

# ON THE EFFECT OF INPUT PERTURBATIONS FOR GRAPH NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The expressive power of a message passing graph neural network (MPGNN) depends on its architecture and the input node attributes. In this work, we study how this interplay is affected by input perturbations. First, perturbations of node attributes may act as noise and hinder predictive power. But, perturbations can also aid expressiveness, by making nodes more identifiable. Recent works show that unique node IDs are necessary to represent certain functions with MPGNNs. Our results relate properties of the noise, smoothness of the model and the geometry of the input graphs and task. In particular, we take the perspective of lower bounding smoothness for achieving discrimination: how much output variation is needed for exploiting random node IDs, or for retaining discriminability? Our theoretical results imply constraints on the model for exploiting random node IDs, and, conversely, insights into the tolerance of a given model class for retaining discrimination with perturbations of node attributes.

## 1 INTRODUCTION

Graph Neural Networks (GNNs) (Scarselli et al., 2008) are the state-of-the-art representation learning algorithms on graphs. Important variants of GNNs include Graph Attention Networks (GATs) Veličković et al. (2018) and Graph Convolutional Networks (GCNs) Kipf & Welling (2017). They have shown to be quite successful in practice, e.g., in physics Shlomi et al. (2020); Battaglia et al. (2016), predicting drug interactions Lim et al. (2019), protein design Strokach et al. (2020); Ingraham et al. (2021), studying molecules Jin et al. (2018); Duvenaud et al. (2015), quantum chemistry Gilmer et al. (2017), and drug discovery Sun et al. (2020) (see Wu et al. (2020) for a survey on GNNs, and Zhou et al. (2020) for a review on applications).

Despite having numerous successful empirical results, it is observed that GNNs are sensitive to perturbation in both graph structure and node attributes (see e.g., Zugner & Günnemann (2019)), and they are vulnerable to over-smoothing: having deep networks results in poor representations (see e.g., Kipf & Welling (2017); Oono & Suzuki (2019)). Intuitively, deeper networks allow to fit more smooth functions to the given data, along with having hidden representations depending on large neighborhoods. These issues have been addressed in several papers, and subsequently, people also tried to provide a consistent theory for these observation.

The discrimination power or the expressive power of GNNs is also another well-studied topic, and many works notice that having similar initial node attributes at different nodes may affect the power of network to distinguish distinct graphs, as well as counting subgraphs, while node identifiers (IDs), which are typically randomly chosen, can resolve the issues. But how much smoothness/non-smoothness of the model can affect the benefits of those (random) node IDs? Here by the smoothness of the model we means (lower/upper) bounds involving (perhaps higher-order) derivatives of the representation with respect to hidden representations, where here these derivatives assumed to exist.

The goal of this paper is to unify these observations/questions in one theory and relate the expressive power of GNN to the smoothness of the model and the sensitivity of it to the random perturbation of initial node attributes (or having (random) node IDs<sup>1</sup>). Our approach is restrict the class of functions

<sup>1</sup>Also, initial node attributes can be considered random in cases that are task independent (while quite being informative). For example, a GNN running on a social network can be initialized with individuals gender (as initial attributes) while the learning task is independent from the information.

that GNNs can approximate (i.e., the expressive power of GNNs), if the model is robust against random perturbations of initial attributes, and satisfy upper/lower bounds on its (perhaps higher-order) derivatives with respect to hidden representations.

Our main finding is a *quantitative* lower bound relating the following important notions about GNNs to each other: expressive power, geometry of the noiseless, smoothness of the model, and sensitivity to node attributes. In particular, we observe that non-smoothness representations are essential to express many graph classifiers with GNNs while still remaining robust against random perturbation of attributes. This scenario can also happen when node attributes are still informative, but independent of the learning task<sup>2</sup>. In addition, since the class of graph classifiers achieved with smooth and non-sensitive GNNs is completely restricted, we conclude that we should either be non-robust against random perturbations, or we should accept non-smooth representation to break the bound.

In short, in this paper, we make the following contributions:

- We present a theory explaining why restriction the model to smooth GNNs and still demanding robustness against random perturbations of node attributes can dramatically decrease the expressive power of the model. The result is a *quantitative* lower bounds which allows to observe the effects of the discussed notions on the representation in one mathematical equation, which we refer to it as an uncertainty principle for GNNs too.
- We introduce the notion of *discrepancy* of graph classifiers for studying GNNs, study its properties, and relate it to the smoothness, sensitivity, and the expressive power of GNNs.
- Our bound also applies to the expressive power of GNN with random (or task independent) node attributes, which can be of independent interest. In particular, we study how smoothness can affect the expressive power of GNN with random node attributes and the results suggest that while node IDs can improve the expressive power in general, with more smooth GNNs can approximate a restricted class of graph classifiers.

## 2 RELATED WORKS

**Theory of GNNs.** The expressive power of GNNs is limited to the so-called Weisfeiler-Lehman (WL) graph isomorphism test (Xu et al., 2019), and under specific conditions, they are universal approximator (Scarselli et al., 2009) (see (Azizian & marc lelarge, 2021; Sato, 2020) for more on expressive power). Note that function approximation and graph isomorphism testing with GNNs are equivalent (Chen et al., 2019). Beyond graph isomorphism testing, GNNs are known to be unable to count substructures (see for instance (Chen et al., 2020b; You et al., 2019)). To improve the expressive power of GNNs, it is proposed to use random attributes (Sato et al., 2021) or coloring nodes (Dasoulas et al., 2020). We notice that just being able to represent a functions does not mean that the representation is achievable with smooth functions. For example, in the context of learning on sets, it is shown that imposing regularity conditions on the underlying functions results in limitations on representing functions (Wagstaff et al., 2019). In addition, GNNs suffer from the problem of over-smoothing (Oono & Suzuki, 2019). For example, with dropping edges is an instance of defeating against over-smoothing problem (Rong et al., 2020). Recently, generalization and stability bounds for GNNs have been also derived, e.g., see (Garg et al., 2020; Verma & Zhang, 2019).

We note that although GNNs with infinite width and node identifiers are Turing-complete (Loukas, 2020a), depth/width lower bounds show that GNNs should scale to several graph learning problems (Loukas, 2020b;a). These papers however does not relate the expressive power of GNNs to the smoothness and they also work only for discrete spaces. For the two-layer neural networks, the expressive power and smoothness are conjectured to satisfy a quantitative law of robustness (Bubeck et al., 2021).

**Other related works.** Adversarial attacks on graph representations have been well studied, see for instance (Bojchevski & Günnemann, 2019a) which is about edge attack on random walk methods. In general, GNNs are known to be sensitive to adversarial perturbations (Zugner & Günnemann, 2019), and heuristic solutions based on task and application exist, e.g., (Zhu et al., 2019; Jin et al., 2020b). For more information about adversarial attacks/defenses on graphs, see reviews (Jin et al., 2020a;

<sup>2</sup>For example, in social networks, node attributes while carry information about nodes, may be independent of a specific learning task and can be considered as noisy inputs.

Sun et al., 2018; Xu et al., 2020), and for solutions see (Zügner & Günnemann, 2020; Liu et al., 2020; Tang et al., 2020; Zügner & Günnemann, 2019; Dai et al., 2018; Bojchevski & Günnemann, 2019b; Geisler et al., 2020; Chen et al., 2020a).

### 3 BACKGROUND, MOTIVATION, AND PROBLEM FORMULATION

#### 3.1 GRAPH NEURAL NETWORKS (GNNs)

Let  $\bar{\mathbb{G}}_n$  denote the set of all simple graphs on  $n$  vertices, and let  $\mathbb{G}_n \subseteq \bar{\mathbb{G}}_n$  be an arbitrary subset. Assume that for each element of  $\mathbb{G}_n$ , the maximum degree is upper/lower bounded by  $\Delta_{\max}, \Delta_{\min} \in \mathbb{N}$ . Attached to a simple graph  $G = ([n], \mathcal{E}_G)$ , there are initial attributes  $\mathbf{x}_v \in \mathcal{X}$ ,  $v \in [n]$ , for nodes, and non-zero  $\mathbf{e}_{(u,v)} \in \mathcal{X}$ ,  $(u,v) \in \mathcal{E}_G$ , for edges (by setting zero elsewhere it can also be extended to all pairs), where  $\mathcal{X} \subseteq \mathbb{R}^d$  is a compact convex set. Thus, a graph with initial attributes is just a tensor which lies in  $\mathbb{R}^{n \times n \times d}$ . The initial attributes profile is denoted by  $\text{init}(G) = (\{\mathbf{x}_v\}_{v \in [n]}, \{\mathbf{e}_{(u,v)}\}_{(u,v) \in \mathcal{E}})$ . A graph representation is a function  $h(\cdot; \theta) : \mathbb{R}^{n \times n \times d} \rightarrow \mathbb{R}$ , that is further invariant under permutations of nodes, where  $\theta$  denotes the set of all (learnable) parameters, and  $\theta \in \Theta$  for some compact set  $\Theta$ . A Message Passing Graph Neural Network (MPGNN) is a particular way to represent graphs with specific iterative learnable functions, and it is a special case of Graph Neural Networks (GNNs) (Gilmer et al., 2017). In this paper, we consider MPNNs as graph representation and we use the two terms GNN and MPGNN interchangeably.

We consider the most expressive GNN, i.e., Graph Isomorphism Network (GIN) in our formulation (Xu et al., 2019). However, our analysis and results are not restricted to this model and are applicable to other first order GNNs with minor changes. Let us consider the following general formulation. In a GNN with  $k$  iterations, at iteration  $\ell \in [k]$ , we compute

$$\mathbf{h}_v^{(\ell)} = \Phi_\ell \left( \mathbf{h}_v^{(\ell-1)} + \sum_{u \in \mathcal{N}(v)} \Psi_\ell(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)}) \right), \quad (1)$$

for each node  $v \in [n]$ , where  $\mathcal{N}(v)$  denotes the set of neighbors of node  $v$ , and  $\Psi_\ell, \Phi_\ell, \ell \in [k]$ , are learnable functions with appropriate dimensions, e.g, they can be modeled by Multi-Layer Perceptron (MLP), thanks to the universal approximation property (Hornik et al., 1989; Hornik, 1991). In particular, to update, it first aggregates on the neighbors of a node  $v \in [n]$ , and then combine the result with the current representation of  $v$  to achieve the updated representation. Throughout this paper, we assume that  $\Psi_\ell, \Phi_\ell, \ell \in [k]$ , have continuous derivatives. For example, if the activation function of MLPs is twice differentiable, then the assumption is satisfied.

To initialize, let  $\mathbf{h}_v^{(0)} = \mathbf{x}_v$ , for  $v \in [n]$ . Finally, a function is applied, called READOUT function, to achieve the final representation of the graph  $G$ ,

$$h(G) = \Phi \left( \sum_{v \in [n]} \mathbf{h}_v^{(k)} \right), \quad (2)$$

where  $\Phi : \mathbb{R}^{d_k} \rightarrow \mathbb{R}$  is a continuously differentiable real-valued function. To emphasize, we can also denote the final representation of a graph  $G$  with  $h(G, \text{init}(G); \theta)$ , where  $\theta$  denotes all parameters involving the representation. Note that some GNNs use the concatenation function as READOUT, but here we focus on a permutation-invariant function of graphs which allow to analyze the graph classification problem. It is known that this special READOUT function can model any permutation-invariant function as well (Zaheer et al., 2017).

#### 3.2 RANDOM NODE ATTRIBUTES PERTURBATION

In the most simple case, the initial node attributes may be extractable from  $G$  itself (e.g, degree, or from data attached to the graph). They can also be informative for the underlying learning task. The perturbed representation<sup>3</sup>, which is now a random variable (since the perturbation is random), is denoted by  $h_{\text{pert}}(G)$  for any graph with attributes  $G \in \mathbb{R}^{n \times n \times d}$ , and is achieved by  $\mathbf{x}_v \sim \mathcal{N}(\mu_v, \frac{\sigma_n^2}{d} I_d)$

<sup>3</sup>In order to facilitate the analysis we consider the Gaussian perturbation in the paper but the assumption can be easily replaced to the other other perturbations (e.g., integer perturbation).

and  $\mathbf{e}_{(u,v)} \sim \mathcal{N}(\rho_{(u,v)}, \frac{\sigma_e^2}{d} I_d)$  where  $\sigma_n^2, \sigma_e^2$  denote the variance of perturbations, and  $\mu_v, v \in [n]$ , are graph/node dependent noiseless node attributes. Still,  $h(\mathbf{G})$  denotes the representation without randomness, achieved by noiseless initial attributes. We also consider the following bounds for the noiseless initial attributes:  $\|\mu_v\|_2 \leq B_\mu$  for all  $v \in [n]$ , and  $1 \leq \|\rho_{u,v}\|_2 \leq B_\rho$  for all  $(u, v) \in \mathcal{E}_\mathbf{G}$ , and  $\rho_{u,v} = 0$  for other pairs<sup>4</sup>.

### 3.3 GRAPH CLASSIFIERS AND DISCREPANCY

In this paper, a function  $f : \mathbb{R}^{n \times n \times d} \rightarrow \mathbb{R}$  is called a graph function if and only if for any  $\mathbf{G}_1, \mathbf{G}_2 \in \mathbb{R}^{n \times n \times d}$ , we have  $f(\mathbf{G}_1) = f(\mathbf{G}_2)$  whenever  $\mathbf{G}_1 \simeq \mathbf{G}_2$  where  $\simeq$  means the graph isomorphism (with respect to the attached initial attributes). A graph classifier is a graph function whose codomain is  $\{\pm 1\}$ . Let  $\mathbb{P}$  denote an arbitrary but fixed probability measure on  $\mathbb{R}^{n \times n \times d}$  which is supported on the class of graphs with initial attributes<sup>5</sup>. The given dataset of graphs can be considered as i.i.d. samples from  $\mathbb{P}$ . For a graph classifier  $f$ , let  $\eta_f^+$ , (and also  $\eta_f^-$ ) denote the conditional probability measures given  $f(\mathbf{G}) = 1$  (or  $f(\mathbf{G}) = -1$ ), and we refer to them as the measures associated with  $f$ .

For any  $\mathbf{G}, \mathbf{G}' \in \mathbb{R}^{n \times n \times d}$  with initial attributes define their  $\lambda$ -distance as

$$\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda) := \min_{\pi \in \mathcal{S}_n} \left\{ \sum_{v \in [n]} \|\mu_v - \mu'_{\pi(v)}\|_2 + \lambda \sum_{\substack{u,v \in [n] \\ u < v}} \|\rho_{(u,v)} - \rho'_{(\pi(u), \pi(v))}\|_2 \right\}, \quad (3)$$

where  $\mathcal{S}_n$  denotes the set of permutations of  $[n]$ .

For each  $\lambda > 0$  this defines a metric on the space of graphs with initial attributes. In particular, for simple graphs without different initial attributes this reduces to the edit distance on graphs: the minimum number of edges to be added/removed to construct one graph from the other (with respect to the graph isomorphism). Note that  $\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda) = 0$  if and only if  $\mathbf{G}_1 \simeq \mathbf{G}_2$  and always  $\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda) \leq 2B_\mu n + \lambda B_\rho \Delta_{\max} n$  if  $\|\mu_v^i\|_2 \leq B_\mu$  and  $\|\rho_{(u,v)}^i\|_2 \leq B_\rho$  for all nodes and  $i = 1, 2$ .

Discrepancy of a graph classifier is an important quantity in our analysis. Roughly speaking, in the case of simple graphs without initial attributes, it is the average number of edges required to be removed/added to change the output of the classifier, and it is similar to the notion of margin in classical classification problems. First note that we have already defined a metric on the set of graphs with initial attributes, hence we can define the Wasserstein distance for probability measures on this space.

**Definition 1.** For a graph classifier  $f$ , the discrepancy of  $f$  is defined as  $\text{disc}(f) := W_1(\eta_f^+, \eta_f^-)$ , where  $W_1$  denotes the Wasserstein distance, and  $\eta_f^+, \eta_f^-$  are two measures associated with  $f$ . In words, the discrepancy of a graph classifier is the optimal transportation cost of changing the binary label, for two measures induced by random choice of the graphs.

For more technical information about optimal transport in this space, see Appendix A. Also, we compute the discrepancy of well-known graph classifiers in Section 5.

An important problem in deep learning and graph representations learning is to determine the approximation power of a model: the classifiers that can be approximated by some parameters. Typically, we say that a graph classifier can be approximated by a class of graph representations if and only if we can successfully classify the data with a threshold function applied on the representation. Accordingly, under the perturbation of input, the associated probability of error must be small. In what follows, we state this definition rigorously.

**Definition 2.** A graph classifier  $f$  is called to be approximated under perturbations (with the error probability of at most  $\delta$ ) by a class of parameterized GNNs  $h(\cdot; \theta)$ ,  $\theta \in \Theta$ , if and only if

$$\inf_{\theta \in \Theta} \inf_{\gamma \in \mathbb{R}} \mathbb{P} \left( f(\mathbf{G}) \neq \mathbb{1}\{h_{\text{pert}}(\mathbf{G}) \geq \gamma\} \right) \leq \delta, \quad (4)$$

where  $\gamma \in \mathbb{R}$  is a threshold which possibly depends on  $\theta$ , and  $\mathbb{1}\{\cdot\}$  is the sign function.

<sup>4</sup>Clearly, those bounds must hold uniformly on all observable graphs.

<sup>5</sup>We assume that given each graph, the resulting measure is absolutely continuous with respect to the Lebesgue measure, and we assume that each graph classifier is measurable in this space.

## 4 MAIN RESULT

Let  $\mathbb{G}_n^*$  denote the set of observable graphs with initial attributes. Also, let us denote the gradient of the graph representation for  $\mathbf{G} \in \mathbb{G}_n^*$  with respect to initial attributes by  $\partial_v h(\mathbf{G})$ ,  $\partial_{(u,v)} h(\mathbf{G})$  for all  $v \in [n]$ ,  $(u, v) \in \mathcal{E}_{\mathbf{G}}$  (similarly defined for  $h_{\text{pert}}(\mathbf{G})$ ). The main result of this paper is summarized in the following theorem.

**Theorem 1.** *Consider a parametrized class of GNNs  $h(\cdot; \theta)$  which obey the following properties.*

- [polynomial approximation] For any graph with initial attributes  $\mathbf{G} \in \mathbb{G}_n^*$ , there exists a polynomial of initial attributes<sup>6</sup>, of degree at most  $k_{\text{poly}}$ , such that

$$|h_{\text{poly}}(\mathbf{G}) - h(\mathbf{G})| \leq C_{\text{poly}} \sqrt{\text{var}(h_{\text{pert}}(\mathbf{G}))}. \quad (5)$$

This polynomial may depend of  $\mathbf{G}$ , and the RHS, which can be interpreted as the normalized approximation error by polynomials is bounded. The variance of representation is just used to normalize the error in RHS.

- [lower bound on average gradients<sup>7</sup>] For any  $\mathbf{G} \in \mathbb{G}_n^*$ ,  $\|\mathbb{E}[\partial_v h_{\text{pert}}(\mathbf{G})]\|_2 \geq L_{\min}$ ,  $\|\mathbb{E}[\partial_{(u,v)} h_{\text{pert}}(\mathbf{G})]\|_2 \geq L_{\min}$  for all  $v \in [n]$ ,  $(u, v) \in \mathcal{E}_{\mathbf{G}}$ .
- [upper bound on norm of representations] For any  $\mathbf{G} \in \mathbb{G}_n^*$ ,  $v' \in [n]$ , and  $\ell \in [k]$ ,

$$\|\mathbf{h}_{v'}^{(\ell-1)}\|_2 \leq \frac{1}{n} \sum_{v \in [n]} L_h \|\mathbf{x}_v\|_2 + \frac{1}{|\mathcal{E}_{\mathbf{G}}|} \sum_{(u,v) \in \mathcal{E}_{\mathbf{G}}} L_h \|\mathbf{e}_{(u,v)}\|_2 + B_h. \quad (6)$$

This means that when  $\|\mathbf{x}_v\|_2, \|\mathbf{e}_{(u,v)}\|_2 = O(1)$ , then  $\|\mathbf{h}_{v'}^{(\ell-1)}\|_2 = O(2L_h + B_h)$ .

- [upper bound on gradients] For any  $\mathbf{G} \in \mathbb{G}_n^*$ ,  $\|\partial_v h(\mathbf{G})\|_2 \leq L_{\max}$ ,  $\|\partial_{(u,v)} h(\mathbf{G})\|_2 \leq L_{\max}$  for all  $v \in [n]$ ,  $(u, v) \in \mathcal{E}_{\mathbf{G}}$ .
- [smoothness for aggregation] For any  $\mathbf{G} \in \mathbb{G}_n^*$ ,

$$\|\Psi_{\ell}(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)})\|_2 \leq L_{\Psi} \|\mathbf{h}_u^{(\ell-1)}\|_2 + L_{\Psi} \|\mathbf{e}_{(u,v)}\|_2. \quad (7)$$

Then, a graph classifier  $f : \mathbb{G}_n^* \rightarrow \{\pm 1\}$  can be approximated by the class (with the error probability of at most  $\delta$ ), only if

$$\left( \text{disc}(f) + 8.68 \left( 1 + 2n\delta C_{\text{poly}} \left( 1 + 3^{\frac{1}{2}k_{\text{poly}}} \right) \right)^4 (2B_{\mu} + \lambda B_{\rho} \Delta_{\max}) \right) \quad (8)$$

$$\times L_{\max} \left( 1 + \frac{1}{\lambda} \left( 1 + 2kL_{\Psi} (L_h (B_{\mu} + \sigma_n + B_{\rho} + \sigma_e) + B_h) + 2kL_{\Psi} (1 + \sigma_e) \right) \right) \quad (9)$$

$$\geq \frac{0.17 L_{\min} \sigma_n \sqrt{n} \sqrt{\left( 1 + \frac{\sigma_e^2}{2\sigma_n^2} \Delta_{\min} \right)}}{\left( 1 + 2C_{\text{poly}} \left( 1 + 3^{\frac{1}{2}k_{\text{poly}}} \right) \right)^2}. \quad (10)$$

Let us interpret the result and conditions. Theorem 1 is a *quantitative* bound restricting the expressive power of GNNs with the smoothness of the model, variance of perturbations, and the discrepancy of graph classifiers. A more general form of this theorem is also available in appendices, which does not involve lower bounds on the average gradients. We note that a class of GNNs can satisfy the conditions of Theorem 1, if we have lower/upper bounds on the gradient of graph representation, and also we can locally approximate the normalized versions of it with polynomials, which we interpret them as smoothness conditions. Note that the approximation by polynomials essentially follows from having bounds on higher-order derivatives. However, the converse is not true. Also, the lower bound on the average gradients, for small perturbations, is essentially the gradient without perturbations.

<sup>6</sup>For vectors, consider polynomials of entires.

<sup>7</sup>This condition can be relaxed to average of these quantities too, we just considered the most simple case here and the more general form of the theorem is available in appendices.

**Remark 1.** A graph representation is called smooth (on observable graphs) if and only if<sup>8</sup> the bounds  $B_\mu, B_\rho, B_h, L_h, kL_\Psi, C_{poly}, k_{poly} = O(1)$ , and  $\sigma_n = \sigma_e = \Theta(1)$  (bounded variance perturbation), and  $\lambda = \Theta(1)$ . In this case, the condition in the theorem can be written as

$$\left( \text{disc}(f) + O(n\delta\Delta_{\max}) \right) \gtrsim \frac{L_{\min}}{L_{\max}} \sigma \sqrt{n}. \quad (11)$$

We interpret  $L_{\min}/L_{\max}$  as a notion of inverse condition number for smooth GNNs. The above condition means that if a smooth GNN further satisfies the condition  $L_{\min}/L_{\max} = \Omega(1)$  (i.e., having small condition number), and has small probability of error  $\delta \ll (n\Delta_{\max})^{-1}$ , and we can approximate the graph classifier  $f$  with it (under perturbations), then the discrepancy of classifier must be lower bounded by  $O(\sqrt{n})$ . However, we will see in next section that this is quite restricting, and the class of those classifiers is not so large. Note that discrepancy is an intrinsic value assigned to the classifier, and it independent for the representation.

Let us also notice an important fact about the results: the conditions are obtained via a statistical formulation of the problem, hence the computational tractability is not considered here. Indeed, the bounds even hold for arbitrary large networks as long as the smoothness conditions are satisfied.

For better interpretation of the result, we introduce a new notion in the following subsection.

#### 4.1 SEPARABILITY OF (NOISELESS) NODE ATTRIBUTES

As one can expect, the effect of random perturbation to node attributes depends on how noiseless attributes are separated. In other words, if noiseless attributes are well-separated from each other, probably having random perturbation could not change the final representation much, and the model is more expressive (see Definition 2). In order to measure this separation, we define *separability* of the noiseless node attributes as the maximum change when we add an edge to a graph.

**Definition 3.** The separability of noiseless node attributes  $(\mu_v)_{v \in [n]}$  is defined as

$$\text{sep}(\mu) := \max_{\substack{\mathbf{g}, \mathbf{g}' \in \mathbb{G}_n^* \\ \text{dist}(\mathbf{g}, \mathbf{g}'; \lambda) = 1}} \sum_{v \in [n]} \|\mu_v - \mu'_v\|_2. \quad (12)$$

In words, the separability of node attributes means how much adding an edge can change the initial node attributes. It is worth mentioning that the higher separability can make the learning task easier.

**Example 1.** In the absence of noiseless node attributes  $\text{sep}(\mu) = 0$ . This may happen when the node attributes are completely random.

**Example 2.** If the initial node attributes are the degrees of nodes (either one-hot encoded or uncoded) then  $\text{sep}(\mu) = 2$ .

**Example 3.** If the initial node attributes are one-hot encoded, and adding an edge can only change the noiseless attributes of at most  $T$  nodes, then  $\text{sep}(\mu) \leq 2T$ . If  $T \ll n$  then  $\text{sep}(\mu) \ll n$ .

In many reasonable scenarios, we can expect  $\text{sep}(\mu) = O(1)$  or at least  $\text{sep}(\mu) = o(n)$ .

#### 4.2 UNCERTAINTY AND GNNs

Let us assume there is no edge attribute. Note that from the definition of distance function between graphs, we can conclude (if  $\lambda = 1$ ) for a given set of observable graphs and a graph classifier  $f$ ,

$$\text{disc}(f) \leq \widetilde{\text{disc}}(f) \times (\text{sep}(\mu) + O(1)) \quad (13)$$

where  $\widetilde{\text{disc}}(f)$  is the discrepancy of  $f$  with respect to the edit-distance.

<sup>8</sup>We bound  $kL_\Psi$  because the linear approximations of the network accumulate, and if the output is bounded, it makes sense to have this bound too. Indeed,  $k$  in our model can be arbitrary large. Also, one can rescale  $\Psi_\ell$  functions and change the mode a little bit to come up with other cases, and the assumption leads to no loss of generality, other than boundedness of output along with other conditions.

In this case, the result in Theorem 1 can be written as

$$\left(\widetilde{\text{disc}}(f) + O(1)\right) \left(\text{sep}(\mu) + O(1)\right) \gtrsim \frac{L_{\min}}{L_{\max}} \sigma \sqrt{n}. \quad (14)$$

This can be considered as an *uncertainty* principle for smooth GNNs: the product of the discrepancy of the target graph classifier (which measures how simple is that graph classifier to be approximated by a GNN), and the separability of the noiseless node attributes (which measures how noiseless node attributes are capable to help the classification) is lower bounded by  $\sqrt{n}$ . Additionally, if also  $\text{sep}(\mu) = O(1)$ , as we argued before, we have a strong lower bound:  $\widetilde{\text{disc}}(f) \gtrsim \sqrt{n}$ . As we will see later in Section 5, this lower bound completely restricts the class of functions GNN is capable to approximate and consequently if also gives lower bounds on the estimation error of estimating graph parameters. Note that here the new bound is with respect to the edit-distance.

**Remark 2.** *The above bound also shows that if we want to approximate a "complex" graph classifier ( $\widetilde{\text{disc}}(f) = O(1)$ ) then one of the following cases must happen:*

- $\text{sep}(\mu) \gtrsim \sqrt{n}$ . *In this case, the noiseless node attributes can be quite informative.*
- *The model is non-smooth: one of the smoothness conditions are violated.*
- *The model has large probability of error (at least  $O(n^{-1})$ ).*
- *The model is not robust against random perturbation of node attributes: the resulting representation is sensitive.*

### 4.3 MULTI-CLASS CLASSIFICATION

In this part, we extend Theorem 1 to the multi-class graph classifiers. A function  $f : \mathbb{G}_n^* \rightarrow [p]$  is called a multi-class graph classifier, where  $p \in \mathbb{N}$ . The case  $p = 2$  reduces to the aforementioned class of (binary) graph classifiers.

First we define approximation of multi-class graph classifiers with GNNs.

**Definition 4.** *A multi-class graph classifier  $f : \mathbb{G}_n^* \rightarrow [p]$  can be approximated (with error probability of at most  $\delta$ ) by a class of parameterized GNNs  $h(\cdot; \theta)$ ,  $\theta \in \Theta$ , if and only if*

$$\inf_{\theta \in \Theta} \inf_{\zeta \in \mathcal{Z}} \sup_{\mathbb{P} \in \{\mathbb{P}_{p'} \otimes \mathcal{N} : p' \in [p]\}} \mathbb{P} \left( f(\mathbb{G}) \neq \zeta(h_{\text{pert}}(\mathbb{G})) \right) \leq \delta, \quad (15)$$

where  $\mathcal{Z}$  is the set of piecewise constant functions  $\zeta : \mathbb{R} \rightarrow [p]$  with at most  $p - 1$  discontinuity points. Also,  $\mathbb{P}_{p'}$  is the conditional probability measure on  $f^{-1}(p')$ , for each  $p' \in [p]$ , and  $\mathcal{N}$  denotes the Gaussian probability measure corresponded to the perturbations (see 3.2 for more explanations).

**Remark 3.** *The motivation of this definition is that if  $f(\mathbb{G}) \approx h_{\text{pert}}(\mathbb{G})$  then there is confidence intervals allowing to infer  $f(\mathbb{G})$  from  $h_{\text{pert}}(\mathbb{G})$  with high probability.*

Let us define the discrepancy of multi-class graph classifiers as follows.

**Definition 5.** *For a multi-class graph classifier  $f : \mathbb{G}_n^* \rightarrow [p]$ , the discrepancy of  $f$  is defined as*

$$\text{disc}(f) := \min_{\substack{p_1, p_2 \in [p] \\ p_1 \neq p_2}} \text{disc}(f_{p_1, p_2}), \quad (16)$$

where  $f_{p_1, p_2}(\mathbb{G}) : f^{-1}(p_1) \cup f^{-1}(p_2) \rightarrow \{\pm 1\}$  is defined as

$$f_{p_1, p_2}(\mathbb{G}) := \begin{cases} +1 & \text{iff } f(\mathbb{G}) = p_1 \\ -1 & \text{iff } f(\mathbb{G}) = p_2 \end{cases} \quad (17)$$

To define the discrepancy of multi-class graph classifiers, we consider all pairs of differently labeled graphs according to  $p_1, p_2$ , and take the minimum discrepancy of the corresponded graph classifier. Note that the set of observable graphs here is  $f^{-1}(p_1) \cup f^{-1}(p_2)$  for fixed  $p_1 \neq p_2$ .

The main result of this part is the following proposition which directly follows from Theorem 1.

**Proposition 1.** *Theorem 1 holds for multi-class graph classifiers  $f : \mathbb{G}_n^* \rightarrow [p]$ .*

Note that to satisfy the smoothness conditions for multi-class graph classifiers, we can normalize the resulting representations  $f_{p_1, p_2}$ , since Theorem 1 holds for each  $p_1, p_2$ .

Table 1: Approximation bounds for graph parameters

| $f : \mathbb{G}_n^* \rightarrow [0, T]$ |                                      |  |                 |
|---|--------------------------------------|--|-----------------|
| Parameter                               | $\text{disc}(f_p)$                   | Discrepancy & additive error   | $\kappa$        |
| Number of edges                         | $\frac{n\Delta_{\max}}{2^p} + O(1)$  | $\text{disc}(f_p) = \epsilon_{\text{additive}} + O(1)$               | 1               |
| Number of triangles                     | $\frac{n}{p} + O(1)$                 | $\text{disc}(f_p) = 3\epsilon_{\text{additive}} + O(1)$              | 3               |
| Number of subgraphs $H$                 | $\frac{p}{n m_H} + O(1)$             | $\text{disc}(f_p) = m_H \epsilon_{\text{additive}} + O(1)$           | $m_H$           |
| Diameter                                | $\frac{p}{n-1} + O(1)$               | $\text{disc}(f_p) = \epsilon_{\text{additive}} + O(1)$               | 1               |
| Max-cut                                 | $\frac{n\Delta_{\max}}{2^p} + O(1)$  | $\text{disc}(f_p) = \epsilon_{\text{additive}} + O(1)$               | 1               |
| Min-cut                                 | $\frac{n\Delta_{\max}}{2^p} + O(1)$  | $\text{disc}(f_p) = \epsilon_{\text{additive}} + O(1)$               | 1               |
| Clique number                           | $\frac{\Delta_{\max}^2}{2^p} + O(1)$ | $\text{disc}(f_p) = \Delta_{\max} \epsilon_{\text{additive}} + O(1)$ | $\Delta_{\max}$ |
| Independence number                     | $\frac{\Delta_{\max}}{p} + O(1)$     | $\text{disc}(f_p) = \Delta_{\max} \epsilon_{\text{additive}} + O(1)$ | $\Delta_{\max}$ |
| Number of connected components          | $\frac{n}{p} + O(1)$                 | $\text{disc}(f_p) = \epsilon_{\text{additive}} + O(1)$               | 1               |

## 5 APPLICATIONS OF THE DERIVED BOUND

In this part, we give some applications of the proved bound by computing the discrepancy of a number of well-known graph functions. First consider a graph property as  $f : \mathbb{G}_n^* \rightarrow [0, T]$ , e.g., the number of edges, diameter, etc. To approximate it with a GNN, let's consider the discrepancy of its quantization to a  $p$ -level function:

$$f_p(\mathbb{G}) := p' \quad \text{iff} \quad \frac{T(p' - 1)}{p} \leq f(\mathbb{G}) < \frac{Tp'}{p}, \quad (18)$$

for any  $p' \in [p]$ . If  $f$  as a multi-class graph classifier can be approximated by GNNs, then  $f_p$  can be approximated by GNNs, and the estimation error for  $f$  is at most  $\epsilon_{\text{additive}} = \frac{T}{p}$ . More precisely, this quantity is the expectation of the absolute value of the error  $f - f_{\text{GNN}}$ .

Hence, to achieve a lower bound for representing a graph parameter (function) from the main result, we need to compute the discrepancy of  $f_p$ . We expect that  $\text{disc}(f_p)$  is decreasing with  $p$ , meaning that having more resolution applies more complexity for representation. We will make this observation more clear in the following proposition and achieve bounds on the estimation error  $\epsilon_{\text{additive}}$ .

**Proposition 2.** *For any graph parameter in Table 1, there is an appropriate set of graphs and  $T \in \mathbb{R}$  such that we have the provided bounds on the table on the discrepancy of their quantization. Also, the additive error of approximation is related to the discrepancy in the third column.*

Note that the number of subgraphs can be interpreted as induced or not necessarily induced, and the bounds are achieved for both cases. Also,  $n_H$  and  $m_H$  denote the number of nodes and the number of edges of a small graph  $H$ , and typically for counting tasks we have  $n_H, m_H = O(1)$ .

Now using the above table and the result on the approximation of the multi-class graph classifiers, we can derive the following bound on the estimation error of the above graph parameters.

**Remark 4.** *The proposition gives a quantitative bound based on the variance of perturbation, smoothness of the model, separability of noiseless attributes, and the additive estimation error of representation. In particular, under the assumptions in Section 4.2, we have the following necessary condition:*

$$\left(\kappa \times \epsilon_{\text{additive}} + O(1)\right) \left(\text{sep}(\mu) + O(1)\right) \gtrsim \sqrt{n}. \quad (19)$$

We can immediately conclude the following proposition.

**Proposition 3.** *Consider the case that  $\text{sep}(\mu) = O(1)$ . Then, for counting subgraph  $H$  (either induced or not necessarily induced) with GNNs, we have*

$$\epsilon_{\text{additive}} \gtrsim \sqrt{n}. \quad (20)$$

*This means that roughly speaking smooth GNNs with node attributes perturbation are unable to count subgraphs with the additive estimation error  $o(\sqrt{n})$ .*

The lower bounds in previous propositions suggest that smooth GNNs with node attributes perturbations are subject to high approximation error for the graph parameters in Table 1, and that is why we can argue that those graph representations are highly restricted to large error.

## 6 CONCLUSION

In this paper, we theoretically analyze GNNs with smooth representations, along with perturbed node attributes. It is observed that having smoothness highly restricts the model, and imposes several lower bounds on the approximation error for graph parameters. Our bound also allows to achieve compute a trade-off relating the expressive power of GNNs to the smoothness, separability of noiseless attributes, and the variance of perturbation.

## REFERENCES

- Waiss Azizian and marc lelarge. Expressive power of invariant and equivariant graph neural networks. In *International Conference on Learning Representations*, 2021.
- Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. In *Advances in neural information processing systems*, pp. 4502–4510, 2016.
- Aleksandar Bojchevski and Stephan Günnemann. Adversarial attacks on node embeddings via graph poisoning. In *International Conference on Machine Learning*, pp. 695–704. PMLR, 2019a.
- Aleksandar Bojchevski and Stephan Günnemann. Certifiable robustness to graph perturbations. *arXiv preprint arXiv:1910.14356*, 2019b.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Sébastien Bubeck, Yuanzhi Li, and Dheeraj M Nagaraj. A law of robustness for two-layers neural networks. In *Conference on Learning Theory*, pp. 804–820. PMLR, 2021.
- Xu Chen, Ya Zhang, Ivor Tsang, and Yuangang Pan. Learning robust node representation on graphs. *arXiv preprint arXiv:2008.11416*, 2020a.
- Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. In *Advances in Neural Information Processing Systems*, pp. 15894–15902, 2019.
- Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count substructures? *arXiv preprint arXiv:2002.04025*, 2020b.
- Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack on graph structured data. In *International conference on machine learning*, pp. 1115–1124. PMLR, 2018.
- George Dasoulas, Ludovic Dos Santos, Kevin Scaman, and Aladin Virmaux. Coloring graph neural networks for node disambiguation. In Christian Bessiere (ed.), *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pp. 2126–2132. International Joint Conferences on Artificial Intelligence Organization, 7 2020.
- David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems*, pp. 2224–2232, 2015.
- Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits of graph neural networks. In *Int. Conference on Machine Learning (ICML)*, pp. 5204–5215. 2020.
- Simon Geisler, Daniel Zügner, and Stephan Günnemann. Reliable graph neural networks via robust aggregation. *Advances in Neural Information Processing Systems*, 33, 2020.

- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1263–1272, 2017.
- Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 1991.
- Kurt Hornik, Maxwell Stinchcombe, Halbert White, et al. Multilayer feedforward networks are universal approximators. *Neural networks*, 1989.
- John Ingraham, Vikas Kamur Garg, Regina Barzilay, and Tommi S Jaakkola. Generative models for graph-based protein design. 2021.
- Wei Jin, Yaxin Li, Han Xu, Yiqi Wang, and Jiliang Tang. Adversarial attacks and defenses on graphs: A review and empirical study. *arXiv preprint arXiv:2003.00653*, 2020a.
- Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph structure learning for robust graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 66–74, 2020b.
- Wengong Jin, Kevin Yang, Regina Barzilay, and Tommi Jaakkola. Learning multimodal graph-to-graph translation for molecular optimization. *arXiv preprint arXiv:1812.01070*, 2018.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- FD Lesley and VI Rotar. Some remarks on lower bounds of chebyshev’s type for half-lines. *JIPAM*, 4(5):96, 2003.
- Jaechang Lim, Seongok Ryu, Kyubyong Park, Yo Joong Choe, Jiyeon Ham, and Woo Youn Kim. Predicting drug–target interaction using a novel graph neural network with 3d structure-embedded graph representation. *Journal of chemical information and modeling*, 59(9):3981–3988, 2019.
- Yang Liu, Xianzhuo Xia, Liang Chen, Xiangnan He, Carl Yang, and Zibin Zheng. Certifiable robustness to discrete adversarial perturbations for factorization machines. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 419–428, 2020.
- Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International Conference on Learning Representations*, 2020a.
- Andreas Loukas. How hard is to distinguish graphs with graph neural networks? In *Advances in Neural Information Processing Systems*, 2020b.
- Ryan O’Donnell. *Analysis of boolean functions*. Cambridge University Press, 2014.
- Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node classification. *arXiv preprint arXiv:1905.10947*, 2019.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*, 2020.
- Ryoma Sato. A survey on the expressive power of graph neural networks. *arXiv preprint arXiv:2003.04078*, 2020.
- Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Random features strengthen graph neural networks. In *Proceedings of the 2021 SIAM International Conference on Data Mining (SDM)*, pp. 333–341. SIAM, 2021.
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. Computational capabilities of graph neural networks. *IEEE Transactions on Neural Networks*, 20(1):81–102, 2009.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.

- Jonathan Shlomi, Peter Battaglia, and Jean-Roch Vlimant. Graph neural networks in particle physics. *Machine Learning: Science and Technology*, 2(2):021001, 2020.
- Alexey Strokach, David Becerra, Carles Corbi-Verge, Albert Perez-Riba, and Philip M Kim. Fast and flexible protein design using deep graph neural networks. *Cell Systems*, 11(4):402–411, 2020.
- Lichao Sun, Yingdong Dou, Carl Yang, Ji Wang, Philip S Yu, Lifang He, and Bo Li. Adversarial attack and defense on graph data: A survey. *arXiv preprint arXiv:1812.10528*, 2018.
- Mengying Sun, Sendong Zhao, Coryandar Gilvary, Olivier Elemento, Jiayu Zhou, and Fei Wang. Graph convolutional networks for computational drug development and discovery. *Briefings in bioinformatics*, 21(3):919–935, 2020.
- Xianfeng Tang, Yandong Li, Yiwei Sun, Huaxiu Yao, Prasenjit Mitra, and Suhang Wang. Transferring robustness for graph neural network against poisoning attacks. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pp. 600–608, 2020.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- Saurabh Verma and Zhi-Li Zhang. Stability and generalization of graph convolutional neural networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1539–1548, 2019.
- Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
- Edward Wagstaff, Fabian Fuchs, Martin Engelcke, Ingmar Posner, and Michael A Osborne. On the limitations of representing functions on sets. In *International Conference on Machine Learning*, pp. 6487–6494, 2019.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 2020.
- Han Xu, Yao Ma, Hao-Chen Liu, Debayan Deb, Hui Liu, Ji-Liang Tang, and Anil K Jain. Adversarial attacks and defenses in images, graphs and text: A review. *International Journal of Automation and Computing*, 17(2):151–178, 2020.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. In *International Conference on Machine Learning*, pp. 7134–7143. PMLR, 2019.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.
- Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.
- Dingyuan Zhu, Ziwei Zhang, Peng Cui, and Wenwu Zhu. Robust graph convolutional networks against adversarial attacks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1399–1407, 2019.
- Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 246–256, 2019.
- Daniel Zügner and Stephan Günnemann. Certifiable robustness of graph convolutional networks under structure perturbations. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1656–1665, 2020.
- Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta learning. In *International Conference on Learning Representations*, 2019.

## A OPTIMAL TRANSPORTATION DUALITY

Let  $\mathbb{G}_n^* \subseteq \mathbb{R}^{n \times n \times d}$  denote the set of all graphs with initial attributes which is equipped with the  $\lambda$ -distance. Let  $\mathcal{P}(\mathbb{G}_n^*)$  denote the set of all probability measures on the set of graphs with initial attributes that are absolutely continuous with respect to the Lebesgue measure. The Wasserstein distance between two probability measures  $\mathbb{P}_1, \mathbb{P}_2 \in \mathcal{P}(\mathbb{G}_n^*)$  is defined as

$$W_1(\mathbb{P}_1, \mathbb{P}_2) := \inf_{\mathbb{P} \in \mathcal{M}(\mathbb{P}_1, \mathbb{P}_2)} \mathbb{E}_{(\mathbf{G}_1, \mathbf{G}_2) \sim \mathbb{P}}[\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda)], \quad (21)$$

where  $\mathcal{M}(\mathbb{P}_1, \mathbb{P}_2)$  denotes the set of all couplings of  $\mathbb{P}_1, \mathbb{P}_2$ . Note that  $(\mathcal{P}(\mathbb{G}_n^*), W_1)$  is a metric space. From the duality principle, we have the following alternative definition for the Wasserstein distance,

$$W_1(\mathbb{P}_1, \mathbb{P}_2) = \sup_{f \in \mathcal{F}_1} \left\{ \mathbb{E}_{\mathbf{G} \sim \mathbb{P}_1} [f(\mathbf{G})] - \mathbb{E}_{\mathbf{G} \sim \mathbb{P}_2} [f(\mathbf{G})] \right\}, \quad (22)$$

where

$$\mathcal{F}_1 := \left\{ f_1 : \mathbb{G}_n^* \rightarrow \mathbb{R} : \sup_{\mathbf{G} \neq \mathbf{G}'} \frac{|f(\mathbf{G}) - f(\mathbf{G}')|}{\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)} \leq 1 \right\}. \quad (23)$$

For more on optimal transport, see (Villani, 2003).

## B PROOF OF THEOREM 1

Since the proof of the theorem is long, let us decompose it into a few subsections.

### B.1 A MORE GENERAL RESULT

Let us restate the main theorem (in a more general case) here.

**Theorem 2.** *Under the conditions of Theorem 1, except the lower bound on average gradients, a graph classifier  $f : \mathbb{G}_n^* \rightarrow \{\pm 1\}$  can be approximated by the class (with the error probability of at most  $\delta$ ), only if*

$$(\text{disc}(f) + 8.68 \left(1 + 2n\delta C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 (2B_\mu + \lambda B_\rho \Delta_{\text{max}})) \quad (24)$$

$$\times L_{\text{max}} \left(1 + \frac{1}{\lambda} (1 + 2kL_\Psi L_h (B_\mu + \sigma_n + B_\rho + \sigma_e) + 2kL_\Psi (1 + \sigma_e))\right) \quad (25)$$

$$\geq \frac{0.17}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2} \max \left\{ \inf_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}, \quad (26)$$

$$\mathbb{E}[\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] - 8.68\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \right\}, \quad (27)$$

**Remark 5.** *This result is more general than that of Theorem 1. In particular, the RHS involves the average variance, if  $\delta$  is small enough, and this means that we are outperforming Theorem 1, which can be concluded from the worst-case infimum variance (details are below).*

We conclude the main result with the following proposition.

**Proposition 4.** *We have that*

$$\inf_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \geq L_{\text{min}} \sigma_n \sqrt{n} \sqrt{\left(1 + \frac{\sigma_e^2}{2\sigma_n^2} \Delta_{\text{min}}\right)}. \quad (28)$$

*Proof.* For any fixed  $\mathbf{G} \in \mathbb{G}_n^*$ ,  $h_{\text{pert}}(\mathbf{G})$  is a deterministic function of a number of Gaussian random variables:  $\mathbf{x}_v = \mu_v + \mathbf{n}_v$  and  $\mathbf{e}_{(u,v)} = \rho_{(u,v)} + \mathbf{n}_{(u,v)}$  such that  $\mathbf{n}_v \sim \mathcal{N}(0, \frac{\sigma_n^2}{d} I_d)$  and  $\mathbf{n}_{(u,v)} \sim$

$\mathcal{N}(0, \frac{\sigma_e^2}{d} I_d)$  for  $v \in [n]$  and  $(u, v) \in \mathcal{E}_G$ . From the theory of Hermit polynomials and Gaussian Hilbert spaces (O'Donnell, 2014), we know that

$$\text{var}(h_{\text{pert}}(\tilde{G})) \geq \sigma_n^2 \sum_{v \in [n]} \|\mathbb{E}[\partial_v h_{\text{pert}}(G)]\|_2^2 + \sigma_e^2 \sum_{(u,v) \in \mathcal{E}_G} \|\mathbb{E}[\partial_{(u,v)} h_{\text{pert}}(G)]\|_2^2. \quad (29)$$

The RHS can be interpreted as the square norm of the average gradient around perturbations. Intuitively, if perturbation is small, this is essentially the gradient without perturbations. From the assumption, we simply get

$$\text{var}(h_{\text{pert}}(\tilde{G})) \geq n\sigma_n^2 L_{\min}^2 + |\mathcal{E}_G| \sigma_e^2 L_{\min}^2 \geq nL_{\min}^2 (\sigma_n^2 + \frac{1}{2} \sigma_e^2 \Delta_{\min}). \quad (30)$$

We are done.  $\square$

## B.2 APPROXIMATING THE GRAPH CLASSIFIER

In this part, we approximate the given graph classifier  $f$  with an auxiliary classifier  $\hat{f}$  which agrees with  $f$  with high probability, we then exploit properties of the new classifier to proceed the proof. Note that if the probability of error in Equation 4 is bounded for any graph in  $\mathbb{G}_n^*$ , instead of random graphs, then we can skip this part since the approximation is not required and all graphs satisfy the conditions below.

Consider is a graph classifier  $f$  which can be approximated by GNNs, meaning that there exist  $\theta, \gamma$  such that

$$\mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\}) \leq \delta. \quad (31)$$

Now consider the following set:

$$\mathcal{A} := \left\{ \tilde{G} \in \mathbb{G}_n^* : \mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\} | G = \tilde{G}) \leq \frac{1}{C_1} \right\}, \quad (32)$$

where  $C_1$  is an arbitrary fixed constant to be determined later. For simplicity, let us define

$$\mathbb{P}_{\tilde{G}}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\}) := \mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\} | G = \tilde{G}). \quad (33)$$

Note that we have

$$\mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\}) = \mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\} | \mathcal{A}) \mathbb{P}(\mathcal{A}) \quad (34)$$

$$+ \mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\} | \mathcal{A}^c) \mathbb{P}(\mathcal{A}^c) \quad (35)$$

$$\geq \mathbb{P}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\} | \mathcal{A}^c) (1 - \mathbb{P}(\mathcal{A})) \quad (36)$$

$$\geq \inf_{\tilde{G} \in \mathcal{A}^c} \mathbb{P}_{\tilde{G}}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\}) (1 - \mathbb{P}(\mathcal{A})) \quad (37)$$

$$\geq \frac{1}{C_1} (1 - \mathbb{P}(\mathcal{A})). \quad (38)$$

Thus,  $\mathbb{P}(\mathcal{A}) \geq 1 - C_1 \delta$ .

Now, consider the function  $\hat{f} : \mathcal{A} \subseteq \mathbb{G}_n^* \rightarrow \{\pm 1\}$  which is the restriction of  $f$  to the set  $\mathcal{A}$ . We write

$$\text{disc}(\hat{f}) = W_1(\eta_{\hat{f}}^+, \eta_{\hat{f}}^-) \quad (39)$$

$$\stackrel{(a)}{\leq} W_1(\eta_{\hat{f}}^+, \eta_f^+) + W_1(\eta_{\hat{f}}^+, \eta_f^-) + W_1(\eta_{\hat{f}}^-, \eta_f^-) \quad (40)$$

$$= \text{disc}(f) + W_1(\eta_{\hat{f}}^+, \eta_f^+) + W_1(\eta_{\hat{f}}^-, \eta_f^-), \quad (41)$$

where (a) holds by the triangle inequality. Recall that  $\mathbb{G}_n^*$  is the set of graphs with initial attributes. To bound the other terms, consider the following coupling between the measures  $\eta_{\hat{f}}^\kappa, \eta_f^\kappa$  for  $\kappa \in \{\pm\}$ . Let  $\eta^\kappa$  denote the following probability measure on  $\mathbb{G}_n^* \times \mathbb{G}_n^*$ : first we choose an element  $G_1$  with respect to the probability measure  $\eta_f^\kappa$ , and then if it belongs to  $\mathcal{A}$  then the second element  $G_2$  is

the same as  $G_1$ . Otherwise, we draw a random element of  $\mathcal{A}$  for  $G_2$ , with respect to the probability measure  $\eta_{\hat{f}}^{\kappa}$ . This is clearly a coupling between  $\eta_{\hat{f}}^{\kappa}, \eta_f^{\kappa}$ , and so we conclude

$$W_1(\eta_{\hat{f}}^{\kappa}, \eta_f^{\kappa}) \leq \mathbb{E}_{\eta^{\kappa}}[\text{dist}(G_1, G_2; \lambda)] \quad (42)$$

$$\leq (2B_{\mu}n + \lambda B_{\rho}\Delta_{\max}n)(1 - \mathbb{P}(\mathcal{A})) \quad (43)$$

$$\leq C_1\delta n(2B_{\mu} + \lambda B_{\rho}\Delta_{\max}), \quad (44)$$

since the trivial bound  $\text{dist}(G_1, G_2; \lambda) \leq 2B_{\mu}n + \lambda B_{\rho}\Delta_{\max}n$  holds for all graphs.

Combining the above result with (41) shows that

$$\text{disc}(\hat{f}) \leq \text{disc}(f) + 2C_1\delta n(2B_{\mu} + \lambda B_{\rho}\Delta_{\max}). \quad (45)$$

### B.3 DUALITY PRINCIPLE AND APPROXIMATING GNNs

In this part, we use the optimal transportation duality, and then we approximate the GNN representation with a polynomial, along with we using mean-value theorem to bound the approximation error. This approximation allows us to compute useful bounds for the given graph representation.

Now let us consider associated measures of  $\hat{f}$  and note that  $\mathbb{E}[h_{\text{pert}}(\tilde{G})] : \mathbb{G}_n^* \rightarrow \mathbb{R}$  is a function on the space  $(\mathbb{G}_n^*, \text{dist}(\cdot, \cdot; \lambda))$ , thus

$$\text{disc}(\hat{f}) = W_1(\eta_{\hat{f}}^+, \eta_{\hat{f}}^-) \quad (46)$$

$$= \sup_{f_1 \in \mathcal{F}_1} \left\{ \mathbb{E}_{G \sim \eta_{\hat{f}}^+}[f_1(G)] - \mathbb{E}_{G \sim \eta_{\hat{f}}^-}[f_1(G)] \right\} \quad (47)$$

$$\geq \frac{1}{\mathcal{L}(h)} \left( \mathbb{E}_{G \sim \eta_{\hat{f}}^+}[\mathbb{E}[h_{\text{pert}}(G)]] - \mathbb{E}_{G \sim \eta_{\hat{f}}^-}[\mathbb{E}[h_{\text{pert}}(G)]] \right), \quad (48)$$

where

$$\mathcal{L}(h) := \sup_{\substack{G, G' \in \mathbb{G}_n^* \\ G \neq G'}} \frac{|\mathbb{E}[h_{\text{pert}}(G)] - \mathbb{E}[h_{\text{pert}}(G')]|}{\text{dist}(G, G'; \lambda)}. \quad (49)$$

Note that the function  $\hat{f}$  is supported on  $\mathcal{A}$ , and thus we have

$$\mathbb{P}_{\tilde{G}}(f(G) \neq \mathbb{1}\{h_{\text{pert}}(G) \geq \gamma\}) \leq \frac{1}{C_1}, \quad (50)$$

for all  $\tilde{G} \in \mathcal{A}$  (see (33) for the definition). Fix  $G \in \mathcal{A}$ , and let us approximate the above probability from below. First let us consider the case that  $\hat{f}(G) = -1$ . Using the definition, we have

$$\mathbb{P}_{\tilde{G}}(h_{\text{pert}}(\tilde{G}) \geq \gamma) \leq \frac{1}{C_1}. \quad (51)$$

But  $h_{\text{pert}}(\tilde{G})$  is a deterministic function of some jointly Gaussian random variables. In particular, let us recall that  $h_{\text{pert}}(\tilde{G})$  is the graph representation for  $\tilde{G}$ , when we have perturbed initial attributes  $\mathbf{x}_v \sim \mathcal{N}(\mu_v, \frac{\sigma_v^2}{d} I_d)$  and  $\mathbf{e}_{(u,v)} \sim \mathcal{N}(\rho_{(u,v)}, \frac{\sigma_{(u,v)}^2}{d} I_d)$  where  $v \in [n]$  and  $(u, v) \in \mathcal{E}_{\tilde{G}}$ . Now we need the following lemma for the rest of analysis.

**Lemma 1** (Lesley & Rotar (2003)). *Let  $X$  be a random variable with zero mean, unit variance, and  $\mathbb{E}[X^4] = \zeta \geq 1$ . Then, we have*

$$\mathbb{P}(X \geq t) \geq \frac{0.46}{\zeta} - \frac{1.4}{\sqrt{\zeta}}t + \frac{0.02}{\zeta^{1.5}}t, \quad (52)$$

for any  $t \geq 0$ .

The above lemma can be proved by carefully finding an upper bound for the indicator function  $\mathbb{1}\{\cdot \leq t\}$  in terms of a degree four polynomial, and do some calculation to optimize the bound. Consequently, under the conditions of the lemma,

$$\forall t \geq 0 : \mathbb{P}(X \geq t) \geq \frac{0.46}{\zeta} - \frac{1.4}{\sqrt{\zeta}}t \implies \mathbb{P}(X \geq \frac{0.17}{\sqrt{\zeta}}) \geq \frac{0.23}{\zeta}. \quad (53)$$

Let  $X_{\tilde{G}} := \frac{h_{\text{pert}}(\tilde{G}) - \mathbb{E}[h_{\text{pert}}(\tilde{G})]}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}}$ . Clearly,  $X_{\tilde{G}}$  is a random variable with zero mean and unit variance, and that randomness comes merely from perturbations, since  $\tilde{G}$  is fixed. Thus, by the above lemma, we conclude that

$$\mathbb{P}_{\tilde{G}}(h_{\text{pert}}(\tilde{G}) \geq \gamma) = \mathbb{P}_{\tilde{G}}\left(X_{\tilde{G}} \geq \frac{\gamma - \mathbb{E}[h_{\text{pert}}(\tilde{G})]}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}}\right) \leq \frac{1}{C_1}. \quad (54)$$

Now using (53), we conclude that

$$C_1 \geq \frac{\mathbb{E}[X_{\tilde{G}}^4]}{0.23} \implies \frac{\gamma - \mathbb{E}[h_{\text{pert}}(\tilde{G})]}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \geq \frac{0.17}{\sqrt{\mathbb{E}[X_{\tilde{G}}^4]}}. \quad (55)$$

Now we need to upper bound the fourth moment of  $X_{\tilde{G}}$ . To do that, we use the following result on Gaussian hypercontractivity.

**Lemma 2** (Boucheron et al. (2013)). *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a multivariate polynomial of degree at most  $k$ , and  $(X_1, X_2, \dots, X_n) \sim \mathcal{N}(0, I)$  be a sequence of i.i.d. Gaussian random variables. Then, we have*

$$(\mathbb{E}[|f(X_1, X_2, \dots, X_n)|^q])^{1/q} \leq (q-1)^{k/2} (\mathbb{E}[|f(X_1, X_2, \dots, X_n)|^2])^{1/2}. \quad (56)$$

Note that by assumption there is a polynomial  $h_{\text{poly}}(\tilde{G})$  of degree at most  $k_{\text{poly}}$  such that

$$|h_{\text{poly}}(\tilde{G}) - h(\tilde{G})| \leq C'_{\text{poly}} := C_{\text{poly}} \sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}, \quad (57)$$

holds for any initial attributes, even if we consider the perturbed case  $h_{\text{pert}}(\tilde{G})$ . The polynomial, however, may depend on  $G$ . Thus, we obtain that

In particular, we can apply the above lemma to  $X_{\tilde{G}}$  since it is a real-valued random variable, which is a deterministic function of Gaussian random variables (i.e., perturbations). Thus, setting  $q = 4$  in the lemma results in

$$(\mathbb{E}[X_{\tilde{G}}^4])^{1/4} = \left(\mathbb{E}\left[\left(\frac{h_{\text{pert}}(\tilde{G}) - \mathbb{E}[h_{\text{pert}}(\tilde{G})]}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}}\right)^4\right]\right)^{1/4} \quad (58)$$

$$= \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(\mathbb{E}[(h_{\text{pert}}(\tilde{G}) - \mathbb{E}[h_{\text{pert}}(\tilde{G})])^4]\right)^{1/4} \quad (59)$$

$$= \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(\mathbb{E}[(h_{\text{pert}}(\tilde{G}) - h_{\text{poly}}(\tilde{G})) + (h_{\text{poly}}(\tilde{G}) - \mathbb{E}[h_{\text{poly}}(\tilde{G})])\right. \quad (60)$$

$$\left. + (\mathbb{E}[h_{\text{poly}}(\tilde{G})] - \mathbb{E}[h_{\text{pert}}(\tilde{G})])^4]\right)^{1/4} \quad (61)$$

$$\stackrel{(a)}{\leq} \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(2C'_{\text{poly}} + (\mathbb{E}[(h_{\text{poly}}(\tilde{G}) - \mathbb{E}[h_{\text{poly}}(\tilde{G})])^4])^{1/4}\right) \quad (62)$$

$$\stackrel{(b)}{\leq} \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(2C'_{\text{poly}} + 3^{\frac{1}{2}k_{\text{poly}}} (\mathbb{E}[(h_{\text{poly}}(\tilde{G}) - \mathbb{E}[h_{\text{poly}}(\tilde{G})])^2])^{1/2}\right) \quad (63)$$

$$= \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(2C'_{\text{poly}} + 3^{\frac{1}{2}k_{\text{poly}}} (\mathbb{E}[(h_{\text{poly}}(\tilde{G}) - h_{\text{pert}}(\tilde{G}))\right. \quad (64)$$

$$\left. + (h_{\text{pert}}(\tilde{G}) - \mathbb{E}[h_{\text{pert}}(\tilde{G})]) + (\mathbb{E}[h_{\text{pert}}(\tilde{G})] - \mathbb{E}[h_{\text{poly}}(\tilde{G})])^2]\right)^{1/2} \quad (65)$$

$$\stackrel{(c)}{\leq} \frac{1}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{G}))}} \left(2C'_{\text{poly}} + 3^{\frac{1}{2}k_{\text{poly}}} (2C'_{\text{poly}} + (\mathbb{E}[(h_{\text{pert}}(\tilde{G}) - \mathbb{E}[h_{\text{pert}}(\tilde{G})])^2])^{1/2})\right) \quad (66)$$

$$= 1 + \frac{2C'_{\text{poly}}}{\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}} \left(1 + 3^{\frac{1}{2}k_{\text{poly}}}\right) \quad (67)$$

$$= 1 + 2C_{\text{poly}} \left(1 + 3^{\frac{1}{2}k_{\text{poly}}}\right), \quad (68)$$

where (a) and (c) hold by Minkowski inequality, and (b) hold by Lemma 2. To proceed, we combine the above lemma with (55) and obtain the optimal constant  $C_1$ , thus if

$$C_1 = 4.34 \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4, \quad (69)$$

then

$$\gamma \geq \mathbb{E}[h_{\text{pert}}(\tilde{\mathbf{G}})] + 0.17 \frac{\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2}, \quad (70)$$

holds for any  $\tilde{\mathbf{G}} \in \mathcal{A}$  such that  $\hat{f}(\mathbf{G}) = -1$ . Thus, taking the expectation from both sides with respect to  $\mathbf{G} \sim \eta_{\tilde{f}}^-$  results in

$$\gamma \geq \mathbb{E}_{\mathbf{G} \sim \eta_{\tilde{f}}^-} [\mathbb{E}[h_{\text{pert}}(\tilde{\mathbf{G}})]] + 0.17 \frac{\mathbb{E}_{\eta_{\tilde{f}}^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}]}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2}, \quad (71)$$

while again we emphasize that  $h_{\text{pert}}(\tilde{\mathbf{G}})$  is a random variable which depends on on perturbations, given any graph  $\mathbf{G} \in \mathcal{A}$ .

In a similar way, for any  $\mathbf{G} \in \mathcal{A}$  such that  $f(\mathbf{G}) = 1$  it can be shown that

$$\gamma \leq \mathbb{E}_{\mathbf{G} \sim \eta_{\tilde{f}}^+} [\mathbb{E}[h_{\text{pert}}(\tilde{\mathbf{G}})]] - 0.17 \frac{\mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}]}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2}, \quad (72)$$

with the same  $C_1$  defined in (69). Combining (71) and (72) shows that

$$\mathbb{E}_{\mathbf{G} \sim \eta_{\tilde{f}}^+} [\mathbb{E}[h_{\text{pert}}(\tilde{\mathbf{G}})]] - \mathbb{E}_{\mathbf{G} \sim \eta_{\tilde{f}}^-} [\mathbb{E}[h_{\text{pert}}(\tilde{\mathbf{G}})]] \geq 0.17 \frac{\mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] + \mathbb{E}_{\eta_{\tilde{f}}^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}]}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2}. \quad (73)$$

Now using the duality principle (see 48) we conclude that

$$\text{disc}(\hat{f}).\mathcal{L}(h) \geq 0.17 \frac{\mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] + \mathbb{E}_{\eta_{\tilde{f}}^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}]}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2}. \quad (74)$$

Note that

$$\mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] = \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}})) | \mathcal{A}}] \mathbb{P}(\mathcal{A}) + \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}})) | \mathcal{A}^c}] \mathbb{P}(\mathcal{A}^c) \quad (75)$$

$$\leq \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}})) | \mathcal{A}}] + \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}})) | \mathcal{A}^c}] \mathbb{P}(\mathcal{A}^c) \quad (76)$$

$$\leq \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \quad (77)$$

$$+ 4.34\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}})) | \mathcal{A}^c}] \quad (78)$$

$$\leq \mathbb{E}_{\eta_{\tilde{f}}^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \quad (79)$$

$$+ 4.34\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \quad (80)$$

$$(81)$$

Similarly, one can obtain

$$\mathbb{E}_{\eta_f^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \leq \mathbb{E}_{\eta_f^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \quad (82)$$

$$+ 4.34\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \quad (83)$$

$$(84)$$

Hence<sup>9</sup>,

$$\mathbb{E}[\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \leq \mathbb{E}_{\eta_f^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] + \mathbb{E}_{\eta_f^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \quad (85)$$

$$\leq \mathbb{E}_{\eta_f^+} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] + \mathbb{E}_{\eta_f^-} [\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] \quad (86)$$

$$+ 8.68\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \quad (87)$$

Now combining with 74 results in

$$\text{disc}(\hat{f})\mathcal{L}(h) \geq \frac{0.17}{\left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^2} \max \left\{ \inf_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}, \quad (88)$$

$$\mathbb{E}[\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] - 8.68\delta \left(1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}})\right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \right\}. \quad (89)$$

#### B.4 ANALYSIS OF THE LIPSCHITZ CONSTANT

Now let us analyze  $\mathcal{L}(h)$ . By the definition

$$\mathcal{L}(h) = \sup_{\substack{\mathbf{G}, \mathbf{G}' \in \mathbb{G}_n^* \\ \mathbf{G} \neq \mathbf{G}'}} \frac{|\mathbb{E}[h_{\text{pert}}(\mathbf{G})] - \mathbb{E}[h_{\text{pert}}(\mathbf{G}')]|}{\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)} \quad (90)$$

$$\stackrel{(a)}{\leq} \sup_{\substack{\mathbf{G}, \mathbf{G}' \in \mathbb{G}_n^* \\ \mathbf{G} \neq \mathbf{G}'}} \frac{\mathbb{E}_\pi[|h_{\text{pert}}(\mathbf{G}) - h_{\text{pert}}(\mathbf{G}')|]}{\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)}, \quad (91)$$

where in (a) we couple the two perturbations into one according to the transportation map  $\pi \in \mathcal{S}_n$  that achieves  $\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)$ . (see 3). Fix  $\mathbf{G}, \mathbf{G}' \in \mathbb{G}_n^*$  and the coupled perturbations, and without loss of generality, assume that  $\pi = \text{id}$ . Note that the initial attributes considered for  $h_{\text{pert}}(\mathbf{G})$  are  $\mathbf{x}_v = \mu_v + \mathbf{n}_v$  and  $\mathbf{e}_{(u,v)} = \rho_{(u,v)} + \mathbf{n}_{(u,v)}$  for  $v \in [n]$  and  $(u, v) \in \mathcal{E}_{\mathbf{G}}$ , and similarly for  $h_{\text{pert}}(\mathbf{G}')$  they are  $\mathbf{x}'_v = \mu'_v + \mathbf{n}_v$  and  $\mathbf{e}'_{(u,v)} = \rho'_{(u,v)} + \mathbf{n}_{(u,v)}$  for  $v \in [n]$  and  $(u, v) \in \mathcal{E}_{\mathbf{G}'}$ . Let us define a few auxiliary graphs as follows. Define  $\mathbf{G}_1$  the same as  $\mathbf{G}$ , except with new initial node attributes  $\mathbf{x}'_v = \mu'_v + \mathbf{n}_v$  for  $v \in [n]$ . Also, define  $\mathbf{G}_2$  the same as  $\mathbf{G}_1$ , except with new edge attributes  $\mathbf{e}_{2,(u,v)} = \rho_{(u,v)} + \mathbf{n}_{(u,v)}$  for all  $(u, v) \in \mathcal{E}_{\mathbf{G}} \cap \mathcal{E}_{\mathbf{G}'}$ , and  $\mathbf{e}_{2,(u,v)} = \rho'_{(u,v)} + \mathbf{n}_{(u,v)}$  for all  $(u, v) \in \mathcal{E}_{\mathbf{G}'} \setminus \mathcal{E}_{\mathbf{G}}$ , and zero elsewhere.

One can write

$$|h_{\text{pert}}(\mathbf{G}) - h_{\text{pert}}(\mathbf{G}')| \leq |h_{\text{pert}}(\mathbf{G}) - h(\mathbf{G}_1)| + |h(\mathbf{G}_1) - h(\mathbf{G}_2)| + |h(\mathbf{G}_2) - h_{\text{pert}}(\mathbf{G}')|. \quad (92)$$

For the first term by the mean-value theorem,

$$|h_{\text{pert}}(\mathbf{G}) - h(\mathbf{G}_1)| \stackrel{(a)}{=} \left| \sum_{v \in [n]} \langle \partial_v h_{\text{pert}}(\mathbf{G})(\xi_v), (\mu_v + \mathbf{n}_v) - (\mu'_v + \mathbf{n}_v) \rangle \right| \quad (93)$$

<sup>9</sup>This expectation is with respect to the probability measure  $\mathbb{P}$  which is the law of graph generation.

$$\stackrel{(b)}{\leq} \sum_{v \in [n]} \|\partial_v h_{\text{pert}}(\mathbf{G})(\xi_v)\|_2 \|\mu_v - \mu'_v\|_2 \quad (94)$$

$$\leq L_{\max} \sum_{v \in [n]} \|\mu_v - \mu'_v\|_2 \quad (95)$$

$$\leq L_{\max} \text{dist}(\mathbf{G}, \mathbf{G}'; \lambda), \quad (96)$$

where (a) holds for some  $\xi \in [\mu_v + \mathbf{n}_v, \mu'_v + \mathbf{n}_v]$ , meaning that it belongs to the line segment between two endpoints, and also (b) holds by the Cauchy–Schwarz inequality.

For the second term, consider  $\mathbf{G}_1$ , and for instance, assume that  $\mathbf{G}_2$  differs from  $\mathbf{G}_1$  by only adding one edge. Observe that adding an edge such as  $(u, v)$  to  $\mathbf{G}_1$ , without changing node attributes, is equivalent to adding  $\Psi_\ell(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)})$  to node  $u$  in  $\ell$ th iteration (i.e., to  $\mathbf{h}_u^{(\ell-1)}$ ) and the same for  $v$ . Note that each  $\mathbf{h}_u^{(\ell-1)}$  in above is computed for  $\mathbf{G}_1$ , so the perturbations are also considered there. Thus, using the mean-value theorem,

$$|h(\mathbf{G}_1) - h(\mathbf{G}_2)| \stackrel{(a)}{=} \left| \sum_{\ell \in [k]} \langle \Psi_\ell(\mathbf{h}_v^{(\ell-1)}, \mathbf{e}_{(u,v)}), \partial_{(u,\ell-1)} h(\mathbf{G}_1)(\xi_{u,\ell}) \rangle \right| \quad (97)$$

$$+ \left| \sum_{\ell \in [k]} \langle \Psi_\ell(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)}), \partial_{(v,\ell-1)} h(\mathbf{G}_1)(\xi_{v,\ell}) \rangle \right| \quad (98)$$

$$\stackrel{(b)}{\leq} \sum_{\ell \in [k]} \|\Psi_\ell(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)})\|_2 \|\partial_{(u,\ell-1)} h(\mathbf{G}_1)(\xi_{u,\ell})\|_2 \quad (99)$$

$$+ \sum_{\ell \in [k]} \|\Psi_\ell(\mathbf{h}_v^{(\ell-1)}, \mathbf{e}_{(u,v)})\|_2 \|\partial_{(v,\ell-1)} h(\mathbf{G}_1)(\xi_{v,\ell})\|_2 \quad (100)$$

$$\leq L_{\max} \sum_{\ell \in [k]} \|\Psi_\ell(\mathbf{h}_u^{(\ell-1)}, \mathbf{e}_{(u,v)})\|_2 + L_{\max} \sum_{\ell \in [k]} \|\Psi_\ell(\mathbf{h}_v^{(\ell-1)}, \mathbf{e}_{(u,v)})\|_2 \quad (101)$$

$$\leq L_{\max} L_\Psi \sum_{\ell \in [k]} (\|\mathbf{h}_u^{(\ell-1)}\|_2 + \|\mathbf{h}_v^{(\ell-1)}\|_2) + 2k L_{\max} L_\Psi \|\mathbf{e}_{(u,v)}\|_2 \quad (102)$$

$$\leq 2k L_{\max} L_\Psi \left( \frac{1}{n} \sum_{v \in [n]} L_h \|\mathbf{x}_v\|_2 + \frac{1}{|\mathcal{E}_G|} \sum_{(u,v) \in \mathcal{E}_G} L_h \|\mathbf{e}_{(u,v)}\|_2 + B_h \right) \quad (103)$$

$$+ 2k L_{\max} L_\Psi \|\mathbf{e}_{(u,v)}\|_2 \quad (104)$$

$$\leq 2k L_{\max} L_\Psi \left( \frac{1}{n} \sum_{v \in [n]} (L_h \|\mu_v\|_2 + L_h \|\mathbf{n}_v\|_2 + n B_h) \right) \quad (105)$$

$$+ \frac{1}{|\mathcal{E}_G|} \sum_{(u,v) \in \mathcal{E}_G} L_h (\|\rho_{(u,v)}\|_2 + \|\mathbf{n}_{(u,v)}\|_2) \quad (106)$$

$$+ 2k L_{\max} L_\Psi (\|\rho_{(u,v)}\|_2 + \|\mathbf{n}_{(u,v)}\|_2), \quad (107)$$

where (a) holds for some  $\xi_{v,\ell}, \xi_{u,\ell}$  for  $\ell \in [k]$ , and (b) follows from the Cauchy–Schwarz inequality. Thus, if we take expectation from both sides with respect to the perturbation, we have the bound

$$\mathbb{E}_\pi |h(\mathbf{G}_1) - h(\mathbf{G}_2)| \leq 2k L_{\max} L_\Psi \left( \frac{1}{n} \sum_{\ell \in [k]} (L_h \|\mu_v\|_2 + L_h \sigma_n) \right) \quad (108)$$

$$+ \frac{1}{|\mathcal{E}_G|} \sum_{(u,v) \in \mathcal{E}_G} (L_h \|\rho_{(u,v)}\|_2 + L_h \sigma_e) + n B_h \quad (109)$$

$$+ 2k L_{\max} L_\Psi (\|\rho_{(u,v)}\|_2 + \sigma_e) \quad (110)$$

$$\leq 2k L_{\max} L_\Psi (L_h B_\mu + L_h B_\rho + L_h \sigma_n + L_h \sigma_e + B_h) \quad (111)$$

$$+ 2k L_{\max} L_\Psi (\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda) / \lambda + \sigma_e). \quad (112)$$

Note that in this case,  $\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda) \geq \lambda \|\rho_{(u,v)}\|_2 \geq \lambda$ . Similarly, if we consider more than one edge difference between  $\mathbf{G}_1, \mathbf{G}_2$ , the same result holds as

$$\mathbb{E}_\pi |h(\mathbf{G}_1) - h(\mathbf{G}_2)| \quad (113)$$

$$\leq \left( 2kL_{\max}L_{\Psi}(L_hB_{\mu} + L_h\sigma_n + L_hB_{\rho} + L_h\sigma_e + B_h) + 2kL_{\max}L_{\Psi}(1 + \sigma_e) \right) \frac{\text{dist}(\mathbf{G}_1, \mathbf{G}_2; \lambda)}{\lambda}. \quad (114)$$

Finally and similar to the previous arguments, we write

$$|h(\mathbf{G}_2) - h_{\text{pert}}(\mathbf{G}')| \leq \sum_{(u,v) \in \mathcal{E}_{\mathbf{G}_2}} \|\partial_{\mathbf{e}_{(u,v)}} h(\mathbf{G})(\xi_v)\|_2 \|(\rho_{(u,v)} + \mathbf{n}_{(u,v)}) - (\rho'_{(u,v)} + \mathbf{n}_{(u,v)})\|_2 \quad (115)$$

$$\leq L_{\max} \sum_{(u,v) \in \mathcal{E}_{\mathbf{G}_2}} \|\rho_{(u,v)} - \rho'_{(u,v)}\|_2 \quad (116)$$

$$\leq L_{\max} \frac{\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)}{\lambda}. \quad (117)$$

Combining (96), (114) and (117) shows that

$$\begin{aligned} \mathcal{L}(h) &\leq \sup_{\substack{\mathbf{G}, \mathbf{G}' \in \mathbb{G}_n^* \\ \mathbf{G} \neq \mathbf{G}'}} \frac{\mathbb{E}_{\pi} [|h_{\text{pert}}(\mathbf{G}) - h_{\text{pert}}(\mathbf{G}')|]}{\text{dist}(\mathbf{G}, \mathbf{G}'; \lambda)} \quad (118) \\ &\leq L_{\max} \left( 1 + \frac{1}{\lambda} (1 + 2kL_{\Psi}(L_hB_{\mu} + L_h\sigma_n + L_hB_{\rho} + L_h\sigma_e + B_h) + 2kL_{\Psi}(1 + \sigma_e)) \right). \quad (119) \end{aligned}$$

## B.5 FINAL STEP OF THE PROOF

Finally, we combine the above result and use (89) to achieve

$$\text{disc}(\hat{f}) \times L_{\max} \left( 1 + \frac{1}{\lambda} (1 + 2kL_{\Psi}(L_hB_{\mu} + L_h\sigma_n + L_hB_{\rho} + L_h\sigma_e + B_h) + 2kL_{\Psi}(1 + \sigma_e)) \right) \quad (120)$$

$$\geq \frac{0.17}{\left( 1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}}) \right)^2} \max \left\{ \inf_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}, \quad (121)$$

$$\mathbb{E}[\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] - 8.68\delta \left( 1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}}) \right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \right\}. \quad (122)$$

Now according to (45) we obtain

$$(\text{disc}(f) + 8.68 \left( 1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}}) \right)^4 \delta n (2B_{\mu} + \lambda B_{\rho} \Delta_{\max})) \quad (123)$$

$$\times L_{\max} \left( 1 + \frac{1}{\lambda} (1 + 2kL_{\Psi}(L_hB_{\mu} + L_h\sigma_n + L_hB_{\rho} + L_h\sigma_e + B_h) + 2kL_{\Psi}(1 + \sigma_e)) \right) \quad (124)$$

$$\geq \frac{0.17}{\left( 1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}}) \right)^2} \max \left\{ \inf_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}, \quad (125)$$

$$\mathbb{E}[\sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))}] - 8.68\delta \left( 1 + 2C_{\text{poly}}(1 + 3^{\frac{1}{2}k_{\text{poly}}}) \right)^4 \sup_{\mathbf{G} \in \mathbb{G}_n^*} \sqrt{\text{var}(h_{\text{pert}}(\tilde{\mathbf{G}}))} \right\}, \quad (126)$$

which completes the proof.

## C PROOF OF PROPOSITION 2

To prove Proposition 2, we introduce a natural class of simple graphs  $\mathbb{G}_n$  that simply allows to couple graphs with different labels.

- *Number of edges.* In this case, simply consider the set of graphs  $\mathbb{G}_n = \{G_0, G_1, \dots, G_m\}$  where  $m = \frac{n\Delta_{\max}}{2}$  constructed as follows. First let  $G_0$  be the empty graph and for each  $G_i$ , add an arbitrary edge to  $G_{i-1}$  to construct it. For the number of edges function, denoted by  $f : \mathbb{G}_n \rightarrow [0, \frac{n\Delta_{\max}}{2}]$ , each interval in the definition of  $f_p$  has length  $\frac{n\Delta_{\max}}{2p}$ . To couple two intervals, we just shift the index:  $i \rightarrow i + \frac{n\Delta_{\max}}{2p} \pm 1$  (since its an integer). This means that the discrepancy can be upper bounded by  $\frac{n\Delta_{\max}}{2p} + O(1)$ . The next part also follows from  $\epsilon_{\text{additive}} = \frac{n\Delta_{\max}}{2p}$  which is the length of interval.
- *Number of triangles, number of subgraphs  $H$ .* Let us just prove the desired result for the latter. For any graph  $H$ , construct the set  $\mathbb{G}_n = \{G_0, G_1, \dots, G_m\}$  where  $m = \frac{n}{n_H}$  constructed as follows. Start with empty graph  $G_0$ , and add independent copies of  $H$  (i.e., without overlapping nodes) sequentially to find  $G_i, i \in [m]$ . Now, for the subgraph count function  $f : \mathbb{G}_n \rightarrow [0, \frac{n}{n_H}]$ , we have the length of interval  $\epsilon_{\text{additive}} = \frac{n}{pn_H}$ . Similar to the previous part, to achieve a coupling, use index shift, and to increase the function by  $\frac{n}{pn_H}$  we need to add  $\frac{nm_H}{pn_H} \pm 1$  edges, which completes the proof.
- *Diameter.* Consider path graphs  $P_i$  for  $i \in \{0, 1, \dots, n-1\}$  as  $\mathbb{G}_n$ , and similar to the number of edges, construct the coupling. The proof is the same.
- *Max-cut, min-cut.* The two case are dual of each other if we compute the completion of the graphs (discrepancy bounds are invariant with respect to this operation). Thus, we only prove the max-cut case, and construct bipartite graph on  $\frac{n}{2} \pm 1$  nodes (since the number is integer), and add edges to greedily construct  $G_0, G_1, \dots, G_m$  for  $m = \frac{n\Delta_{\max}}{2} \pm 1$ . Then, since adding edges will directly change the function value, we conclude the result similar to the number of edges case.
- *Clique number, independence number.* By duality we only prove the result for the clique number. Start from  $G_0$  as empty set, and greedily enlarge a clique to achieve  $G_m$  which has a clique of size  $\Delta_{\max}$ . Then,  $\epsilon_{\text{additive}} = \frac{\Delta_{\max}}{p}$ , and to add this number we need to add at most  $\Delta_{\max} \times \frac{\Delta_{\max}}{p}$  new edges. The proof is thus complete.
- *Number of connected components.* This case essentially the same as the analysis for the diameter of graphs.