f is (ϵ, γ) -robust with:

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right)$$

547 *Proof.* Let f be a T -layers neural network. We additionally assume that its corresponding activation
548 functions are 1-Lipschitz. Let x (with h its hidden representation) be an input vector and x' (corresp.
549 h') its corresponding crafted adversarial input (corresp. hidden representation). For an adversarial
550 attack with budget ϵ , we have the following:

$$\begin{aligned} \forall x' \in \mathcal{X} : \|x - x'\| \leq \epsilon, \|f(x) - f(x')\| &= \|h^{(l)} - h'^{(l)}\| \\ &= \|\phi^{(l)}(W^{(l)}h^{(l-1)} + b^{(l)}) - \phi^{(l)}(W^{(l)}h'^{(l-1)} + b^{(l)})\| \\ &\leq \|W^{(l)}\| \|h^{(l-1)} - h'^{(l-1)}\| \end{aligned}$$

551 Recurrently, we find the final result as:

$$\sup_{x' \in \mathcal{X} : \|x - x'\| \leq \epsilon} \|f(x) - f(x')\| \leq \prod_{l=1}^T \|W^{(l)}\| \epsilon \quad (10)$$

552 Note that similar results and analysis have been provided in previous work [5, 2]. By using the result
553 derived in Equation 3, we have:

$$\|W_t^{(i)}\| \leq 2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \quad (11)$$

554 By merging these two inequalities, and applying the Markov Inequality, we find the following
555 upper-bound:

$$\gamma = \epsilon \prod_{i=1}^T \left(2^t \|W_0^{(i)}\| + 2^{t+1} \|W_*^{(i)}\| \right)$$

556

□

557 F On the case of strong-convexity - Proof of Lemma 7

558 **Lemma.** Let $f : \mathcal{X} \subseteq \mathbf{R}^{in} \rightarrow \mathcal{Y} \subseteq \mathbf{R}^{out}$ be a T -layers neural network trained with a μ -strongly
559 convex and L -smooth loss function. Let $W_0^{(i)}$ denote the initial weight matrix of the i -th layer. When
560 subject to adversarial attacks, with a budget ϵ , we have that f is (ϵ, γ) -robust with:

$$\gamma = \epsilon \prod_{i=1}^T \left((1 - \mu/L)^t \|W_0^{(i)}\| + 2 \|W_*^{(i)}\| \right)$$

561 *Proof.* We consider f to be a T -layers neural network (following the same propagation as equation
562 the one presented in Section 5). From Section E, we have the following:

$$\|f(x) - f(x')\| \leq \prod_{l=1}^T \|W^{(l)}\| \epsilon$$

563 In addition to the previous assumption of L -smoothness of the loss function, we consider that its
564 μ -strongly convex. Hence, for the layer (l) , we have the following result:

$$\|W_t^{(l)}\| \leq (1 - \mu/L)^t \|W_0^{(l)} - W_*^{(l)}\| + \|W_*^{(l)}\| \quad (12)$$

$$\leq (1 - \mu/L)^t \|W_0^{(l)}\| + 2 \|W_*^{(l)}\| \quad (13)$$

565 When subject to adversarial attacks, we can use the previous result from E, specifically from Equa-
566 tion10:

$$\sup_{x' \in \mathcal{X} : \|x - x'\| \leq \epsilon} \|f(x) - f(x')\| \leq \prod_{l=1}^T \|W^{(l)}\| \epsilon \quad (14)$$

567 Hence, by merging the two previous results, we deduce that:

$$\gamma = \epsilon \prod_{i=1}^T \left((1 - \mu/L)^t \|W_0^{(i)}\| + 2\|W_*^{(i)}\| \right) \quad (15)$$

568

□

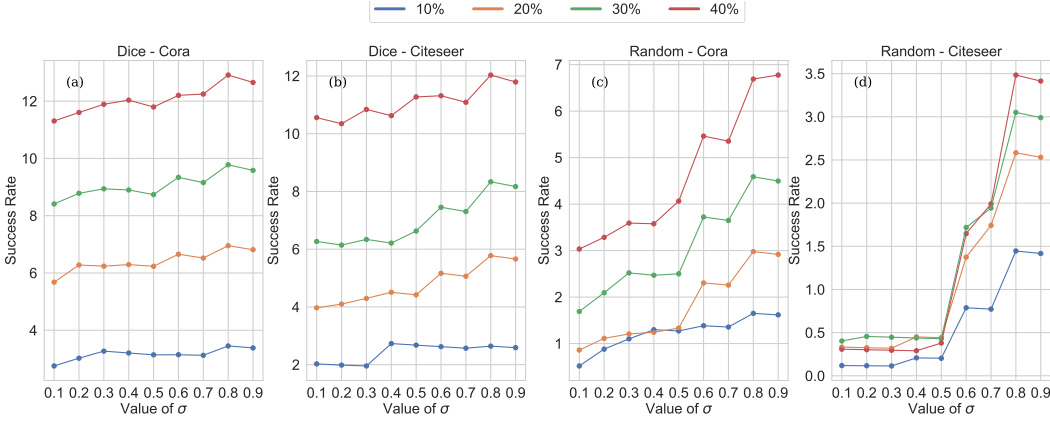


Figure 5: Effect of the variance on the model’s robustness in the case of Gaussian Initialization when subject to DICE (a,b) and Random Attacks (c,d) for both Cora and CiteSeer.

569 G Additional Results

570 G.1 Adversarial Robustness of Deep Neural Networks

571 We consider the general family of neural networks for which the computation during layer l , using an
 572 activation function $\phi^{(l)}$, can be written as :

$$h^{(l)} = \phi^{(l)}(W^{(l)}h^{(l-1)} + b^{(l)}).$$

573 with $W^{(l)} \in \mathbb{R}^{n_{l-1}, n_l}$ being the weight matrix and $b_l \in \mathbb{R}^{n_l}$ the bias of the l^{th} layer.

574 In this perspective, let $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}$ be a neural network n_0 being the input dimension. The
 575 adversarial task in this case consists of finding a perturbed input \tilde{x} for which the prediction differs
 576 from the original prediction $f(x)$. The perturbed input \tilde{x} should hence adhere to the similarity
 577 constraints defined by a perturbation budget ϵ . Let’s consider the ℓ_2 norm within both the input space
 578 \mathbb{R}^{n_0} and the output space \mathbb{R} , we can hence define the set of valid adversarial perturbation as:

$$579 \quad B(x; \epsilon) = \{\tilde{x} : \|x - \tilde{x}\| \leq \epsilon\}$$

580 Similar to Section 3, we can introduce the adversarial risk of a DNN within the input’s neighborhood
 581 defined by the budget ϵ as the following:

$$\mathcal{R}_\epsilon[f] = \mathbb{E}_{\substack{x \sim \mathcal{D} \\ \tilde{x} \in B(x; \epsilon)}} [|(f(\tilde{x}) - f(x))|]. \quad (16)$$

582 From this adapted adversarial risk, we can introduce the notion of a DNN’s adversarial robustness

583 **Definition 8.** (DNN - Adversarial Robustness). The neural network $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}$ is said to be
 584 (ϵ, γ) - robust if its adversarial risk is upper-bounded by γ , i. e., $\mathcal{R}_\epsilon[f] \leq \gamma$.

585 **G.2 Additional Adversarial Attacks**

586 In addition to the previously reported Mettack and PGD adversarial attack, we consider two additional
 587 adversarial attacks. Notably, we first consider "DICE" which involves iteratively perturbing a graph's
 588 structure by adding or removing edges while ensuring connectivity, and then adjusting the perturbation
 589 based on the gradient of the graph neural network's loss function to generate an adversarial example.
 590 The process aims to find a minimal perturbation that misleads the network's predictions while keeping
 591 the perturbation size small. We additionally consider a "Random" attack which consists of randomly
 592 perturbing the adjacency matrix by dropping or adding edges. Figure 5 shows the adversarial accuracy
 593 results on the Cora and CiteSeer dataset when subject to DICE and Random attacks for different
 594 values of σ of the Gaussian initialization. Similarly, Figure 6 shows the effect of scaling both a
 595 uniform initialization and an Orthogonal one as previously explained in Section 6.

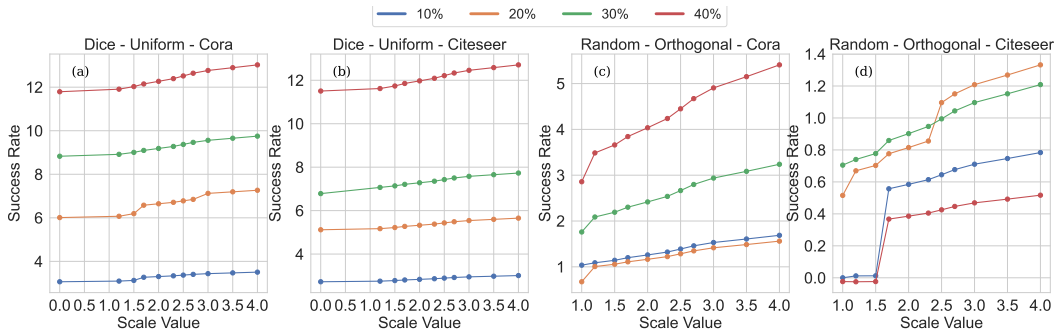


Figure 6: Effect of Uniform and Orthogonal Initialization on the model's robustness in the case of DICE Attack on Cora (a,c) and CiteSeer (b,d).

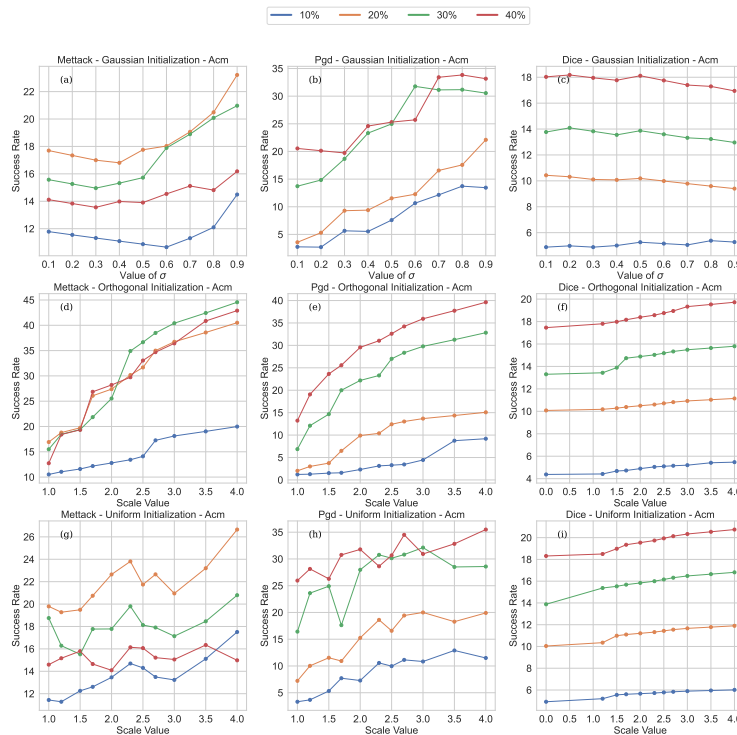


Figure 7: Effect of the Gaussian (a; b; c), Orthogonal (d; e; f) and Uniform (g;h;i) Initialization on the ACM dataset

596 **G.3 Additional Datasets**

597 We additionally extend the results to the ACM Dataset [30] within the node classification setting.
 598 Figure 7 presents the results using the Mettack, PGD and DICE for the ACM dataset for the Gaussian
 599 initialization (effect of σ), the Uniform and Orthogonal initialization.

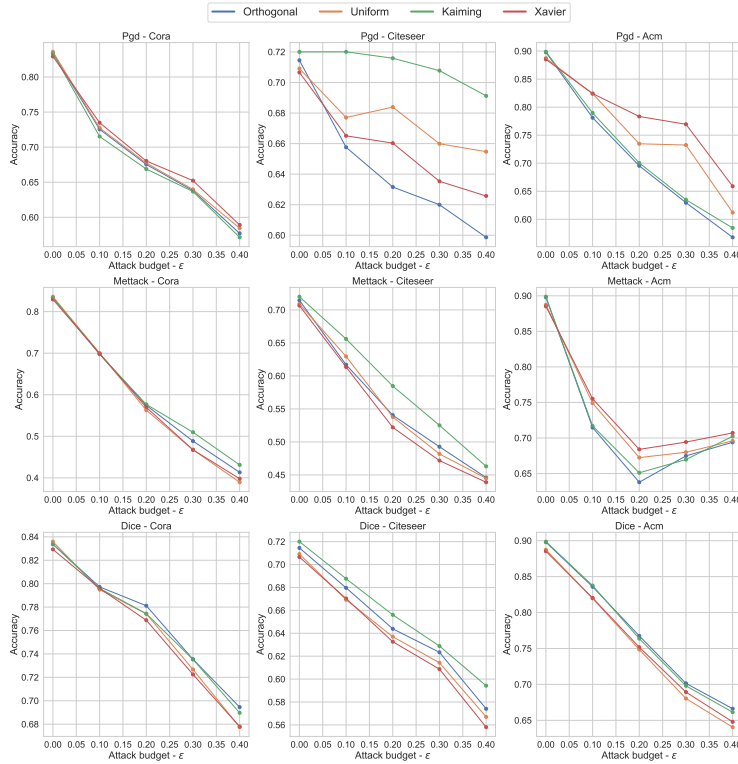


Figure 8: Effect of the initial distribution on RGCN’s robustness and performance when subject to structural adversarial attacks.

600 **G.4 Additional Models**

601 As previously explained in Section 5, while our theoretical analysis primarily focuses on GCN,
 602 GIN, and DNN models, the derived insights extend to other models as well. To illustrate this point,
 603 we examine the effect of initialization distribution on the performance of defense methodologies.
 604 Specifically, we first consider RGCN [38], which employs Gaussian distributions in its hidden
 605 layers to mitigate the effects of adversarial attacks. We additionally consider GCN-Jaccard [31]
 606 which preprocesses the network by eliminating edges that connect nodes with jaccard similarity of
 607 features smaller than a certain level. We use various initialization schemes, similar to those in our
 608 previous experiments, and evaluate against the same adversarial attacks (PGD, Mettack, and DICE).
 609 Figure 8 (resp. Figure 9) presents the adversarial accuracy and defense performance of RGCN (resp.
 610 GCN-Jaccard) on the Cora, CiteSeer, and ACM datasets. Although the performance gap is not very
 611 pronounced for Cora, it is clearly observed for CiteSeer and ACM. This demonstrates the broader
 612 applicability of our insights across different models but also defense methods.

613 **H Datasets and Implementation details**

614 **Datasets** Characteristics and information about the node classification datasets used in our experi-
 615 mental study are presented in Table 1. As outlined in the main paper, we conduct experiments on a
 616 set of citation networks, including Cora, CiteSeer (in the main paper), and ACM dataset (Appendix
 617 G) [30]. For all these datasets, we adhere to the train/valid/test splits provided by with the dataset.

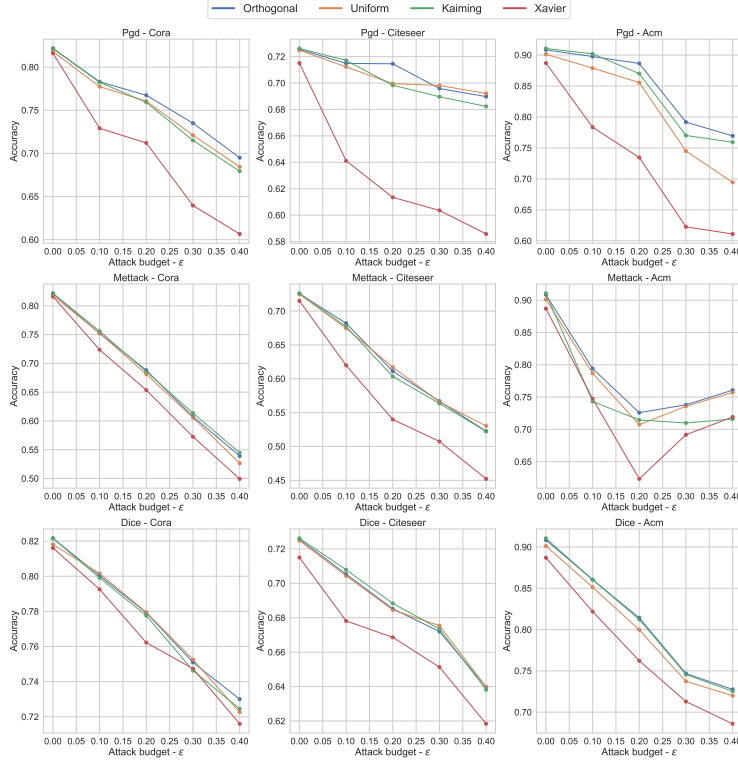


Figure 9: Effect of the initial distribution on GCN-Jaccard’s robustness and performance when subject to structural adversarial attacks.

618 **About the architectures.** In all of the experiments, the models employed a 2-layer convolutional
 619 architecture (consisting of two iterations of message passing and updating) stacked with a Multi-Layer
 620 Perception (MLP) as a readout. The intent was to compare the models in an iso-architectural setting,
 621 to ensure a fair evaluation of their robustness. We maintained the same hyperparameters, including a
 622 learning rate of $1e-2$, 300 epochs, and a hidden feature dimension of 16 have been. To account for the
 623 impact of random initialization, each experiment was repeated 10 times.

624 **Reproducibility of the experiments.** We emphasize that all experiments should be easily repro-
 625 ducible by directly using the provided code. The archive contains a ReadMe file containing a small
 626 documentation on how to run the experiments.

Table 1: Statistics of the node classification datasets used in our experiments.

DATASET	#FEATURES	#NODES	#EDGES	#CLASSES
CORA	1433	2708	5208	7
CITSEER	3703	3327	4552	6

627 **On the adversarial attacks.** For the PGD attack on the MNIST dataset, we used a step-size of
 628 0.1 and we set the number of iterations to 100 (which was observed to be enough for the attack
 629 convergence). Note that we set these parameters for all the considered initializations in Figure 4 as
 630 our aim is to compare the effect of the different distribution on the final robustness.

631 **Implementation details.** Our implementation is available in the supplementary materials (and will
 632 be publicly available afterwards). It is built using the open-source library *PyTorch Geometric* (PyG)
 633 under the MIT license [11]. We used the publicly available implementation of the adversarial attacks
 634 provided in the DeepRobust package (<https://github.com/DSE-MSU/DeepRobust>). For RGCN, we
 635 used the implementation from the same package. The experiments have been run on both a NVIDIA
 636 A100 GPU where training a GCN takes around $1.2(\pm 0.2)$ s.

637 **NeurIPS Paper Checklist**

638 **1. Claims**

639 Question: Do the main claims made in the abstract and introduction accurately reflect the
640 paper's contributions and scope?

641 Answer: [\[Yes\]](#)

642 Justification: In addition to stating the novelty of our proposed approach, we used our
643 abstract and introduction to summarize our main findings and contributions related to the
644 effect of initialization on the adversarial robustness (as theoretically justified and empirically
645 tested in the following sections).

646 Guidelines:

- 647 • The answer NA means that the abstract and introduction do not include the claims
648 made in the paper.
- 649 • The abstract and/or introduction should clearly state the claims made, including the
650 contributions made in the paper and important assumptions and limitations. A No or
651 NA answer to this question will not be perceived well by the reviewers.
- 652 • The claims made should match theoretical and experimental results, and reflect how
653 much the results can be expected to generalize to other settings.
- 654 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
655 are not attained by the paper.

656 **2. Limitations**

657 Question: Does the paper discuss the limitations of the work performed by the authors?

658 Answer: [\[Yes\]](#)

659 Justification: Together with our conclusion, we presented the set of limitations of work.
660 Specifically, we stated that while our work is innovative, we didn't provide a solution to
661 the initialization problem from an adversarial defense perspective. We also discussed in
662 the "problem setup" section our different theoretical choices (the smoothness of the loss
663 function) and how realistic they are.

664 Guidelines:

- 665 • The answer NA means that the paper has no limitation while the answer No means that
666 the paper has limitations, but those are not discussed in the paper.
- 667 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 668 • The paper should point out any strong assumptions and how robust the results are to
669 violations of these assumptions (e.g., independence assumptions, noiseless settings,
670 model well-specification, asymptotic approximations only holding locally). The authors
671 should reflect on how these assumptions might be violated in practice and what the
672 implications would be.
- 673 • The authors should reflect on the scope of the claims made, e.g., if the approach was
674 only tested on a few datasets or with a few runs. In general, empirical results often
675 depend on implicit assumptions, which should be articulated.
- 676 • The authors should reflect on the factors that influence the performance of the approach.
677 For example, a facial recognition algorithm may perform poorly when image resolution
678 is low or images are taken in low lighting. Or a speech-to-text system might not be
679 used reliably to provide closed captions for online lectures because it fails to handle
680 technical jargon.
- 681 • The authors should discuss the computational efficiency of the proposed algorithms
682 and how they scale with dataset size.
- 683 • If applicable, the authors should discuss possible limitations of their approach to
684 address problems of privacy and fairness.
- 685 • While the authors might fear that complete honesty about limitations might be used by
686 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
687 limitations that aren't acknowledged in the paper. The authors should use their best
688 judgment and recognize that individual actions in favor of transparency play an impor-
689 tant role in developing norms that preserve the integrity of the community. Reviewers
690 will be specifically instructed to not penalize honesty concerning limitations.

691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: For each Theorem, Lemma and theoretical claim, we provide the proof in the Appendix and point out to the corresponding section in the main paper. We also stated all the assumptions and analytical choices in the Preliminaries (Section 3.1)

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: In addition to providing the code as supplementary materials, we have provided all the implementations details that are sufficient to reproduce the results. These details include the used hyper-parameters (the architecture, learning rate . . .) and also for the used adversarial attacks we provide the different parameters used. We also point out the dataset that we used (which are public) and that we used the same public folds as the one provided with the datasets.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- 743 (c) If the contribution is a new model (e.g., a large language model), then there should
744 either be a way to access this model for reproducing the results or a way to reproduce
745 the model (e.g., with an open-source dataset or instructions for how to construct
746 the dataset).
- 747 (d) We recognize that reproducibility may be tricky in some cases, in which case
748 authors are welcome to describe the particular way they provide for reproducibility.
749 In the case of closed-source models, it may be that access to the model is limited in
750 some way (e.g., to registered users), but it should be possible for other researchers
751 to have some path to reproducing or verifying the results.

752 5. Open access to data and code

753 Question: Does the paper provide open access to the data and code, with sufficient instruc-
754 tions to faithfully reproduce the main experimental results, as described in supplemental
755 material?

756 Answer: [Yes]

757 Justification: We provide the anonymized code following the Neurips guidelines. Specifi-
758 cally, we submitted the code with the supplementary material section and we clearly state
759 the steps to run it using a ReadMe file. Please note that for this question, we consider "open
760 source" as providing the code to the reviewers and making it public afterwards for the public.

761 Guidelines:

- 762 • The answer NA means that paper does not include experiments requiring code.
- 763 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
764 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 765 • While we encourage the release of code and data, we understand that this might not be
766 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
767 including code, unless this is central to the contribution (e.g., for a new open-source
768 benchmark).
- 769 • The instructions should contain the exact command and environment needed to run to
770 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
771 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 772 • The authors should provide instructions on data access and preparation, including how
773 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 774 • The authors should provide scripts to reproduce all experimental results for the new
775 proposed method and baselines. If only a subset of experiments are reproducible, they
776 should state which ones are omitted from the script and why.
- 777 • At submission time, to preserve anonymity, the authors should release anonymized
778 versions (if applicable).
- 779 • Providing as much information as possible in supplemental material (appended to the
780 paper) is recommended, but including URLs to data and code is permitted.

781 6. Experimental Setting/Details

782 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
783 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
784 results?

785 Answer: [Yes]

786 Justification: We provided all the details about the architecture, the used hyper-parameters
787 for the considered models (Section H of the Appendix) and all the hyper-parameters used
788 for our adversarial attacks. Note that our work's goal is to provide comprehensive overview
789 of the effect of initialization on the robustness, hence making sure that the same choice of
790 hyper-parameters is enough to ensure the fairness of the experiments.

791 Guidelines:

- 792 • The answer NA means that the paper does not include experiments.
- 793 • The experimental setting should be presented in the core of the paper to a level of detail
794 that is necessary to appreciate the results and make sense of them.
- 795 • The full details can be provided either with the code, in appendix, or as supplemental
796 material.

797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We reproduce each experiment 10 times to take into account the factor of randomization and we report the mean value. Note that since we use mainly figures (which are appropriate for our setting – given the different attack budgets we are using), this seemed as the perfect approach. For the train/test folds, we use the public folds provided with each dataset and hence reducing the effect of randomization.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We reported the details of implementation in Section H of the Appendix, where we specified the GPU that was used and the average time to do the experiments. Note that while we have chosen to use a GPU, our experiments can be easily done using a CPU.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We follow the guidelines of the Neurips Code of Ethics.

849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We provided overview on the harm that adversarial attacks can have on the applications of Deep Learning models. The main goal of our paper is to identify new potential factors related to adversarial attacks and hence should rather have a positive impact on the society.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: In this work, we study the theoretical effect of initialization on the adversarial robustness. We don't provide any new pre-trained model nor new datasets.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

901 • We recognize that providing effective safeguards is challenging, and many papers do
902 not require this, but we encourage authors to take this into account and make a best
903 faith effort.

904 12. Licenses for existing assets

905 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
906 the paper, properly credited and are the license and terms of use explicitly mentioned and
907 properly respected?

908 Answer: [Yes]

909 Justification: We made sure to cite the papers that are relevant to our work and that were used
910 to justify some theoretical or empirical insights. For the different code implementations, we
911 cited clearly the license and the owner of the used function/code.

912 Guidelines:

- 913 • The answer NA means that the paper does not use existing assets.
- 914 • The authors should cite the original paper that produced the code package or dataset.
- 915 • The authors should state which version of the asset is used and, if possible, include a
916 URL.
- 917 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 918 • For scraped data from a particular source (e.g., website), the copyright and terms of
919 service of that source should be provided.
- 920 • If assets are released, the license, copyright information, and terms of use in the
921 package should be provided. For popular datasets, `paperswithcode.com/datasets`
922 has curated licenses for some datasets. Their licensing guide can help determine the
923 license of a dataset.
- 924 • For existing datasets that are re-packaged, both the original license and the license of
925 the derived asset (if it has changed) should be provided.
- 926 • If this information is not available online, the authors are encouraged to reach out to
927 the asset's creators.

928 13. New Assets

929 Question: Are new assets introduced in the paper well documented and is the documentation
930 provided alongside the assets?

931 Answer: [Yes]

932 Justification: We have provided the implementation code together with all the experimental
933 details to reproduce our work. We also clearly justify the use of the packages and their license.
934 Note that the code have been anonymized and provided as a supplementary materials.

935 Guidelines:

- 936 • The answer NA means that the paper does not release new assets.
- 937 • Researchers should communicate the details of the dataset/code/model as part of their
938 submissions via structured templates. This includes details about training, license,
939 limitations, etc.
- 940 • The paper should discuss whether and how consent was obtained from people whose
941 asset is used.
- 942 • At submission time, remember to anonymize your assets (if applicable). You can either
943 create an anonymized URL or include an anonymized zip file.

944 14. Crowdsourcing and Research with Human Subjects

945 Question: For crowdsourcing experiments and research with human subjects, does the paper
946 include the full text of instructions given to participants and screenshots, if applicable, as
947 well as details about compensation (if any)?

948 Answer: [NA]

949 Justification: There is no crowdsourcing nor research with human subjects in our case.

950 Guidelines:

- 951 • The answer NA means that the paper does not involve crowdsourcing nor research with
952 human subjects.

953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977

- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: There is no crowdsourcing nor research with human subjects in our case.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.