

Model-Agnostic Graph Dataset Compression with the Tree Mover’s Distance

Anonymous Authors¹

Abstract

Graph neural networks have demonstrated remarkable success across a variety of domains. However, the acquisition and management of large-scale graph datasets poses several challenges. Acquiring graph-level labels can be prohibitively costly, especially for applications in the biosciences and combinatorial optimization. Storage and privacy constraints can pose additional challenges. In this work, we propose an approach for data subset selection for graph datasets, which downsamples graphs and nodes based on the Tree Mover’s Distance. We provide new efficient methods for computing the TMD in our setting; empirical results showing our approach outperforms other node and graph sampling methods; and theoretical results bounding the decrease in accuracy caused by training on the downsampled graphs. Surprisingly, we find that with our method, we can subsample down to 1% of the number of graphs and 10% of the number of nodes on some datasets, with minimal degradation in model accuracy.

1. Introduction

Graph neural networks (GNNs) are a popular architecture for learning over graph-structured data. Despite their widespread popularity and demonstrated success (Bongini et al., 2022; Zhou et al., 2020; Peng et al., 2021; Fan et al., 2019; Peng et al., 2021; Zhang & Chen, 2018), training GNNs remains challenging because the datasets can be massive, both in terms of the size (number of nodes) of each graph, as well as the number of graphs in the dataset.

Large-scale datasets can pose many challenges to training GNNs. For instance, training may be computationally intractable (Hashemi et al., 2024; Yan et al., 2020; Zeng et al., 2021). Likewise, storing and loading massive graphs into memory can be slow and/or expensive (Huang et al., 2024;

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the Workshop on Advancing Neural Network Training at International Conference on Machine Learning (WANT@ICML 2024). Do not distribute.

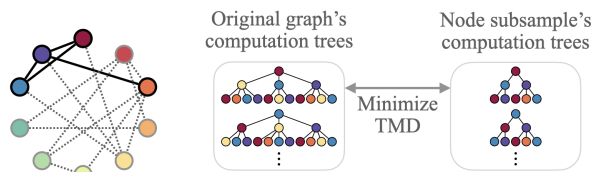


Figure 1. Node subsampling. Given node budget k , we want to find a subgraph of k nodes with minimum TMD from the original.

Hamilton et al., 2017; Da San Martino et al., 2012) or may raise privacy concerns (Wu et al., 2022). In addition, in many applications, labeling many graphs or labeling graphs with many nodes can be expensive. This is the case, for example, when GNNs are used to model molecular properties for computational drug discovery (Igarashi et al., 2024; Shen et al., 2020; Jiang et al., 2021). GNNs have also been employed to circumvent expensive simulation, e.g., in molecular dynamics (Gasteiger et al., 2021; Park et al., 2021; Li et al., 2022), where the cost of labeling larger instances grows quickly with graph size. Similarly, GNNs have been applied for disease prediction (Sun et al., 2020; Zheng et al., 2022), in which case labeling a single graph may require expensive clinical or diagnostic procedures. Additionally, there is growing interest in using GNNs to solve challenging or costly optimization problems, in which case the runtime to collect training labels grows with the graph size (e.g., Bengio et al., 2021; Cappart et al., 2023; Vesselinova et al., 2020).

Subsampling is a natural approach to tackle these multifaceted challenges. In node and graph subsampling, we aim to reduce the number of nodes or graphs, respectively, in a graph dataset while preserving the performance of a downstream GNN trained on the dataset. Existing work on node and graph subsampling focuses primarily on node-level tasks (e.g., node classification), and often makes strong assumptions about the GNN model architecture—and can even assume access to the final model or the ability to approximate the training trajectories on the original, unsampled dataset.

In this paper, we propose and theoretically motivate new approaches for (1) subsampling nodes from each training graph and (2) subsampling entire graphs in the training set for graph-level learning tasks (e.g., graph classification

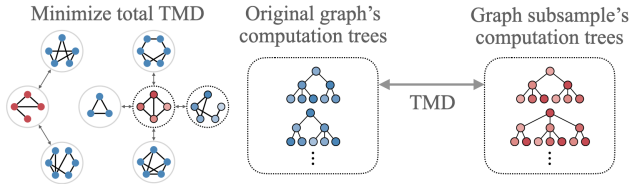


Figure 2. Graph subsampling. Given a graph budget k , the goal is to find a sample of k graphs that minimizes the total distance (TMD) from every original graph to its closest sample.

and regression). We must ensure that training a GNN on the subsampled data does not significantly reduce accuracy compared to the original training set. Moreover, in practice, we often do not know *a priori* (i.e., before collecting, compressing, and labeling our graph dataset) the architecture or hyperparameters of the GNN. Consequently, it is crucial that the original and subsampled datasets maintain *model-agnostic similarity*, i.e., they produce similar model behavior regardless of the downstream model.

For many GNNs, specific *structural* similarities—tied to the GNN’s computations—are essential to successfully transfer a GNN out of distribution (Yehudai et al., 2021; Chuang & Jegelka, 2022; Ruiz et al., 2021; Levie et al., 2021; Le & Jegelka, 2023). A GNN uses *message-passing* to compute predictions about graphs. Initially, each node has an embedding, which it updates by aggregating its neighbors’ embeddings and processing the result through a neural network. This process repeats over L iterations—or *layers*—and the nodes’ final embeddings are used to compute the GNN’s predictions. The L -hop neighbors involved in a node’s final embedding can be visualized via the node’s *computation tree*, illustrated in Figure 3. Computation trees play a critical role in the transfer and stability of GNNs (Yehudai et al., 2021; Chuang & Jegelka, 2022; Georgiev et al., 2022). Two nodes with similar computation trees will likely have similar embeddings. To capture this intuition, Chuang & Jegelka (2022) introduced a graph distance—the *Tree Mover’s Distance (TMD)*—as an optimal transport distance between sets of computation trees, and demonstrated its correlation with the performance and stability of GNNs under data perturbations, including generalization bounds.

Based on the strong theoretical foundations provided by TMD, our primary goal is to ensure that the TMD between the original and subsampled graphs is small. However, achieving this goal requires computing the TMD between many graphs while subsampling, which poses a significant challenge. TMD is a combinatorially complex function: the dynamic programming algorithm developed in prior work for computing the TMD between graphs with n nodes has a runtime of $\mathcal{O}(Ln^4)$ (Chuang & Jegelka, 2022). Given our objective to reduce the size of large graphs, this runtime is prohibitively large and presents a computational challenge

that we overcome in our methodological contributions.

Contributions. In this work, we make the following contributions: We present an approach for node subsampling that, given a budget k and a training graph, returns a subgraph on k nodes with low TMD from the original graph (Figure 1). Moreover, we present an approach for graph subsampling that, given a training set of graphs, returns k graphs such that the total TMD between each original graph and its closest sampled graph is small (Figure 2).

We show that if h is a GNN with minimal loss on the training set subsampled using our methods, then h also has minimal loss on the original training set. When subsampling nodes in each training graph, we bound the resulting increase in loss by the average TMD between the original and subsampled graphs. When subsampling entire graphs, we bound the increase in loss by the average TMD between each training graph and the nearest graph in the subsampled set.

Given a set of candidate subgraphs, we develop a faster algorithm for selecting a subgraph with minimum TMD from a given graph. It computes the depth- L TMD between a graph $G = (V, E)$ and any of its subgraphs in just $\mathcal{O}(L|E|)$ time. This algorithm leverages structural properties of the TMD that we uncover, which may be of independent interest.

Our experiments demonstrate our methods outperform or perform competitively with other standard approaches on real-world graph learning benchmarks. For graph sampling, our method consistently outperforms the baselines on four benchmark datasets, achieving over 90% of the test accuracy of the full training set for three of the datasets using 1% of the graphs. For node subsampling, it is consistently optimal or near-optimal across all datasets, unlike other competitive baselines, whose performance varies significantly with the fraction of deleted data.

2. Related work

Graph reduction, graph sparsification, and graph condensation approaches aim to reduce the size of a graph while minimizing the impact on test loss. Existing approaches focus on node subsampling (e.g., Hashemi et al., 2024; Tanaka et al., 2020, and references therein), primarily for *node classification* and graph signal processing. These approaches involve subsampling nodes, and do not extend to graph subsampling, a critical focus of this work, and a setting which remains largely unstudied in the graph reduction literature. Two existing approaches to *model-agnostic* graph reduction for graph classification are Herding (Welling, 2009) and k -Centers (Farahani & Hekmatfar, 2009), which subsample k nodes by clustering the nodes’ features into k clusters and selecting the centers. A key advantage of our approach is that, unlike these approaches, we leverage for the node features *and* the graph structure.

Gradient-matching approaches (Jin et al., 2022; Zhang et al., 2024; Fang et al., 2024) have been used for graph dataset subsampling and ensure that the loss’s gradients remain similar when trained on the small and large graphs. While these are powerful methods that can aid in developing reduced graph datasets that are smaller and more interpretable, they are fundamentally incomparable in our setting because (1) they assume knowledge of a downstream GNN architecture and therefore are not *model-agnostic*; and (2) they require training on the complete, unreduced graphs to generate the gradient trajectories. This obviates several of the advantages of graph reduction with regard to labeling, training, and storage costs.

Georgiev et al. (2022) uses a normalized variant of TMD to select training datasets for neural algorithmic reasoning tasks. However, this is not agnostic to the downstream model and task and is not directly applicable to graph or node subsampling.

3. Notation and background

Notation. An undirected, unweighted graph $G = (V, E, f)$ has a node set V , edges E , and node features $f \in \mathbb{R}^{p \times |V|}$, where f^v is the feature vector of $v \in V$. We use $N_G(v)$ to denote the neighborhood of $v \in V$ and \mathcal{G} to denote a set of graphs. An r -rooted tree is denoted by $T = (V, E, f, r)$. We use $\mathbf{1}, \mathbf{0}, \mathbf{I}$ to denote the all ones and zeros vectors and identity matrix, respectively. For norms, $\|\cdot\|$ denotes any norm over \mathbb{R}^p and $\|\cdot\|_1$ is the ℓ_1 norm.

Message-passing neural networks (GNNs). We analyze *message-passing graph neural networks (GNNs)*, which operate on graphs $G = (V, E, f)$. For simplicity, we define *Graph Isomorphism Networks (GINs)* (Xu et al., 2019) here; however, the results of Chuang & Jegelka (2022) and, therefore, our results also apply to other GNNs. A GIN is defined by L (learnable) functions $\phi^{(\ell)} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, such as (shallow) feed-forward neural networks, and a fixed weight $\eta > 0$. Each node begins with an initial embedding $z_v^{(0)} = f^v$ and updates its embedding using the following message-passing routine:

Message passing layers: for $\ell \in [L - 1]$,

$$z_v^{(\ell)} = \phi^{(\ell)} \left(z_v^{(\ell-1)} + \eta \sum_{u \in N(v)} z_u^{(\ell-1)} \right),$$

Graph readout: $h(G) = \phi^{(L)} \left(\sum_{v \in V} z_v^{(L-1)} \right).$

Other GNN architectures may, for instance, replace the summations with an average (as in the case of Graph Convolutional Neural Networks (GCNs) (Kipf & Welling, 2016)) or other aggregation functions (Hamilton, 2020).

Tree Mover’s Distance (TMD). The TMD was introduced by Chuang & Jegelka (2022) and has found applica-

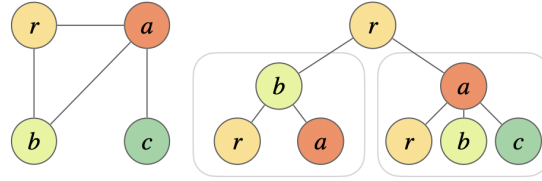


Figure 3. A graph and the depth-2 computation tree T of node r . The set $\mathcal{T}_r(T)$ consists of the two boxed subtrees.

tions for graph learning under distribution shifts (Yu et al., 2023; Georgiev et al., 2022) and adversarially robust graph learning (Schuchardt et al., 2024). TMD is a pseudometric on graphs based on the optimal transport distance. Like a message passing GNN, it views a graph as a set of *computation trees*. A node’s computation tree is constructed by adding the node’s neighbors to the tree level by level, as in Figures 1 and 3.

Definition 3.1 (Computation tree). Given $G = (V, E, f)$ and a node $v \in V$, let $T_v^1(G) = (v, \emptyset, \{f^v\}, v)$ be a singleton v -rooted tree that consists of only the node v , no edges, and the node feature vector f^v . Inductively, we define $T_v^L(G)$ to be the depth- L computation tree of node v . That is, $T_v^L(G)$ is a tree rooted at v obtained as follows. For each leaf ℓ of $T_v^{L-1}(G)$ and each neighbor $u \in N_G(\ell)$, we add a new node ℓ_u to $T_v^L(G)$ and add an edge from ℓ to ℓ_u . The multiset of depth- L computation trees defined by G is denoted by $\mathcal{T}_G^L := \{T_v^L(G)\}_{v \in V}$.

We can compare graphs by comparing their computation trees. TMD does this via optimal transport.

Definition 3.2 (OT distance). Let $X = \{x_i\}_{i=1}^q$ and $Y = \{y_i\}_{i=1}^q$ be two multisets of q elements each. Given a distance metric $d : X \times Y \rightarrow \mathbb{R}$, let $C \in \mathbb{R}^{q \times q}$ be a matrix where $C_{i,j} = d(x_i, y_j)$ is the *transportation cost* between x_i and y_j . Moreover, let the set of all *transportation plans* between X and Y be $\Gamma(X, Y) := \{\gamma \in \mathbb{R}_+^{q \times q} \mid \gamma \mathbf{1} = \gamma^\top \mathbf{1} = \mathbf{1}\}$. The *transportation cost* of $\gamma \in \Gamma$ is $\langle \gamma, C \rangle := \sum_{i,j} \gamma_{i,j} C_{i,j}$. The (unnormalized) Wasserstein distance OT_d is defined as $\text{OT}_d(X, Y) := q \min_{\gamma \in \Gamma(X, Y)} \langle C, \gamma \rangle$. We say γ^* is an *optimal transport (OT) plan* if $\gamma^* \in \text{argmin}_{\gamma \in \Gamma(X, Y)} \langle C, \gamma \rangle$.

We need a metric between trees to compute the OT distance between sets of computation trees. TMD compares two trees T_v, T_u by comparing their roots f_v, f_u and recursively comparing their subtrees.

Definition 3.3 (Tree multisets). Given an r -rooted tree $T = (V, E, f, r)$, let $\ell = \text{depth}(T)$. We use $\mathcal{T}_r(T)$ to denote the multiset of depth $(\ell - 1)$ computation trees of all neighbors $u \in N_T(r)$, illustrated in Figure 3. (In other words, if T_u is the subtree of T rooted at u , then

$\mathcal{J}_r(T) := \cup_{u \in N_T(r)} T_u^{(\ell-1)}(T_u)$, where the union denotes a multiset union).

TMD compares $\mathcal{J}_u(T)$ and $\mathcal{J}_v(T')$ via an OT distance. However, the number of subtrees could differ in $\mathcal{J}_u(T)$ and $\mathcal{J}_v(T')$. Chuang & Jegelka (2022) therefore augment the smaller set with *blank trees*. A blank tree T_0 consists of a single node with node features $\mathbf{0} \in \mathbb{R}^d$ and no edges. Given two tree multisets $\mathcal{J}_v(T')$ and $\mathcal{J}_u(T)$, the function $\rho(\mathcal{J}_v(T'), \mathcal{J}_u(T))$ returns two multisets of the same size, where the smaller is padded with blank trees (see Appendix A.3 for more details).

We are now ready to define a distance metric between trees.

Definition 3.4 (Tree distance). Let $T_a = (V_a, E_a, f_a, r_a)$ and $T_b = (V_b, E_b, f_b, r_b)$ be two trees with $\ell := \max(\text{depth}(T_a), \text{depth}(T_b))$. Moreover, let $w : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ be a depth-dependent weighting function. The *tree distance* (TD) between T_a and T_b is defined as

$$\text{TD}_w(T_a, T_b) := \begin{cases} \|f_a^{r_a} - f_b^{r_b}\| \\ + \begin{cases} w(\ell) \cdot \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_{r_a}(T_a), \mathcal{J}_{r_b}(T_b))), & \text{if } \ell > 1 \\ 0, & \text{otherwise.} \end{cases} \end{cases}$$

Here, OT_{TD_w} is the OT distance with metric TD_w . We now define the tree mover’s distance (TMD), which calculates the cost of the optimal transport plan between two graphs’ computation trees.

Definition 3.5 (Tree mover’s distance (TMD) (Chuang & Jegelka, 2022)). Given graphs G and G' , weight function $w : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$, and a depth parameter $L > 0$, we define $\text{TMD}_w^L(G, G') := \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_G^L, \mathcal{J}_{G'}^L))$. The *tree norm* is $\|G\|_{\text{TN}_w^L} := \text{TMD}_w^L(G, \emptyset)$ where we overload \emptyset to denote the *empty graph*.

4. Node subsampling

In this section, we present an approach to subsampling the nodes of a graph dataset while preserving the performance of a downstream GNN trained on the subsampled dataset. Full proofs for all results in this section can be found in Appendix B. To this end, suppose we are given a graph dataset $X = \{G_1, \dots, G_n\}$ where $G_i = (V_i, E_i)$ together with a node budget k . Given a subset of k nodes $S_i \subseteq V_i$, we use $G_i[S_i]$ to denote the subgraph of G_i induced by S_i . Our goal is to select these subsets so that a GNN obtains similar readouts on the original training dataset X and on the induced subgraphs $X' = \{G_1[S_1], \dots, G_n[S_n]\}$. Although quantifying the stability of the GNN readouts with respect to node deletion directly is analytically intractable, it is possible to bound the Lipschitz constant of a message-

passing GNN’s readouts in terms of the TMD.¹

Theorem 4.1 (Informal, Theorem 8 in (Chuang & Jegelka, 2022), restated). *There exists a weight function w such that given an $(L - 1)$ -layer message-passing GNN with readout function $h : \mathcal{G} \mapsto \mathbb{R}^d$ and ℓ -th layer Lipschitz constants ϕ_ℓ , for any two graphs G_a, G_b , $\|h(G_a) - h(G_b)\| \leq \left(\prod_{\ell=1}^{L-1} \phi_\ell\right) \cdot \text{TMD}_w^L(G_a, G_b)$.*

Theorem 4.1 implies that if the Lipschitz constants of the GNN layers and the distances $\text{TMD}(G_i[S_i], G_i)$ are small, then minimizing a Lipschitz loss over the subsampled training set X' yields a nearly optimal hypothesis over X . The following corollary summarizes this observation.

Corollary 4.2. *Let \mathcal{H} be a hypothesis class of $(L - 1)$ -layer GNNs $h : \mathcal{G} \rightarrow \mathbb{R}^d$ with ℓ -th layer Lipschitz constant upper-bounded by Φ_ℓ . Let $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$ be an M -Lipschitz loss function. Let $X = (G_1, \dots, G_n)$ and (y_1, \dots, y_n) be a set of (training) graphs and their labels. For each $i \in [n]$, let $G_i = (V_i, E_i)$ and $S_i \subset V_i$. Suppose that $M \left(\prod_{\ell \in [L-1]} \Phi_\ell\right) \leq c$ and $\text{TMD}_w^L(G_i, G_i[S_i]) \leq \varepsilon_i$ for all $i \in [n]$. Finally, let $\hat{h} \in \mathcal{H}$ be a hypothesis with minimum loss over the subsampled training set:*

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i \in [n]} \mathcal{L}(h(G_i[S_i]); y_i).$$

Then \hat{h} has nearly optimal loss over the original training set as well:

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) + \frac{2c}{n} \sum_{i \in [n]} \varepsilon_i.$$

Corollary 4.2 illustrates that the loss incurred by training on X' as opposed to X is bounded by an error proportional to the average of the distances $\text{TMD}_w^L(G_i, G_i[S_i])$. Motivated by this bound, our ideal goal would be to compute the optimal solution to the following optimization problem:

$$\min_{S \subset V : |S| \leq k} \text{TMD}_w^L(G, G[S]). \quad (1)$$

We face two key challenges in computing (approximately) optimal solutions to Equation (1).

Challenge 1: exponentially large feasible set. The first challenge is that the number of candidate subsets $\{S \subset V : |S| \leq k\}$ grows exponentially with k . Therefore, we restrict Equation (1) to an appropriately chosen feasible set \mathcal{S} . We discuss a natural heuristic for selecting \mathcal{S} in Section 4.3.

¹While Theorem 8 as stated in (Chuang & Jegelka, 2022) is specific to GINs, Chuang & Jegelka (2022) generalize it to other GNNs.

Definition 4.3 (Relaxed node subsampling problem). Let $G = (V, E, f)$ be a graph, k be a node budget, and $\mathcal{S} \subset \{S \subset V : |S| \leq k\}$ be a set of candidate node subsets. The *relaxed node subsampling problem* is defined as: $\min_{S \in \mathcal{S}} \text{TMD}_w^L(G, G[S])$.

Challenge 2: computing $\text{TMD}_w^L(G, G[S])$ is computationally expensive. Unfortunately, even this relaxed node subsampling problem is computationally expensive, as the TMD is expensive to compute even for a single candidate $S \in \mathcal{S}$. The original dynamic-programming algorithm for computing $\text{TMD}_w^L(G, G[S])$ requires $\mathcal{O}(|V|^4)$ -time in general (Chuang & Jegelka, 2022) (regardless of $|S|$), bringing the total runtime to $\mathcal{O}(|V|^4 L |\mathcal{S}|)$. Since we aim to reduce the size of large graphs, this runtime is likely prohibitive in practice. To address this challenge, we prove that, surprisingly, Definition 4.3 is equivalent to a much simpler optimization problem which *can* be solved efficiently.

Theorem 4.4. *Let $G = (V, E, f)$ be a graph, k be a node budget, and $\mathcal{S} \subset \{S \subset V : |S| \leq k\}$. Then S^* is an optimal solution to the relaxed node subsampling problem if and only if*

$$S^* \in \operatorname{argmax}_{S \in \mathcal{S}; |S|=k} \|G[S]\|_{\text{TN}_w^L}. \quad (2)$$

We prove Theorem 4.4 in Section 4.1. Along the way, we prove several structural properties of the TMD, which may be of independent interest. Finally, since the algorithm of Chuang & Jegelka (2022) would take $\mathcal{O}(L |S|^4)$ -time to compute the tree norm of $G[S]$, we provide a new linear-time algorithm for computing tree norms—Algorithm 1—that better leverages the graph’s sparsity.

Theorem 4.5. *Given a graph $G = (V, E, f)$, Algorithm 1 computes $\|G\|_{\text{TN}_w^L}$ in $\mathcal{O}(|E| L)$ time.*

We discuss the algorithm and proof sketch in Section 4.2; however, we first point out that this implies a more tractable algorithm for the relaxed node subsampling problem.

Remark 4.6. There is an algorithm that solves the relaxed node subsampling problem in $\mathcal{O}(k^2 |S|)$ time. It computes $\|G[S]\|_{\text{TN}_w^L}$ for each $S \in \mathcal{S}$ and outputs the set S with the largest tree norm.

4.1. Analysis of TMD between graphs and subgraphs

In this section, we sketch our proof of Theorem 4.4. Along the way, we prove several properties of the TMD, which may be of independent interest given the broad applications of TMD to out-of-distribution generalization. Our analysis crucially relies on the following technical lemma.

Lemma 4.1. *Let $w : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ be a weight function and $G = (V, E, f)$ be a graph. Then, the identity transportation*

plan \mathbf{I} that maps $T_v(G)$ to $T_v(G[S])$ for $v \in S$, and $T_u(G)$ to \emptyset for $u \notin S$ is an optimal transport plan for the OT problem

$$\text{TMD}_w^L(G, G[S]) = \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{F}_G^L, \mathcal{F}_{G[S]}^L \right) \right). \quad (3)$$

Consequently, we can decompose the right-hand-side of (3) into three terms: a sum over the feature norms of deleted nodes $v \notin S$; the cost of transporting the deleted nodes’ computation trees (that is, $\{\mathcal{F}_v(T_v^L(G))\}_{v \notin S}$) to the empty set; and the cost of transporting the computation trees of the remaining nodes $v \in S$ (that is, $\{\mathcal{F}_v(T_v^L(G))\}_{v \in S}$) to their computation trees in $G[S]$. Formally,

$$\begin{aligned} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{F}_G^L, \mathcal{F}_{G[S]}^L \right) \right) &= \sum_{v \notin S} \|f^v\| \\ &+ w(L-1) \sum_{v \notin S} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{F}_v(T_v^L(G)), \emptyset \right) \right) \\ &+ w(L-1) \sum_{v \in S} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{F}_v(T_v^L(G)), \mathcal{F}_v(T_v^L(G[S])) \right) \right). \end{aligned}$$

Proof sketch. We prove the statement by induction on L . In the base case ($L = 1$), the statement is equivalent to Definition 3.5. This is because by the definition of ρ and the tree distance TD, $\text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_G^1, \mathcal{F}_{G[S]}^1))$ is equal to the sum of features of the nodes that are in G but not in $G[S]$. That is, $\text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_G^1, \mathcal{F}_{G[S]}^1)) = \sum_{v \notin S} \|f^v\| = \langle C, \mathbf{I} \rangle$. This shows that the lemma holds in the base case. We prove the inductive step by using the recursive formulation of the TD to reduce OT problems of depth L to OT problems of depth $L - 1$, which allows us to apply the inductive hypothesis. ■

Using Lemma 4.1, we can precisely characterize the TMD induced by dropping a sequence of nodes from the graph. First, we prove the following chain rule.

Lemma 4.2. *For $T \subset S \subset V$, $\text{TMD}_w^L(G, G[S \setminus T]) = \text{TMD}_w^L(G, G[S]) + \text{TMD}_w^L(G[S], G[S \setminus T])$.*

Proof sketch. To prove this lemma, we apply the expression in Lemma 4.1 inductively over L to decompose $\text{TMD}_w^L(G, G[S \setminus T])$ into $\text{TMD}_w^L(G, G[S])$ and $\text{TMD}_w^L(G[S], G[S \setminus T])$. To highlight the intuition, we discuss the base case here and discuss the inductive step in the full proof. If $L = 1$, then the second and third terms in the expression of Lemma 4.1 are always 0, because \mathcal{F}_G^1 and $\mathcal{F}_{G[S]}^1$ are multisets of depth 1, so $\mathcal{F}_G(T_v^1(G))$ and $\mathcal{F}_{G[S]}(T_v^1(G[S]))$ are empty sets. Thus,

$$\begin{aligned} \text{TMD}_w^1(G, G[S \setminus T]) &= \sum_{v \notin (S \setminus T)} \|f^v\| \\ &= \sum_{v \notin S} \|f^v\| + \sum_{v \in T} \|f^v\|, \end{aligned}$$

Algorithm 1: TreeNorm(G, L, w)

Input: Graph $G = (V, E, f)$ with adjacency matrix A , weights $w : \{1, \dots, L - 1\} \rightarrow \mathbb{R}_+, L \geq 1$

- 1 Compute $x \in \mathbb{R}^{|V|}$ to be the vector such that $x_v = \|f^v\|$;
- 2 Initialize $z^{(0)} = x$;
- 3 **for** $\ell \in [L - 1]$ **do**
- 4 $z^{(\ell)} \leftarrow Az^{(\ell-1)}$;
- 5 $b \leftarrow z^{(0)} + \sum_{\ell=1}^{L-1} \left(\prod_{t=1}^{\ell} w(L - t) \right) \cdot z^{(\ell)}$;

return: $\|b\|_1$

(since $T \subset S, \{v \in V : v \notin (S \setminus T)\} = \{v \in V : v \notin S\} \cup \{v \in T\}$). Similarly applying Lemma 4.1 to $\text{TMD}_w^L(G, G[S \setminus T])$ and $\text{TMD}_w^1(G[S], G[S \setminus T])$ implies that the base case holds:

$$\begin{aligned} \text{TMD}_w^1(G, G[S]) &= \sum_{v \notin S} \|f^v\|, \text{ and} \\ \text{TMD}_w^1(G[S], G[S \setminus T]) &= \sum_{v \in T} \|f^v\|. \end{aligned}$$

The intuition is similar for $L > 1$ but requires care to handle the recursive OT terms. ■

Finally, we can prove Theorem 4.4.

Proof of Theorem 4.4. Let $S \in \mathcal{S}$. Applying Lemma 4.2 with $T = V \setminus S$, we have that $\|G\|_{\text{TN}_w^L} = \text{TMD}_w^L(G, G[S]) + \|G[S]\|_{\text{TN}_w^L}$. Consequently,

$$\begin{aligned} & \max_{S \in \mathcal{S}; |S|=k} \|G[S]\|_{\text{TN}_w^L} \\ & \equiv \max_{S \in \mathcal{S}; |S|=k} \|G\|_{\text{TN}_w^L} - \text{TMD}_w^L(G, G[S]) \\ & \equiv \min_{S \in \mathcal{S}; |S|=k} \text{TMD}_w^L(G, G[S]). \end{aligned} \quad \blacksquare$$

4.2. Faster algorithm for computing tree norms

We next leverage Lemma 4.1 to obtain Algorithm 1, a faster algorithm for computing $\|G\|_{\text{TN}_w^L}$ for $G = (V, E, f)$. The original implementation would require $O(L|V|^4)$ time (Chuang & Jegelka, 2022), whereas ours has runtime $O(L|E|)$. Algorithm 1 is based on the fact that by Lemma 4.1, is that $\|G\|_{\text{TN}_w^L}$ is essentially a weighted sum of the number of vertices in all of its L -hop computation trees. For each tree, nodes in level ℓ receive weight $w(L - \ell + 1)$. Algorithm 1 computes this weighted sum. Its runtime is dominated by the cost of L matrix-vector multiplies with the adjacency matrix, which requires $O(|E|)$ -time. The proof of correctness and runtime (Theorem 4.5) is in Appendix B.

4.3. Heuristic for selecting \mathcal{S}

We now present our heuristic for constructing the set of candidate subsets \mathcal{S} in our experiments. Most real-world

networks are scale-free and small-world, and hence can be well-modeled by random graphs, such as preferential attachment, configuration models, and inhomogeneous random graphs (Bollobás & Riordan, 2004; Newman & Watts, 1999). These graphs converge *locally* to graph limits (Van Der Hofstad, 2024, Vol. 2, Ch. 2). By a “transitive” argument, we can view real-world networks as parts of sequences converging to graph limits in the same sense. Alimohammadi et al. (2023) showed that for such graphs, sampling nodes’ local breadth-first search (BFS) trees asymptotically preserves critical motifs of the graph limit (see Appendix A). Thus, we define the set of candidate subsets \mathcal{S} as a set of BFS trees rooted at the graph’s nodes.

Definition 4.7 (k -BFS subset). Given $G = (V, E, f)$ and $v \in V$, let ℓ_k be the deepest level such that the v -rooted BFS tree of depth ℓ_k has at most k nodes (breaking ties in a fixed but arbitrary way). The k -BFS subset of v , denoted $S_{\text{BFS}(v;k)}$, is the set of nodes at distance at most ℓ_k from v .

We define the set of candidate subsets $\mathcal{S} = \{S_{\text{BFS}(v;k)}\}_{v \in V}$ as the nodes’ k -BFS subsets.

5. Graph subsampling

We now present an approach that uses TMD to subsample the number of graphs in a dataset while preserving the performance of a downstream GNN trained on the subsampled dataset. Omitted proofs are in Section B. Let $X = \{G_1, \dots, G_n\}$ be a graph dataset. Given a budget k , the goal is to identify a subset $\mathcal{G} \subset [n]$ of size k such that a GNN obtains similar readouts and loss on the original training set X and on the subsampled set $\{G_i\}_{i \in \mathcal{G}}$. We must ensure that the selected set \mathcal{G} is representative of the full dataset. To this end, we select \mathcal{G} by optimizing the medoids objective, which captures how well each graph $G_i \in X$ is represented by \mathcal{G} .

Definition 5.1 (Medoids objective). Let $X = \{G_1, \dots, G_n\}$ be a graph dataset and $\mathcal{G} \subset [n]$ with $|\mathcal{G}| \leq k$. The medoids objective value of \mathcal{G} on X with respect to TMD_w^L is defined

$$f_{\text{TMD}_w^L}(\mathcal{G}; X) = \frac{1}{|X|} \sum_{i \in [n]} \min_{j \in \mathcal{G}} \text{TMD}_w^L(G_i, G_j).$$

For each $j \in \mathcal{G}$, let τ_j be the number of graphs in X closest to G_j (i.e., the medoid’s cluster size): $\tau_j := |\{i \in [n] : \text{TMD}_w^L(G_i, G_j) < \text{TMD}_w^L(G_i, G_k), \forall k \neq j \in \mathcal{G}\}|$, breaking ties arbitrarily.

The following lemma shows how we can bound the difference in a GNN’s loss on the subsampled dataset $\mathcal{G} \subset X$ and the original dataset X in terms of this medoids objective.

Lemma 5.1. Let \mathcal{H} be a hypothesis class of $(L - 1)$ -layer GNNs $h : \mathcal{G} \rightarrow \mathbb{R}^d$ with ℓ -th layer Lipschitz constant upper-bounded by Φ_ℓ . Let $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$ be an M -Lipschitz loss

function. Let $y = (y_1, \dots, y_n)$ be labels for X and \mathcal{G} be a subset of $[n]$. Then for any GNN $h \in \mathcal{H}$ and $G \in \mathcal{G}$,

$$\frac{1}{n} \left| \sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(h(G_i); y_i) - \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) \right| \leq \quad (4)$$

$$M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) f_{\text{TMD}_w^L}(\mathcal{G}; X). \quad (5)$$

Proof sketch. We use the Lipschitz constant M of \mathcal{L} and Theorem 4.1 to bound the average deviation in the loss on X and the (weighted) loss on \mathcal{G} . ■

Lemma 5.1 implies that if the right-hand-side of Equation (5) is small and we minimize loss over the subsampled training set (the left-most term in Equation (5)), we will obtain a hypothesis with nearly optimal loss over the entire training set. We summarize this observation as follows.

Corollary 5.2. *Suppose that $M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \leq c$ and $f_{\text{TMD}_w^L}(\mathcal{G}; X) \leq \varepsilon$. Let $\hat{h} \in \mathcal{H}$ be the hypothesis that minimizes loss on the subsampled graphs $\sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(h(G_i); y_i)$. Then*

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) + 2c\varepsilon.$$

Corollary 5.2 illustrates that the additional loss incurred by training on the subsampled graphs as opposed to X is bounded by an error proportional to $f_{\text{TMD}_w^L}(\mathcal{G}; X)$. This motivates us to formulate the graph subsampling problem as the problem of computing \mathcal{G} that minimizes $f_{\text{TMD}_w^L}(\mathcal{G}; X)$.

Definition 5.3 (Graph subsampling problem). Let $X = \{G_1, \dots, G_n\}$ be a set of graphs and k be a graph budget. In the graph subsampling problem, we must select a subset $\mathcal{G} \subset X$ of k graphs that minimizes the k -medoids objective value of $f_{\text{TMD}_w^L}(\mathcal{G}; X)$.

The k -medoids clustering problem is NP-hard (Kazakovtsev & Rozhnov, 2020), but there are efficient approximation algorithms. We use Python’s k -medoids implementation in the sklearn package (Buitinck et al., 2013).

6. Experiments

We evaluate our sampling methods on four graph benchmark datasets from the TUDataset graph benchmark library (Morris et al., 2020). All GNNs were trained for 50 epochs with an 80%-train set and 20%-test split on CPUs. For MUTAG, COX2, and BZX we used a 3-layer Graph Isomorphism Network (Xu et al., 2019). For IMDB-BINARY we used a 3-layer Graph Convolutional Neural Network (GCN) (Kipf & Welling, 2016). Hyper-parameter details are in

Appendix A.4. In all figures, the line plots and shadows indicate the means and 95% confidence intervals across 256 random trials. Figure 4 compares our TMD-based k -medoids graph subsampling procedure from Section 5 against several alternative benchmarks and against training on the full dataset.

We are unaware of any previous work that studies graph subsampling, but we nonetheless compare against two clustering approaches based on Herding (Welling, 2009) and k -Centers (Farahani & Hekmatfar, 2009). These approaches are typically used for node subsampling and coresets selection in other domains (see Section 2). We adapt them to graph subsampling by defining distances between graphs using three popular graph kernels: Weisfeiler Lehman kernel (WL) (Shervashidze et al., 2011), the WL optimal assignment kernel (WL-OT) (Kriege et al., 2016), graphlet sampling (GS) kernel (Borgwardt & Kriegel, 2005), and the shortest paths (SP) kernel. Figure 4 shows that our method consistently outperforms uniform random sampling, Herding (with respect to all kernels), and k -Centers (with respect to all kernels) on all four benchmarks. Our approach achieves over 90% of the test accuracy of the full training set for MUTAG, COX2, and BZR using just 1% of the graphs. Where labeling is costly, this could be a powerful approach.

Node subsampling. Figure 5 compares our TMD-based node subsampling procedure from Section 4 against several alternative benchmarks and against training on the full dataset. For node subsampling, we only need to compute tree norms, so we use Algorithm 1. We compare against model-agnostic node selection policies Herding (Welling, 2009) and k -Centers (Farahani & Hekmatfar, 2009), and approaches by Salha et al. (2022) and Razin et al. (2023). While other approaches exhibit significant variability across datasets and subsampling levels, our approach consistently has the strongest or near-strongest performance across all datasets, sometimes requiring as little as 10% of the original number of nodes to achieve similar test accuracy.

7. Discussion

We proposed new approaches for compressing graph datasets that preserve the performance of a downstream GNN. Our methods apply to both node subsampling and graph subsampling. By leveraging the TMD, our methods ensure that the subsampled data maintains structural similarities with the original dataset, and we could bound the loss incurred by subsampling. Our methods are enabled by an algorithm we developed for selecting subgraphs with minimal TMD from a given graph, providing significant runtime speed-up compared to prior work. Our experiments demonstrated that our method outperforms comparable methods.

385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439

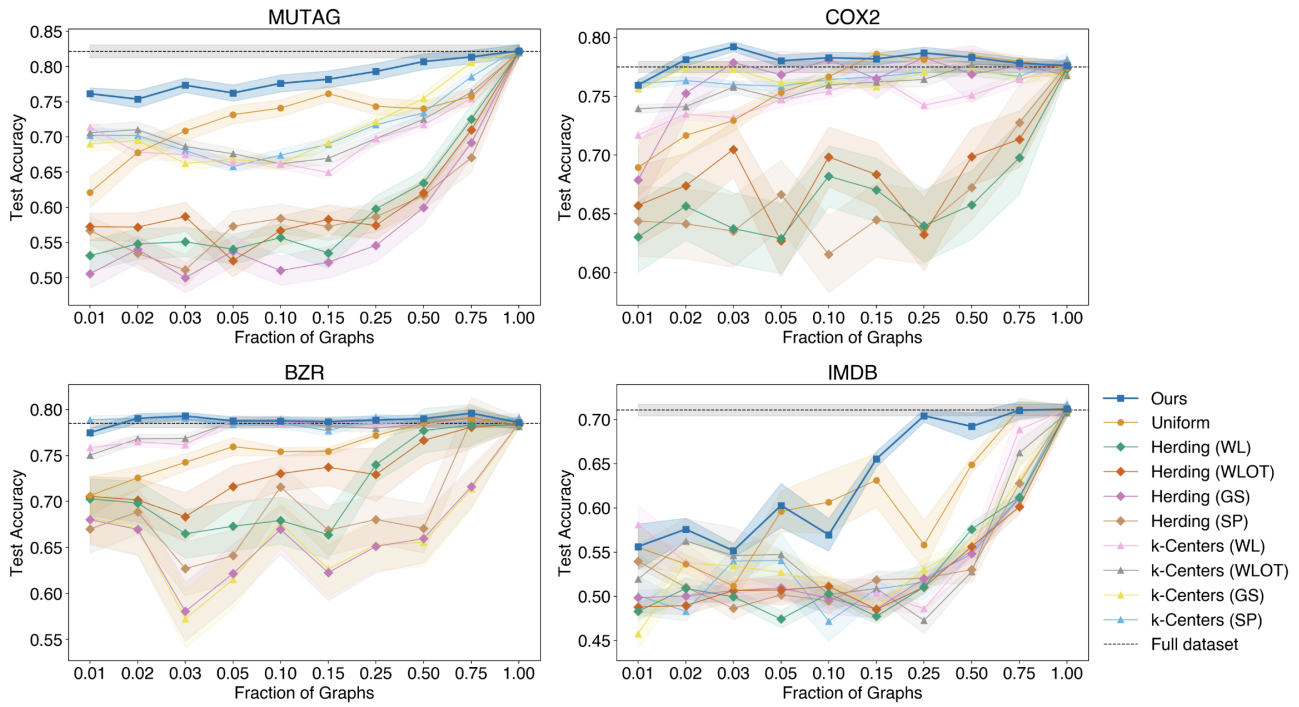


Figure 4. Graph subsampling. Accuracy vs. fraction of graphs in the training set, subsampled with our approach and alternative methods.

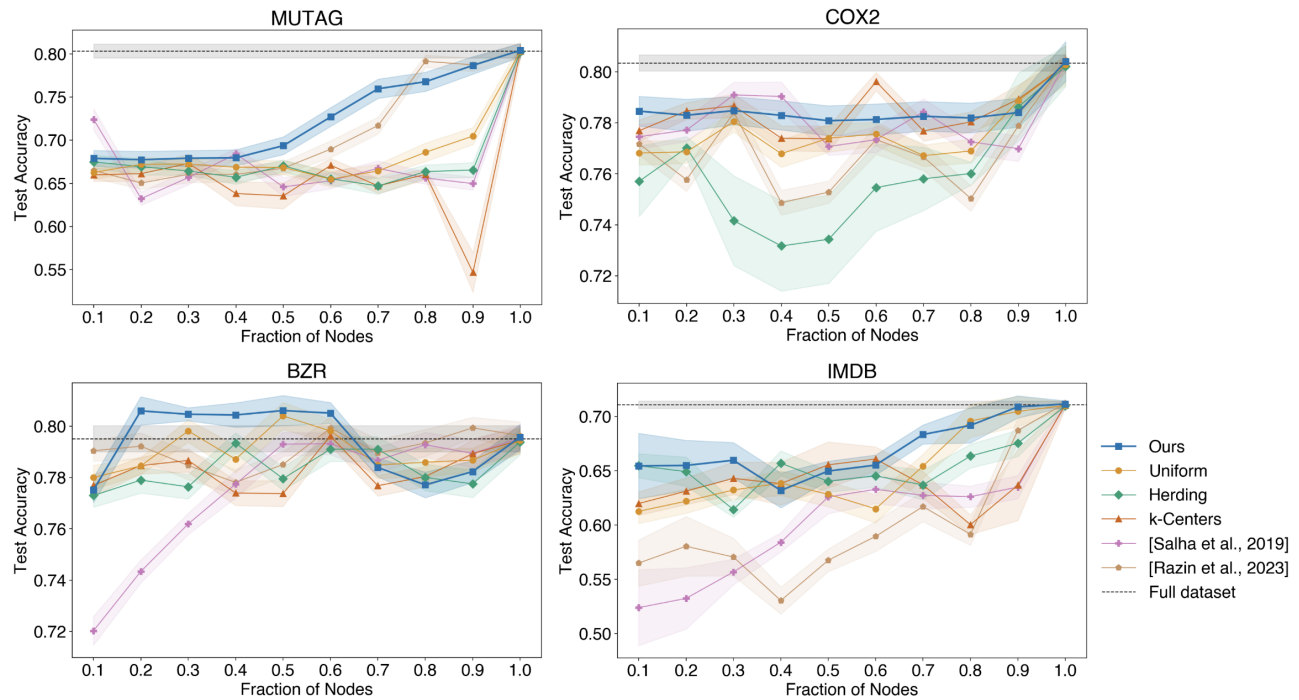


Figure 5. Node subsampling. Accuracy vs. fraction of nodes in the training set, subsampled with our approach and existing model-agnostic methods.

References

Alimohammadi, Y., Ruiz, L., and Saberi, A. A local graph limits perspective on sampling-based gnns. *arXiv preprint*

arXiv:2310.10953, 2023.

- 440 Bengio, Y., Lodi, A., and Prouvost, A. Machine learning
441 for combinatorial optimization: a methodological tour
442 d’horizon. *European Journal of Operational Research*,
443 290(2):405–421, 2021.
- 444 Bollobás, B. and Riordan, O. Coupling scale-free and classical
445 random graphs. *Internet Mathematics*, 1(2):215–225,
446 2004.
- 447
448 Bongini, P., Pancino, N., Scarselli, F., and Bianchini, M.
449 Biognn: how graph neural networks can solve biological
450 problems. In *Artificial Intelligence and Machine Learning
451 for Healthcare*, pp. 211–231. Springer, 2022.
- 452
453 Borgwardt, K. M. and Kriegel, H.-P. Shortest-path kernels
454 on graphs. pp. 8–pp. IEEE, 2005.
- 455 Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F.,
456 Mueller, A., Grisel, O., Niculae, V., Prettenhofer, P.,
457 Gramfort, A., Grobler, J., Layton, R., VanderPlas, J.,
458 Joly, A., Holt, B., and Varoquaux, G. API design for machine
459 learning software: experiences from the scikit-learn
460 project. In *ECML PKDD Workshop: Languages for Data
461 Mining and Machine Learning*, pp. 108–122, 2013.
- 462
463 Cappart, Q., Chételat, D., Khalil, E. B., Lodi, A., Morris, C.,
464 and Veličković, P. Combinatorial optimization and reasoning
465 with graph neural networks. *Journal of Machine
466 Learning Research*, 24(130):1–61, 2023.
- 467
468 Chuang, C.-Y. and Jegelka, S. Tree mover’s distance: Bridging
469 graph metrics and stability of graph neural networks.
470 In *Conference on Neural Information Processing Systems
471 (NeurIPS)*, volume 35, pp. 2944–2957, 2022.
- 472
473 Da San Martino, G., Navarin, N., and Sperduti, A. A memory
474 efficient graph kernel. In *International Joint Conference
475 on Neural Networks (IJCNN)*. IEEE, 2012.
- 476
477 Fan, W., Ma, Y., Li, Q., He, Y., Zhao, E., Tang, J., and
478 Yin, D. Graph neural networks for social recommendation.
479 In *Proceedings of the International World Wide Web
480 Conference (WWW)*, pp. 417–426, 2019.
- 481
482 Fang, J., Li, X., Sui, Y., Gao, Y., Zhang, G., Wang, K.,
483 Wang, X., and He, X. Exgc: Bridging efficiency and
484 explainability in graph condensation. In *Proceedings of
485 the International World Wide Web Conference (WWW)*,
486 2024.
- 487
488 Farahani, R. Z. and Hekmatfar, M. *Facility location: concepts,
489 models, algorithms and case studies*. Springer
490 Science & Business Media, 2009.
- 491
492 Fey, M. and Lenssen, J. E. Fast graph representation learning
493 with pytorch geometric. *arXiv preprint arXiv:1903.02428*,
494 2019.
- 495
496 Gasteiger, J., Becker, F., and Günnemann, S. Gemnet: Uni-
497 versal directional graph neural networks for molecules.
498 In *Conference on Neural Information Processing Systems
499 (NeurIPS)*, 2021.
- 500
501 Georgiev, D., Lio, P., Bachurski, J., Chen, J., Shi, T., and
502 Giusti, L. Beyond Erdos-Renyi: Generalization in algorithmic
503 reasoning on graphs. In *Learning on Graphs
504 (LoG)*, 2022.
- 505
506 Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation
507 learning on large graphs. 30, 2017.
- 508
509 Hamilton, W. L. *Graph representation learning*. Morgan &
510 Claypool Publishers, 2020.
- 511
512 Hashemi, M., Gong, S., Ni, J., Fan, W., Prakash, B. A.,
513 and Jin, W. A comprehensive survey on graph reduction:
514 Sparsification, coarsening, and condensation. *arXiv
515 preprint arXiv:2402.03358*, 2024.
- 516
517 He, Y., Zhang, X., Huang, J., Rozemberczki, B., Cucuringu,
518 M., and Reinert, G. Pytorch geometric signed directed:
519 A software package on graph neural networks for signed
520 and directed graphs. In *Learning on Graphs Conference*,
521 pp. 12–1. PMLR, 2024.
- 522
523 Huang, K., Jiang, H., Wang, M., Xiao, G., Wipf, D., Song,
524 X., Gan, Q., Huang, Z., Zhai, J., and Zhang, Z. Freshgcn:
525 Reducing memory access via stable historical embeddings
526 for graph neural network training. 17(6):1473–
527 1486, 2024.
- 528
529 Igarashi, Y., Kojima, R., Matsumoto, S., Iwata, H., Okuno,
530 Y., and Yamada, H. Developing a gnn-based ai model to
531 predict mitochondrial toxicity using the bagging method.
532 49(3):117–126, 2024.
- 533
534 Jiang, D., Wu, Z., Hsieh, C.-Y., Chen, G., Liao, B., Wang,
535 Z., Shen, C., Cao, D., Wu, J., and Hou, T. Could graph
536 neural networks learn better molecular representation for
537 drug discovery? a comparison study of descriptor-based
538 and graph-based models. 13:1–23, 2021.
- 539
540 Jin, W., Tang, X., Jiang, H., Li, Z., Zhang, D., Tang, J.,
541 and Yin, B. Condensing graphs via one-step gradient
542 matching. In *International Conference on Knowledge
543 Discovery and Data Mining (KDD)*, pp. 720–730, 2022.
- 544
545 KAWANO, K., KOIDE, S., SHIOKAWA, H., and AMAGASA,
546 T. Multi-dimensional fused gromov wasserstein
547 discrepancy for edge-attributed graphs. pp. 683–693,
548 2024.
- 549
550 Kazakovtsev, L. A. and Rozhnov, I. Application of algorithms
551 with variable greedy heuristics for k-medoids
552 problems. *Informatica*, 44(1), 2020.
- 553
554 Kipf, T. N. and Welling, M. Semi-supervised classification
555 with graph convolutional networks. 2016.

- 495 Kriege, N. M., Giscard, P.-L., and Wilson, R. On valid
496 optimal assignment kernels and applications to graph
497 classification. 29, 2016.
- 498
499 Le, T. and Jegelka, S. Limits, approximation and size trans-
500 ferability for gnns on sparse graphs via graphops. In
501 *Conference on Neural Information Processing Systems*
502 (*NeurIPS*), 2023.
- 503
504 Levie, R., Huang, W., Bucci, L., Bronstein, M. M., and Ku-
505 tyzniok, G. Transferability of spectral graph convolutional
506 neural networks. *Journal of Machine Learning Research*,
507 2021.
- 508
509 Li, Z., Meidani, K., Yadav, P., and Farimani, A. B. Graph
510 neural networks accelerated molecular dynamics. *J.*
511 *Chem. Phys.*, 156, 2022.
- 512
513 Morris, C., Kriege, N. M., Bause, F., Kersting, K., Mutzel,
514 P., and Neumann, M. Tudataset: A collection of bench-
515 mark datasets for learning with graphs. *arXiv preprint*
516 *arXiv:2007.08663*, 2020.
- 517
518 Newman, M. E. and Watts, D. J. Scaling and percolation in
519 the small-world network model. *Physical review E*, 60
520 (6):7332, 1999.
- 521
522 Park, C. W., Kornbluth, M., Vandermause, J., Wolverson,
523 C., Kozinsky, B., and Mailoa, J. P. Accurate and scalable
524 graph neural network force field and molecular dynam-
525 ics with direct force architecture. *npj Computational*
526 *Materials*, 7(73), 2021.
- 527
528 Peng, Y., Choi, B., and Xu, J. Graph learning for combi-
529 natorial optimization: a survey of state-of-the-art. pp.
530 119–141, 2021.
- 531
532 Razin, N., Verbin, T., and Cohen, N. On the ability of graph
533 neural networks to model interactions between vertices.
534 In *Conference on Neural Information Processing Systems*
535 (*NeurIPS*), 2023.
- 536
537 Ruiz, L., Gama, F., and Ribeiro, A. Graph neural networks:
538 Architectures, stability, and transferability. *Proceedings*
539 *of the IEEE*, 109(5):660–682, 2021.
- 540
541 Salha, G., Hennequin, R., Tran, V. A., and Vazirgiannis, M.
542 A degeneracy framework for scalable graph autoencoders,
543 2022.
- 544
545 Schuchardt, J., Scholten, Y., and Günnemann, S. (Prov-
546 able) adversarial robustness for group equivariant tasks:
547 Graphs, point clouds, molecules, and more. 2024.
- 548
549 Shen, X., Liu, Y., Wu, Y., and Xie, L. Molgnn: Self-
supervised motif learning graph neural network for drug
discovery. In *Machine Learning for Molecules Workshop*
at *NeurIPS*, pp. 4, 2020.
- Shervashidze, N., Schweitzer, P., Van Leeuwen, E. J.,
Mehlhorn, K., and Borgwardt, K. M. Weisfeiler-lehman
graph kernels. 12(9), 2011.
- Siglidis, G., Nikolentzos, G., Limnios, S., Giatsidis, C.,
Skianis, K., and Vazirgiannis, M. Grakel: A graph kernel
library in python. *Journal of Machine Learning Research*,
21(54):1–5, 2020.
- Sun, Z., Yin, H., Chen, H., Chen, T., Cui, L., and Yang,
F. Disease prediction via graph neural networks. 25(3):
818–826, 2020.
- Tanaka, Y., Eldar, Y. C., Ortega, A., and Cheung, G. Sam-
pling signals on graphs: From theory to applications.
IEEE Signal Processing Magazine, 37(6):14–30, 2020.
- Van Der Hofstad, R. *Random graphs and complex networks*.
Cambridge University Press, 2024.
- Vesselinova, N., Steinert, R., Perez-Ramirez, D. F., and Bo-
man, M. Learning combinatorial optimization on graphs:
A survey with applications to networking. *IEEE Access*,
8:120388–120416, 2020.
- Welling, M. Herding dynamical weights to learn. In *Pro-
ceedings of the 26th annual international conference on*
machine learning, pp. 1121–1128, 2009.
- Wu, B., Li, J., Yu, J., Bian, Y., Zhang, H., Chen, C., Hou, C.,
Fu, G., Chen, L., Xu, T., et al. A survey of trustworthy
graph learning: Reliability, explainability, and privacy
protection. *arXiv preprint arXiv:2205.10014*, 2022.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How pow-
erful are graph neural networks? In *Proceedings of the*
International Conference on Learning Representations
(ICLR), 2019.
- Yan, B., Wang, C., Guo, G., and Lou, Y. Tinygcn: Learning
efficient graph neural networks. In *International Confer-
ence on Knowledge Discovery and Data Mining (KDD)*,
pp. 1848–1856, 2020.
- Yehudai, G., Fetaya, E., Meirom, E., Chechik, G., and
Maron, H. From local structures to size generalization in
graph neural networks. In *International Conference on*
Machine Learning (ICML), 2021.
- Yu, J., Liang, J., and He, R. Mind the label shift of
augmentation-based graph ood generalization. In *Con-
ference on Computer Vision and Pattern Recognition*
(CVPR), pp. 11620–11630, 2023.
- Zeng, H., Zhou, H., Srivastava, A., Kannan, R., and
Prasanna, V. Accurate, efficient and scalable training
of graph neural networks. *Journal of Parallel and Dis-
tributed Computing*, 147:166–183, 2021.

550 Zhang, M. and Chen, Y. Link prediction based on graph
551 neural networks. 31, 2018.

552 Zhang, T., Zhang, Y., Wang, K., Wang, K., Yang, B., Zhang,
553 K., Shao, W., Liu, P., Zhou, J. T., and You, Y. Two
554 trades is not baffled: Condense graph via crafting rational
555 gradient matching. *arXiv preprint arXiv:2402.04924*,
556 2024.

557
558 Zheng, S., Zhu, Z., Liu, Z., Guo, Z., Liu, Y., Yang, Y.,
559 and Zhao, Y. Multi-modal graph learning for disease
560 prediction. 41(9):2207–2216, 2022.

561
562 Zhou, J., Cui, G., Hu, S., Zhang, Z., Yang, C., Liu, Z., Wang,
563 L., Li, C., and Sun, M. Graph neural networks: A review
564 of methods and applications. 1, 2020.

565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604

605 A. Additional background

606 A.1. Random graph limits

607 As we discussed in Section 4.3, most real networks are scale-free and small-world, and hence well-modeled by random
608 graph models, such as preferential attachment, configuration models and inhomogenous random graphs. As shown by (Van
609 Der Hofstad, 2024), such random graph models produce graphs that can be shown to converge *locally* to a limit (G, o) ,
610 where G is a graph to which we assign a root node o . By a “transitive” argument, it makes sense to see real networks as
611 parts of sequences converging to graph limits in a similar sense.

612 To be precise, let \mathcal{G}_* be the set of all possible rooted graphs. A limit graph is defined as a measure over the space \mathcal{G}_* with
613 respect to the local metric

$$614 d_{loc}((G_1, o_1), (G_2, o_2)) = \frac{1}{1 + \inf_k \{k : B_k(G_1, o_1) \not\cong B_k(G_2, o_2)\}}$$

615 where $B_k(G, v)$ is the k -hop neighborhood of node v , and \cong is the graph isomorphism. A sequence of graphs converging to
616 this limit is defined as follows.

617 **Definition A.1** (Local convergence (Alimohammadi et al., 2023)). Let $G_n = (V_n, E_n)$ denote a finite connected graph. Let
618 (G_n, o_n) be the rooted graph obtained by letting $o_n \in V_n$ be chosen uniformly at random. We say that (G_n, o_n) converges
619 locally to the connected rooted graph (G, o) , which is a (possibly random) element of \mathcal{G}_* having law μ , when, for every
620 bounded and continuous function $h : \mathcal{G}_* \rightarrow \mathbb{R}$,

$$621 \mathbb{E}[h(G_n, o_n)] \rightarrow \mathbb{E}_\mu(G, o)$$

622 where the expectation on the right-hand-side is with respect to (G, o) having law μ , while the expectation on the left-hand-side
623 is with respect to the random vertex o_n .

624 A.2. Other graph kernels

625 TMD is one recently proposed pseudometric on graphs, or graph kernel (Chuang & Jegelka, 2022). It is appealing in
626 our setting because it characterizes the robustness of GNNs (see Theorem 4.1.) Other popular graph kernels include the
627 Weisfeiler Lehman (WL) kernel (Shervashidze et al., 2011), Weisfeiler Lehman Optimal Transport (WL-OT) (Kriege et al.,
628 2016), Shortest Paths kernel ((Borgwardt & Kriegel, 2005)), and FGW (KAWANO et al., 2024). To our knowledge, these
629 kernels do not have comparable robustness guarantees to TMD.

630 A.3. Additional details about the TMD

631 We use T_0^n to denote n disjoint copies of T_0 . Given two tree multisets $\mathcal{J}_u(T), \mathcal{J}_v(T')$, we define ρ to be the following
632 augmentation function, which returns two multi-sets of the same size:

$$633 \rho : (\mathcal{J}_v(T'), \mathcal{J}_u(T)) \mapsto \left(\mathcal{J}_v(T') \cup T_0^{\max(|\mathcal{J}_u(T)| - |\mathcal{J}_v(T')|, 0)}, \mathcal{J}_u(T) \cup T_0^{\max(|\mathcal{J}_v(T)| - |\mathcal{J}_u(T')|, 0)} \right).$$

634 A.4. Additional details about experiments

635 All experiments were run on a 16-core machine with 64 GB of RAM. The GNNs were implemented using Pytorch Geometric
636 (He et al., 2024). The TMD weight function was set according to Pascal’s triangle rule (Chuang & Jegelka, 2022, Theorem
637 8) and our implementation of TMD was based on (Chuang & Jegelka, 2022).² We used the Graph Kernel Library (Siglidis
638 et al., 2020) for the kernel distances in our experiments. We will be making the code for all of our experiments publicly
639 available.

640 For MUTAG, COX2, and BZX we used a GIN model with the following hyperparameters:

- 641 • Number of layers: 3
- 642 • Learning rate: 0.01
- 643 • Batch size: 128

644 ²One consideration with is the runtime required to compute the k -medoids. However, reduced labeling cost, reduced data storage, and
645 privacy are often more significant priorities.

-
- 660 • Number of hidden channels per layer: 32
 - 661 • Aggregation function: global add pool
 - 662 • Weight regularization: N/A
 - 664 • Optimizer: Adam
 - 665 • Loss: BCE with logits
 - 667 • Dropout: N/A

669 For IMDB we used a GCN model with the following hyperparameters:

- 670 • Number of layers: 3
- 672 • Learning rate: 0.01
- 673 • Batch size: 32
- 675 • Number of hidden channels per layer: 32
- 676 • Aggregation function: global add pool
- 678 • Weight regularization: N/A
- 679 • Optimizer: Adam
- 681 • Loss: BCE with logits
- 682 • Dropout: N/A

684 We refrained from using batch normalization in both models. Models were implemented using PytorchGeometric (Fey & Lenssen, 2019).

687 B. Appendix

688 Throughout this section, we use $\deg_G(v)$ to denote the degree of $v \in V$.

691 B.1. Omitted proofs from Section 4

692 **Corollary 4.2.** *Let \mathcal{H} be a hypothesis class of $(L - 1)$ -layer GNNs $h : \mathcal{G} \rightarrow \mathbb{R}^d$ with ℓ -th layer Lipschitz constant upper-bounded by Φ_ℓ . Let $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$ be an M -Lipschitz loss function. Let $X = (G_1, \dots, G_n)$ and (y_1, \dots, y_n) be a set of (training) graphs and their labels. For each $i \in [n]$, let $G_i = (V_i, E_i)$ and $S_i \subset V_i$. Suppose that $M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \leq c$ and $\text{TMD}_w^L(G_i, G_i[S_i]) \leq \varepsilon_i$ for all $i \in [n]$. Finally, let $\hat{h} \in \mathcal{H}$ be a hypothesis with minimum loss over the subsampled training set:*

$$697 \hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i \in [n]} \mathcal{L}(h(G_i[S_i]); y_i).$$

700 Then \hat{h} has nearly optimal loss over the original training set as well:

$$702 \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) \\ 704 + \frac{2c}{n} \sum_{i \in [n]} \varepsilon_i.$$

709 *Proof.* Consider any $h \in \mathcal{H}$. By Theorem 4.1, we have that for each $i \in [n]$,

$$711 \|h(G_i; y_i) - h(G_i[S_i]; y_i)\| \leq \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \text{TMD}_w^L(G_i, G_i[S_i]). \quad (6)$$

Consequently, using the fact that \mathcal{L} is M -Lipschitz and (6),

$$\begin{aligned} |\mathcal{L}(h(G_i); y_i) - \mathcal{L}(h(G_i[S_i])); y_i)| &\leq M \|h(G_i; y_i) - h(G_i[S_i]; y_i)\| \\ &\leq M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \text{TMD}_w^L(G_i, G_i[S_i]) \\ &= c \text{TMD}_w^L(G_i, G_i[S_i]) \leq c \varepsilon_i. \end{aligned}$$

Consequently, by the triangle inequality,

$$\begin{aligned} \left| \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) - \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i[S_i])); y_i) \right| &\leq \frac{1}{n} \sum_{i \in [n]} |\mathcal{L}(h(G_i); y_i) - \mathcal{L}(h(G_i[S_i])); y_i)| \\ &\leq \frac{c}{n} \sum_{i \in [n]} \varepsilon_i. \end{aligned}$$

So, we know that for any $h \in \mathcal{H}$

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) \leq \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i[S_i]); y_i) + \frac{c}{n} \sum_{i \in [n]} \varepsilon_i. \quad (7)$$

Let h^* be the hypothesis that minimizes average loss across the original dataset X :

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i).$$

By definition of \hat{h} , this means that

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i[S_i]); y_i) \leq \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h^*(G_i[S_i]); y_i) + \frac{c}{n} \sum_{i \in [n]} \varepsilon_i.$$

Applying (7) again to h^* , we have that

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i[S_i]); y_i) \leq \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h^*(G_i); y_i) + \frac{c}{n} \sum_{i \in [n]} \varepsilon_i.$$

■

Lemma 4.1. *Let $w : \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ be a weight function and $G = (V, E, f)$ be a graph. Then, the identity transportation plan \mathbf{I} that maps $T_v(G)$ to $T_v(G[S])$ for $v \in S$, and $T_u(G)$ to \emptyset for $u \notin S$ is an optimal transport plan for the OT problem*

$$\text{TMD}_w^L(G, G[S]) = \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{J}_G^L, \mathcal{J}_{G[S]}^L \right) \right). \quad (3)$$

Consequently, we can decompose the right-hand-side of (3) into three terms: a sum over the feature norms of deleted nodes $v \notin S$; the cost of transporting the deleted nodes' computation trees (that is, $\{\mathcal{J}_v(T_v^L(G))\}_{v \notin S}$) to the empty set; and the cost of transporting the computation trees of the remaining nodes $v \in S$ (that is, $\{\mathcal{J}_v(T_v^L(G))\}_{v \in S}$) to their computation trees in $G[S]$. Formally,

$$\begin{aligned} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{J}_G^L, \mathcal{J}_{G[S]}^L \right) \right) &= \sum_{v \notin S} \|f^v\| \\ &+ w(L-1) \sum_{v \notin S} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{J}_v(T_v^L(G)), \emptyset \right) \right) \\ &+ w(L-1) \sum_{v \in S} \text{OT}_{\text{TD}_w} \left(\rho \left(\mathcal{J}_v(T_v^L(G)), \mathcal{J}_v(T_v^L(G[S])) \right) \right). \end{aligned}$$

770 *Proof.* Let $\bar{G}[S]$ be the graph obtained by augmenting $G[S]$ with node features $\mathbf{0}$ for each $v \in V \setminus S$. That is, $\bar{G}[S] =$
771 $(V, E[S], f')$ where $f'^v = f^v$ if $v \in S$ and $\mathbf{0}$ otherwise. By the definition of the augmentation function ρ , we can instead
772 consider the following OT problem, which is equivalent to (3):

$$773 \text{OT}_{\text{TD}_w} \left(\rho(\mathcal{T}_G^L, \mathcal{T}_{\bar{G}[S]}^L) \right).$$

774
775
776
777
778 In the base case where $L = 1$, we see that *any* transportation plan between the nodes of $\bar{G}[S]$ and G must transport the
779 excess node feature mass $\sum_{v \notin S} \|f^v\|$ in G to some nodes in $G[S]$. The transportation plan \mathbf{I} has cost exactly equal to
780 $\sum_{v \in V} \|f^v - f'^v\| = \sum_{v \notin S} \|f^v\|$, so \mathbf{I} must be an optimal transportation plan.

781
782 Now, if $L > 1$, consider the OT problem 3, and let C and Γ be the distance matrix and feasible transportation plans for this
783 OT problem. Let $\gamma \in \Gamma$ be any feasible transportation plan.

784
785 By the definition of TD, for any $i, j \in V$, transporting $\gamma_{i,j}$ mass from i to j incurs two costs:

- 786 (1) $\gamma_{i,j}$ incurs the cost of transporting f^i to f'^j .
- 787 (2) $\gamma_{i,j}$ incurs the cost of taking all of the depth $(L - 1)$ computation trees of all neighbors of i in $T_i^L(G)$, and transporting
788 them to the depth $(L - 1)$ computation trees of all neighbors of j in $T_j^L(G[S])$ (after appropriately augmenting with
789 isolated nodes of feature $\mathbf{0}$ for each $v \notin S$).

790 To quantify these two costs, let us define $\bar{\mathcal{T}}_j(T_j^L(G[S]))$ to be the multiset of computation trees obtained by augmenting
791 $\mathcal{T}_j(T_j^L(G[S]))$ with isolated nodes of feature vectors $\mathbf{0}$ for each $j \notin S$. That is, let $\bar{\mathcal{T}}_j(T_j^L(G[S]))$ be defined as the second
792 output of ρ when applied to $(\mathcal{T}_j(T_j^L(G)), \mathcal{T}_j(T_j^L(G[S])))$:

$$793 (\mathcal{T}_j(T_j^L(G)), \bar{\mathcal{T}}_j(T_j^L(G[S]))) = \rho(\mathcal{T}_j(T_j^L(G)), \mathcal{T}_j(T_j^L(G[S]))) .$$

794
795
796
797 We can now quantify the cost of transportation plan γ (3) as follows. In the following equation, the first term corresponds to
798 item (1) above and the second term corresponds to item (2) discussed above.

$$799 \sum_{i,j \in V} \gamma_{i,j} C_{i,j} = \sum_{i,j \in V} \gamma_{i,j} \|f^i - f'^j\| + w(L - 1) \sum_{i,j \in V} \gamma_{i,j} \text{OT}_{\text{TD}_w}(\mathcal{T}_i(T_i^L(G)), \bar{\mathcal{T}}_j(T_j^L(G[S]))) .$$

800
801
802
803 Now, we can immediately lower bound the cost $\sum_{i,j} \gamma_{i,j} C_{i,j}$ of γ as follows:

$$804 \sum_{i,j \in V} \gamma_{i,j} C_{i,j} \geq \min_{\nu \in \Gamma} \sum_{i,j \in V} \nu_{i,j} \|f^i - f'^j\|$$

$$805 + w(L - 1) \min_{\tau \in \Gamma} \sum_{i,j \in V} \tau_{i,j} \text{OT}_{\text{TD}_w}(\mathcal{T}_i(T_i^L(G)), \bar{\mathcal{T}}_j(T_j^L(G[S]))) .$$

806
807
808 For the first term, it is easy to see that $\nu = \mathbf{I}$ is an optimal solution, by the same argument as in the base case: any
809 transportation plan ν must incur the cost of transporting the excess mass $\{\|f^v\|\}_{v \notin S}$ in the node features of G which is not
810 present in the node features of S .

811 For the second term, we can notice that $\text{OT}_{\text{TD}_w}(\mathcal{T}_i(T_i^L(G)), \bar{\mathcal{T}}_j(T_j^L(G[S])))$ is also an TMD between graphs which contain
812 the corresponding multisets of depth $L - 1$ trees! So, by the inductive hypothesis, $\tau = \mathbf{I}$ is an optimal solution for this OT
813 problem as well.

Consequently, substituting $\tau = \nu = \mathbf{I}$, we can further simplify the lower bound to

$$\begin{aligned}
\sum_{i,j \in V} \gamma_{i,j} C_{i,j} &\geq \sum_{i,j \in V} \mathbf{I}_{i,j} \|f^i - f^j\| + w(L-1) \sum_{i,j \in V} \mathbf{I}_{i,j} \text{OT}_{\text{TD}_w}(\mathcal{J}_i(T_i^L(G)), \bar{\mathcal{J}}_j(T_j^L(G[S]))) \\
&= \sum_{i \in V} \|f^i - f^j\| + w(L-1) \sum_{i \in V} \text{OT}_{\text{TD}_w}(\mathcal{J}_i(T_i^L(G)), \bar{\mathcal{J}}_i(T_i^L(G[S]))) \\
&= \sum_{v \notin S} \|f^v\| + w(L-1) \sum_{v \notin S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_v(T_v^L(G)), \emptyset)) \\
&\quad + w(L-1) \sum_{v \in S} \text{OT}_{\text{TD}_w}(\mathcal{J}_v(T_v^L(G)), \bar{\mathcal{J}}_v(T_v^L(G[S]))) \\
&= \sum_{v \notin S} \|f^v\| + w(L-1) \sum_{v \notin S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_v(T_v^L(G)), \emptyset)) \\
&\quad + w(L-1) \sum_{v \in S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_v(T_v^L(G)), \mathcal{J}_v(T_v^L(G[S]))) .
\end{aligned}$$

Finally, the inequality is tight, because $\gamma = \mathbf{I} \in \Gamma$ is clearly a feasible transportation plan. So, by induction, the lemma holds. \blacksquare

Remark B.1. The terms $\text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_v(T_v^L(G)), \mathcal{J}_v(T_v^L(G[S])))$ in the expression of Lemma 4.1 can themselves be viewed as the TMD between two graphs \bar{G} and \bar{G}' , where \bar{G} is the graph consisting of all elements in $\mathcal{J}_v(T_v^L(G))$ as a disjoint trees, and \bar{G}' is the graph consisting of all elements in $\mathcal{J}_v(T_v^L(G[S]))$ as disjoint trees. In this sense, Lemma 4.1 precisely characterizes the recursive structure of TMD between a graph and its subgraph. We leverage this fact in later proofs.

Lemma 4.2. For $T \subset S \subset V$, $\text{TMD}_w^L(G, G[S \setminus T]) = \text{TMD}_w^L(G, G[S]) + \text{TMD}_w^L(G[S], G[S \setminus T])$.

Proof. Let $Z = \{z \in V : z \notin S\}$. Since $T \subset V$, T and Z are disjoint. So, Lemma 4.1 shows that

$$\begin{aligned}
\text{TMD}_w^L(G, G[S \setminus T]) &= \sum_{z \in Z} \|f^z\| + w(L-1) \sum_{z \in Z} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_z(T_z^L(G)), \emptyset)) \\
&\quad + \sum_{t \in T} \|f^t\| + w(L-1) \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_t(T_t^L(G)), \emptyset)) \\
&\quad + w(L-1) \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G)), \mathcal{J}_u(T_u^L(G[S \setminus T]))) .
\end{aligned}$$

Now, each term in $\text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G)), \mathcal{J}_u(T_u^L(G[S \setminus T])))$ inside the summation is, in fact a depth $L-1$ TMD between two graphs corresponding to the two disjoint forests of computation trees (as we discussed in Remark B.1). Thus, by the inductive hypothesis, we can expand the last term as follows:

$$\begin{aligned}
\sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G)), \mathcal{J}_u(T_u^L(G[S \setminus T]))) &= \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G)), \mathcal{J}_u(T_u^L(G[S]))) \\
&\quad + \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G[S])), \mathcal{J}_u(T_u^L(G[S \setminus T]))) .
\end{aligned}$$

Consequently, we obtain:

$$\begin{aligned}
\text{TMD}_w^L(G, G[S \setminus T]) &= \sum_{z \in Z} \|f^z\| + w(L-1) \sum_{z \in Z} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_z(T_z^L(G)), \emptyset)) \\
&\quad + \sum_{t \in T} \|f^t\| + w(L-1) \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_t(T_t^L(G)), \emptyset)) \\
&\quad + w(L-1) \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G)), \mathcal{J}_u(T_u^L(G[S]))) \\
&\quad + w(L-1) \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{J}_u(T_u^L(G[S])), \mathcal{J}_u(T_u^L(G[S \setminus T]))) .
\end{aligned}$$

Next, we can add and subtract terms as follows.

$$\begin{aligned}
\text{TMD}_w^L(G, G[S \setminus T]) &= \sum_{z \in Z} \|f^z\| + w(L-1) \sum_{z \in Z} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_z(T_z^L(G)), \emptyset)) \\
&+ \sum_{t \in T} \|f^t\| + w(L-1) \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G)), \emptyset)) \\
&+ w(L-1) \sum_{u \in S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathcal{F}_u(T_u^L(G[S]))) \\
&- w(L-1) \sum_{u \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathcal{F}_u(T_u^L(G[S]))) \\
&+ w(L-1) \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G[S])), \mathcal{F}_u(T_u^L(G[S \setminus T])))).
\end{aligned}$$

Now, $\text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathcal{F}_u(T_u^L(G[V \setminus S])))$ is also a depth $L-1$ TMD between two graphs corresponding to two disjoint forests of computation trees (again, see Remark B.1). So, we can apply the inductive hypothesis to see that for each $t \in T$,

$$\begin{aligned}
\text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G)), \emptyset)) &= \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G)), \mathcal{F}_t(T_t^L(G[S]))) \\
&+ \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G[S])), \emptyset)).
\end{aligned}$$

Thus, summing over $t \in T$, we have

$$\begin{aligned}
\sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G[S])), \emptyset)) &= \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G)), \emptyset)) \\
&- \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G)), \mathcal{F}_t(T_t^L(G[S]))) .
\end{aligned}$$

By substituting this into the equation above, we obtain

$$\begin{aligned}
\text{TMD}_w^L(G, G[S \setminus T]) &= \sum_{z \in Z} \|f^z\| + w(L-1) \sum_{z \in Z} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_z(T_z^L(G)), \emptyset)) \\
&+ \sum_{t \in T} \|f^t\| + w(L-1) \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G[S])), \emptyset)) \\
&+ w(L-1) \sum_{u \in S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathcal{F}_u(T_u^L(G[S]))) \\
&+ w(L-1) \sum_{u \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G[S])), \mathcal{F}_u(T_u^L(G[S \setminus T])))).
\end{aligned}$$

Rearranging in the form of Lemma 4.1, we obtain

$$\begin{aligned}
\text{TMD}_w^L(G, G[S \setminus T]) &= \sum_{u \notin S} \|f^u\| + w(L-1) \sum_{u \notin S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathbf{0})) \\
&+ w(L-1) \sum_{u \in S} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_u(T_u^L(G)), \mathcal{F}_u(T_u^L(G[S]))) \\
&+ \sum_{t \in T} \|f^t\| + w(L-1) \sum_{t \in T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G[S])), \mathbf{0})) \\
&+ w(L-1) \sum_{t \in S \setminus T} \text{OT}_{\text{TD}_w}(\rho(\mathcal{F}_t(T_t^L(G[S])), \mathcal{F}_t(T_t^L(G[S \setminus T]))) \\
&= \text{TMD}_w^L(G, G[S]) + \text{TMD}_w^L(G[S], G[S \setminus T]).
\end{aligned}$$

■

Theorem 4.5. Given a graph $G = (V, E, f)$, Algorithm 1 computes $\|G\|_{\text{TN}_w^L}$ in $O(|E|L)$ time.

Proof. Consider Algorithm 1, TreeNorm. The runtime is immediate from the algorithm pseudocode, as matrix-vector products can be computed in $O(|E|)$. In this proof, let A_G denote the adjacency matrix of a graph G and let x_G be the vector in \mathbb{R}^V such that $x_G(v) = \|f^v\|$.

To verify the correctness, we proceed by induction. We just need to show that for all L , $\|G\|_{\text{TN}_w^L} = \|x_G\|_1 + \left\| \sum_{\ell=1}^{L-1} \left(\prod_{t=1}^{\ell} w(L-t) \right) A_G^{\ell} x_G \right\|_1$, as this is the exact computation computed by the algorithm pseudocode. Note that an equivalent expression is $\|G\|_{\text{TN}_w^L} = \left\| \sum_{\ell=0}^{L-1} \left(\prod_{t=1}^{\ell} w(L-t) \right) A_G^{\ell} x_G \right\|_1$, where the product over the emptyset is defined as 1.

If $L = 1$, then $\|G\|_{\text{TN}_w^L}$ is simply the sum of its feature values, so this is trivially true.

Consider the case of a depth- L tree-norm computation. Throughout this proof, all variables are positive, and consequently we can move the ℓ_1 norm in and out of sums freely (i.e., we have a triangle *equality*.)

By taking $S = \emptyset$ in Lemma 4.1, we see that $\|G\|_{\text{TN}_w^L}$ is equal to the norm of its nodes plus $w(L-1) \sum_{v \in V} \|\mathcal{J}_v(T_v^L(G))\|_{\text{TN}_w^{L-1}}$. That is,

$$\|G\|_{\text{TN}_w^L} = \|x_G\|_1 + w(L-1) \sum_{v \in V} \|\mathcal{J}_v(T_v^L(G))\|_{\text{TN}_w^{L-1}}.$$

To interpret this expression, let G'_v be the graph which is a forest of all the trees in $\mathcal{J}_v(T_v^L(G))$ —for each $u \in N_G(v)$, we G'_v will contain a depth- $(L-1)$ tree $G'_{v,u}$. Then, by inductive hypothesis,

$$\begin{aligned} \|\mathcal{J}_v(T_v^L(G))\|_{\text{TN}_w^{L-1}} &= \sum_{u \in N_G(v)} \|(T_v^L(G))\|_{\text{TN}_w^{L-1}} = \|G'_v\|_{\text{TN}_w^{L-1}} = \sum_{u \in N_G(v)} \|G'_{u,v}\|_{\text{TN}_w^{L-1}} \\ &= \sum_{u \in N_G(v)} \left\| \sum_{\ell=0}^{L-2} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) A_{G'_{u,v}}^{\ell} x_{G'_{u,v}} \right\|_1 \\ &= \sum_{u \in N_G(v)} \sum_{\ell=0}^{L-2} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) \|A_{G'_{u,v}}^{\ell} x_{G'_{u,v}}\|_1, \end{aligned}$$

Letting $A_G(v, u)$ be the (v, u) -th entry of A_G ,

$$\begin{aligned} \|\mathcal{J}_v(T_v^L(G))\|_{\text{TN}_w^{L-1}} &= \sum_{u \in N_G(v)} \sum_{\ell=0}^{L-2} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) \|A_{G'_{u,v}}^{\ell} x_{G'_{u,v}}\|_1 \\ &= \left\| \sum_{u \in V} \sum_{\ell=0}^{L-2} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) A_G(u, v) A_{G'_{u,v}}^{\ell} x_{G'_{u,v}} \right\|_1. \end{aligned}$$

Now, by the way we constructed the trees $\mathcal{J}_v(T_v^L(G))$, u has the same $(L-1)$ -hop neighborhood in G as it does in $G'_{u,v}$. Consequently,

$$\begin{aligned} \|G\|_{\text{TN}_w^L} &= \|x_G\|_1 + w(L-1) \sum_{v \in V} \|\mathcal{J}_v(T_v^L(G))\|_{\text{TN}_w^{L-1}} \\ &= \|x_G\|_1 + w(L-1) \left\| \sum_{\ell=0}^{L-2} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) \sum_{v \in V} \sum_{u \in V} A_G(u, v) A_{G'_{u,v}}^{\ell} x_{G'_{u,v}} \right\|_1 \\ &= \|x_G\|_1 + \left\| \sum_{\ell=1}^{L-1} \left(\prod_{t=1}^{\ell} w(L-t-1) \right) \sum_{v \in V} A_G^{\ell+1} x_G \right\|_1 \\ &= \|x_G\|_1 + \left\| \sum_{\ell=1}^{L-1} \left(\prod_{t=1}^{\ell} w(L-t) \right) A_G^{\ell} x_G \right\|_1, \end{aligned}$$

completing the induction and the proof of correctness. ■

B.2. Omitted proofs from Section 5

Lemma 5.1. *Let \mathcal{H} be a hypothesis class of $(L - 1)$ -layer GNNs $h : \mathcal{G} \rightarrow \mathbb{R}^d$ with ℓ -th layer Lipschitz constant upper-bounded by Φ_ℓ . Let $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$ be an M -Lipschitz loss function. Let $y = (y_1, \dots, y_n)$ be labels for X and \mathcal{G} be a subset of $[n]$. Then for any GNN $h \in \mathcal{H}$ and $G \in \mathcal{G}$,*

$$\frac{1}{n} \left| \sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(h(G_i); y_i) - \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) \right| \leq \quad (4)$$

$$M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) f_{\text{TMD}_w^L}(\mathcal{G}; X). \quad (5)$$

Proof. Let $\kappa : [n] \rightarrow \mathcal{G}$ map $i \in [n]$ to the index $j \in \mathcal{G}$ such that G_j is the closest graph in $\{G_i\}_{i \in \mathcal{G}}$ to G_i . That is,

$$\kappa(i) = \operatorname{argmin}_{j \in \mathcal{G}} \text{TMD}_w^L(G_i, G_j),$$

with ties broken arbitrarily (but consistently with the mapping τ .) First, we can rearrange the sums as follows

$$\begin{aligned} \left| \frac{1}{n} \sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(h(G_i); y_i) - \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) \right| &= \left| \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_{\kappa(i)}); y_{\kappa(i)}) - \mathcal{L}(h(G_i); y_i) \right| \\ &\leq \frac{1}{n} \sum_{i \in [n]} |\mathcal{L}(h(G_{\kappa(i)}); y_{\kappa(i)}) - \mathcal{L}(h(G_i); y_i)| \\ &\leq \frac{M}{n} \sum_{i \in [n]} \|\mathcal{L}(h(G_{\kappa(i)}); y_{\kappa(i)}) - \mathcal{L}(h(G_i); y_i)\|, \end{aligned}$$

where in the last inequality, we used that \mathcal{L} is M -lipschitz. Now, by Theorem 4.1, we have

$$\begin{aligned} \frac{M}{n} \sum_{i \in [n]} \|\mathcal{L}(h(G_{\kappa(i)}); y_{\kappa(i)}) - \mathcal{L}(h(G_i); y_i)\| &\leq \frac{M}{n} \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \sum_{i \in [n]} \text{TMD}_w^L(G_{\kappa(i)}, G_i) \\ &= M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) f_{\text{TMD}_w^L}(\mathcal{G}; X). \end{aligned}$$

Corollary 5.2. *Suppose that $M \left(\prod_{\ell \in [L-1]} \Phi_\ell \right) \leq c$ and $f_{\text{TMD}_w^L}(\mathcal{G}; X) \leq \varepsilon$. Let $\hat{h} \in \mathcal{H}$ be the hypothesis that minimizes loss on the subsampled graphs $\sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(h(G_i); y_i)$. Then*

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i) + 2c\varepsilon.$$

Proof. From Lemma 5.1, we know that

$$\frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \frac{1}{n} \sum_{i \in \mathcal{G}} \tau_i \mathcal{L}(\hat{h}(G_i); y_i) + c\varepsilon.$$

1045 Let h^* be the hypothesis that minimizes average loss over the original dataset X :

$$1046 \quad h^* = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h(G_i); y_i).$$

1047
1048
1049
1050 By definition of \hat{h} , this means that

$$1051 \quad \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \frac{1}{n} \sum_{i \in S} \tau_i \mathcal{L}(h^*(G_i); y_i) + c\varepsilon.$$

1052
1053
1054
1055 Applying Lemma 5.1 again, we have that

$$1056 \quad \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(\hat{h}(G_i); y_i) \leq \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(h^*(G_i); y_i) + 2c\varepsilon.$$

1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

■