# A ROBUST STACKING FRAMEWORK FOR DEEP GRAPH NETWORKS WITH MULTIFACETED NODE FEATURES

## **Anonymous authors**

Paper under double-blind review

# Abstract

Graph Neural Networks (GNNs) with numerical node features and graph structure as inputs have demonstrated superior performance on various semi-supervised learning tasks with graph data. However, the numerical node features utilized by GNNs are commonly extracted from raw data which is of text or tabular (numeric/categorical) type in most real-world applications. The best models for such data types in most standard supervised learning settings with IID (non-graph) data are not simple neural network layers and thus are not easily incorporated into a GNN. Here we propose a robust stacking framework that fuses graph-aware propagation with arbitrary models intended for IID data, which are ensembled and stacked in multiple layers. Our layer-wise framework leverages bagging and stacking strategies to enjoy strong generalization, in a manner which effectively mitigates label leakage and overfitting. Across a variety of graph datasets with tabular/text node features, our method achieves comparable or superior performance relative to both tabular/text and graph neural network models, as well as existing state-of-the-art hybrid strategies that combine the two.

## **1** INTRODUCTION

Graph datasets comprise nodes of various data types and modalities linked by edges that encapsulate non-IID conditional dependencies between them. While it is often assumed that graph neural networks (GNN) (Kipf & Welling, 2016; Veličković et al., 2017) are preferable for handling such data relative to models originally designed for IID instances, GNNs are nonetheless subject to various limitations. In particular, the best architecture may be data-set specific and require appropriately setting many attendant structural hyperparameters, e.g., note the complex assortment of GNN architectures that populate the top of the Open Graph Benchmark (OGB) leaderboard (Hu et al., 2020). Moreover, most GNNs implicitly assume that node features are numerical, and may struggle to remain competitive with more complex text, tabular, or composite alternatives.

In fact, with richer node feature sets it has even been observed that models tailored to IID data (which in our setting simply operate on individual node features as though they were independent of the others) can at times outperform GNNs if they are combined with simple graph propagation operations to account for the graph structure (Huang et al., 2020; Chen et al., 2021). Moreover, for graph data with text features, Chien et al. (2021) has demonstrated that leveraging a BERT Transformer in addition to a GNN can greatly improve performance. And beyond these considerations, real-world applications of ML typically involve more than just a single model, GNN or otherwise. Instead they usually require an ML pipeline composed of data preprocessing and training/tuning/aggregation of many models to achieve the best results.

In this paper, we investigate how to adapt ML pipelines designed for supervised learning with IID data (e.g., Transformers for text, gradient boosted decision trees or related for tabular data) to node classification/regression tasks with graph-structured statistical dependencies between node features. We focus on using *k*-fold bagging (Breiman, 1996), i.e. *cross-validation*, to avoid label leakage issues, with stack ensembling methods for maximal flexibility (Wolpert, 1992; Van der Laan et al., 2007). These techniques are particularly effective for achieving high accuracy across diverse IID datasets, and are utilized in many popular AutoML frameworks (Erickson et al., 2020; LeDell & Poirier, 2020; Feurer et al., 2015), but have largely been ignored within the context of graph data.

Within this context, our goal is to design a single architecture that integrates graph propagation or message passing steps and stacked ensembles of arbitrary base models to flexibly accommodate diverse node/instance types within a unified framework. In doing so, our contributions are as follows:

- We propose a framework of stack ensembling with graph propagation called **BestowGNN** for Bagged, Ensembled, Stacked Training Of Well-balanced GNNs (see Figure 1) that can *bestow* arbitrary (non-graph) base models intended for IID data with the capability of producing highly accurate node predictions in the graph (i.e., non-IID) setting.
- Using only a single, unified architecture, our proposed methodology can match or outperform bespoke dataset-specific models that top competitive leaderboards for popular node classification/regression tasks (e.g., on OGB and elsewhere completely different network architectures typically dominate the top positions for different datasets and data types).
- Label leakage is an unavoidable issue for many layer-wise training strategies (SAGN (Sun & Wu, 2021) and GAMLP (Zhang et al., 2021)). To address this potential shortcoming, we formalize how our bagging and stacking framework can effectively mitigate the label leakage issue within the graph setting using analytical tools from differential privacy. This is the first work establishing that bagging with graph-based predictors can be useful for ameliorating label leakage.

# 2 RELATED WORK

# 2.1 FROM SCALABILITY TO LAYER-WISE TRAINING

Currently, GNN training suffers from high computational cost with the number of layers growing. To improve the scalability of GNNs, graph sampling scheme GraphSAGE (Hamilton et al., 2017) is adopted by uniformly sampling a fixed number of neighbours for a batch of nodes. Cluster-GCN (Chiang et al., 2019) uses graph clustering algorithms to sample a block of nodes that form a dense subgraph and runs SGD-based algorithms on these subgraphs.  $L^2$ -GCN (You et al., 2020) proposes a layer-wise training framework by disentangling feature aggregation and feature transformation to reduce time and memory complexity.

SAGN (Sun & Wu, 2021) iteratively trains models in several stages by applying graph structureaware attention mechanisms on node features and also combines the self-training approach with label propagation to further improve performance. GAMLP (Zhang et al., 2021) proposes two attention mechanisms to explore the relation between features with different propagation steps. Both SAGN and GAMLP achieve state-of-the-art performance on two large open graph benchmarks (ogbn-products and ogbn-papers100M), demonstrating the high scalability and efficiency of layer-wise training strategies. However, SAGN and GAMLP suffer from the risk of label leakage: label information is included in the enhanced training set, and can cause performance degredation if the model extracts and relies on these labels. SAGN empirically shows that enough propagation depth can effectively alleviate label leakage, thus they only use label information at one fixed propagation step. Meanwhile, GAMLP passes label information between propagation steps using residual connections. Wang et al. (2021) further randomly masks nodes during every training epoch to mitigate label leakage issue.

# 2.2 GRAPH MODELS WITH MULTIFACETED NODE FEATURES

Traditional GNN models are mostly studied for graphs with homogeneous sparse node features. Leading GNN models fail to achieve competitive results for heterogeneous features with tabular or text node features (Ivanov & Prokhorenkova, 2021; Huang et al., 2020; Chen et al., 2021). To remedy this, Ivanov & Prokhorenkova (2021) jointly train Gradient Boosted Decision Trees (GBDT) and GNN in an end-to-end fashion, demonstrating a significant increase in performance on graph data with tabular node features.

Chen et al. (2021) removes the need for a GNN altogether, proposing a generalized framework for iterating boosting with parameter-free graph propagation steps that share node/sample information across edges connecting related samples.

Correct and Smooth (C&S) (Huang et al., 2020) is a simple post-processing step that applies label propagation to further incorporate graph information into the outputs of a learning algorithm. Chen

et al. (2021) trains Gradient Boosted Decision Trees with label propagation incorporated into the objective function, producing competitive results for graph data with tabular node features.

Because common GNNs take numerical node features as inputs, one must establish a way to extract numerical embeddings from raw data such as text and images. For example, the embeddings of ogbn-arxiv data are computed by running the skip-gram model (Mikolov et al., 2013). Chien et al. (2021) proposes self-supervised learning to fully utilizing correlations between graph nodes, and extracts the embedding of three open graph benchmark datasets (ogbn-arxiv, ogbn-products and ogbnpapers100M). Chien et al. (2021) demonstrates the superior performance of these new embeddings for the Open Graph Benchmark datasets. Lin et al. (2021) proposes BertGCN, which combines the Bert model and transductive learning for text classification in an end-to-end fashion and achieves superior performance on a range of text classification tasks.



not depicted here). The stacking layer repeats the operations depicted between it and the input data.

with  $\{\alpha_m\}$  fitted via Ensemble Selection

#### BACKGROUND 3

Consider an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with  $n = |\mathcal{V}|$  nodes. The node feature matrix is denoted by  $X \in \mathbb{R}^{n \times d}$ , and the corresponding node label matrix is  $Y \in \mathbb{R}^{n \times c}$  with d and c being the dimension of features and labels respectively. The unweighted adjacency matrix is  $A \in \mathbb{R}^{n \times n}$ . For training purposes we only have access to the labels of a subset of nodes  $\{y_i\}_{i \in \mathcal{L}}$ , with  $\mathcal{L} \subset \mathcal{V}$ . Given feature values of all nodes  $\{x_i\}_{i \in \mathcal{V}}$ , label data  $\{y_i\}_{i \in \mathcal{L}}$ , and the connectivity of the graph  $\mathcal{E}$ , the task is to

predict the labels of the unlabeled nodes  $\{y_i\}_{i \in \mathcal{U}}$ , with  $\mathcal{U} = \mathcal{V} \setminus \mathcal{L}$ . We denote the labeled dataset  $\{x_i, y_i\}_{i \in \mathcal{L}}$  as  $D_{\mathcal{L}}$  and the unlabeled dataset  $\{x_i\}_{i \in \mathcal{U}}$  as  $D_{\mathcal{U}}$ .

## 3.1 BAGGING, ENSEMBLING, AND STACKING

For classification/regression with IID (non-graph) data, bagging, ensembling, and stacking represent practical tools that can be combined in various ways to produce more accurate predictions relative to other strategies across diverse tabular and text datasets (Shi et al., 2021; Blohm et al., 2020; Yoo et al., 2020; Fakoor et al., 2020; Bezrukavnikov & Linder, 2021; Feldman, 2021). For example, in each stacking layer of an ensemble-based architecture, bagging simply trains the same types of base models with out-of-fold predictions from the previous layer models (obtained via bagging) as extra predictive features. These base models might include various Gradient Boosted Decision Trees (Ke et al., 2017; Prokhorenkova et al., 2018), fully-connected neural networks (MLP), K Nearest Neighbors (Erickson et al., 2020), or pretrained Electra Transformer models (Clark et al., 2020). For our purposes herein, we adopt the AutoML package AutoGluon (Erickson et al., 2020), which is capable of exploiting these techniques while serving open-source code that we can readily adapt to include graph propagation.

## 3.2 GRAPH-AWARE PROPAGATION LAYERS

Recently there has been a surge of interest GNN architectures with layers defined in one-to-one correspondence with descent iterations that minimize a principled class of graph-regularized energy functions (Klicpera et al., 2018; Ma et al., 2020; Pan et al., 2021; Yang et al., 2021; Zhang et al., 2020; Zhu et al., 2021). IN this way GNN models can benefit from the inductive bias afforded by energy function minimizers (or close approximations thereof) whose specific form can be controlled by trainable parameters.

Following Zhou et al. (2004), one relevant energy function capable of inducing such graph-aware propagation is given by

$$\ell_{Y}(\boldsymbol{Y}) \triangleq (1-\lambda) \|\boldsymbol{Y} - m(\boldsymbol{X};\boldsymbol{\theta})\|_{\mathcal{F}}^{2} + \lambda \operatorname{tr} \left[\boldsymbol{Y}^{\top} \boldsymbol{L} \boldsymbol{Y}\right],$$
(1)

where  $\lambda \in (0, 1)$  is a weight that determines the trade-off between the two terms.  $Y \in \mathbb{R}^{n \times d}$  is a learnable *d*-dimensional embedding across *n* nodes, and  $m(X; \theta)$  denotes a base model (parameterized by  $\theta$ ) that computes an initial target embedding based on the node features X.  $L \in \mathbb{R}^{n \times n}$  is the graph Laplacian of  $\mathcal{G}$ , meaning L = D - A, where D represents the degree matrix.

Intuitively, the first term of (1) encourages Y to be close to initial target embedding, while the second term introduces the smoothness over the whole graph. On the positive side, the closed-form optimal solution of energy function (1) can be easily derived as

$$\widetilde{m}^{*}(\boldsymbol{X};\boldsymbol{\theta}) \triangleq \arg\min_{\boldsymbol{Y}} \ell_{\boldsymbol{Y}}(\boldsymbol{Y}) = \boldsymbol{P}^{*}m(\boldsymbol{X};\boldsymbol{\theta}), \qquad (2)$$

with  $P^* \triangleq (I + \lambda L)^{-1}$ . However, for large graphs the requisite inverse is impractical to compute, and alternatively iterative approximations are more practically-feasible. To this end, we may initialize as  $Y^{(0)} = m(X; \theta)$ , and it follows that Y can be approximated by iterative descent in the direction of the negative gradient. Given that

$$\frac{\partial \ell_{Y}(\boldsymbol{Y})}{\partial \boldsymbol{Y}} = 2\lambda \boldsymbol{L}\boldsymbol{Y} + 2\boldsymbol{Y} - 2m\left(\boldsymbol{X};\boldsymbol{\theta}\right), \qquad (3)$$

the k-th iteration of gradient descent becomes

$$\boldsymbol{Y}^{(k)} = \boldsymbol{Y}^{(k-1)} - \alpha \left[ \left( \lambda \boldsymbol{L} + \boldsymbol{I} \right) \boldsymbol{Y}^{(k-1)} - m \left( \boldsymbol{X}; \boldsymbol{\theta} \right) \right], \tag{4}$$

where  $\frac{\alpha}{2}$  serves as the effective step size. Considering that L is generally sparse, computation of (4) can leverage efficient sparse matrix multiplications, and we may also introduce modifications such as Jacobi preconditioning to speed convergence (Axelsson, 1996; Yang et al., 2021).

Furthermore, based on well-known properties of gradient descent, if k is sufficiently large and  $\alpha$  is small enough, then

$$\widetilde{m}^{*}(\boldsymbol{X};\boldsymbol{\theta}) \approx \widetilde{m}^{(k)}(\boldsymbol{X};\boldsymbol{\theta}) \triangleq \boldsymbol{P}^{(k)}[m(\boldsymbol{X};\boldsymbol{\theta})],$$
(5)

where the operator  $P^{(k)}(\cdot)$  computes k gradient steps via (4). The structure of these propagation steps, as well as related variants based on normalized modifications of gradient descent, equate to principled GNN layers, such as those used by GCN (Kipf & Welling, 2016), APPNP (Klicpera et al., 2018), and many others, which can be trained within a broader bilevel optimization framework as described next.

# 4 STACK ENSEMBLING FOR GRAPH DATA (BESTOWGNN)

For node prediction tasks (either regression or classification), each (non-graph) base model is trained within our BestowGNN framework by simply treating each node and its label as a separate IID training example and fitting the model in the usual manner. Such a model may informatively encode tabular or text features from the nodes, but its predictions will be uniformed by the additional information available in the graph structure. To enhance such models with graph information we utilize graph-aware propagation.

# 4.1 GRAPH-AWARE PROPAGATION

Let  $\hat{Y}_{\mathcal{L}}$ ,  $\hat{Y}_{\mathcal{U}}$  denote the predictions of labeled (i.e. training) nodes and unlabeled (i.e. validation/test) nodes, respectively. In node classification tasks, these may be predicted class probability vectors. Via iterative application of the update in (4), we can apply graph-aware propagation to predictions  $\{\hat{Y}_{\mathcal{L}}, \hat{Y}_{\mathcal{U}}\}$  in order to ensure they reflect statistical dependencies between nodes encoded by the graph structure. We denote  $F^{(0)} \triangleq \{\hat{Y}_{\mathcal{L}}, \hat{Y}_{\mathcal{U}}\}$ , and for each propagation step t we compute the update  $F^{(t)}$  via (4). In our method,  $\hat{Y}$  may actually be predictions from multiple models concatenated together at each node, but the propagation procedure remains identical in this case.

# 4.2 STACK ENSEMBLING

In stack ensembling, the predictions output by individually trained *base* models are concatenated together as features that are subsequently used to train a *stacker* model whose target is still to predict the original labels (Wolpert, 1992; Ting & Witten, 1997). A good stacker model learns how to nonlinearly combine the predictions of base models into an even more accurate prediction. This process can be iterated in multiple layers, a strategy that has been used to win high-profile prediction competitions with IID data (Koren, 2009).

In this work, we closely follow the stacking methodology of Erickson et al. (2020), but adapt it for graphs rather than IID data. We allow stacker models to access the original node features X by concatenating X with the base models' predictions when forming the features used to train each stacker model. To produce a final prediction for each node, we aggregate predictions from the topmost layer models via a simple weighted combination where weights are learned via the efficient Ensemble Selection technique of Caruana et al. (2004). Our base models before the first stacking layer are those which can effectively encode the original tabular or text features observed at the nodes (here we utilize AutoGluon which leverages models like Gradient Boosted Decision Trees for tabular features and Transformers for text features). Our stacker models are simply chosen as the same types of models as the base models.

## 4.3 REPEATED K-FOLD BAGGING TO MITIGATE OVER-FITTING

A problem that arises in the aforementioned stacking strategy is *label leakage*. If a base model is even slightly overfit to its training data such that its predictions memorize parts of the training labels, then subsequent stacker models will have low accuracy due to distribution shift in their features between training and inference time (their features will be highly correlated with the labels during training but not necessarily during inference). This issue is remedied by ensuring stacker models are only trained on features comprised of base model predictions on held-out nodes omitted from the base model's training set.

We achieve this while still being able to train stacker models using all labeled nodes by leveraging k-fold bagging (i.e. cross-validation) of all models (Van der Laan et al., 2007; Parmanto et al., 1996; Erickson et al., 2020). Here the training nodes are partitioned into k disjoint chunks and k copies of

each (non-graph-aware) model m are trained with a different data chunk held-out  $\{X^{-j}, Y^{-j}\}_{j=1}^k$  held out from each copy. After training all k copies of model m, we can produce out-of-fold (OOF) predictions  $\hat{Y}_m^j$  for each chunk  $X^j$  by feeding it into the model copy from which it was previously held-out. Following Erickson et al. (2020), we repeat this k-fold bagging procedure over n different random partitions of the training data to further reduce variance and distribution shift that arises in stack ensembling with bagging. Thus for a labeled training node, the OOF prediction from a model of type m is averaged over n different partitions (this node is held-out from exactly one model copy in each partition):

$$\hat{\boldsymbol{Y}}_{\mathcal{L}} = \left\{ \frac{1}{n} \sum_{i=1}^{n} \hat{\boldsymbol{Y}}_{m,i}^{j} \right\}_{j=1}^{k}.$$
(6)

Since unlabeled (validation/test) nodes were technically held-out from every model copy, we can feed them through any copy without harming stacking performance. For a particular type of model m, we simply make predictions  $\hat{Y}_{\mathcal{U}}$  for unlabeled nodes by averaging over all n bagging repeats and all k copies of the model within each repeat:

$$\hat{\mathbf{Y}}_{\mathcal{U}} = \frac{1}{kn} \sum_{j=1}^{k} \sum_{i=1}^{n} \hat{\mathbf{Y}}_{m,i}^{j}.$$
(7)

For IID data, this stack ensembling procedure with bagging can produce powerful predictors, both in theory (Van der Laan et al., 2007) and in practice (Erickson et al., 2020).

## 4.4 STACKING WITH GRAPH-AWARE PROPAGATION

To extend this methodology to graph data, our proposed training strategy is precisely detailed in Algorithm 1. The main idea is to apply graph-aware propagation on the predictions of models at each intermediate layer of the stack. Different amounts of propagation lead to different characteristics of the data being captured in the resulting prediction (few steps of propagation means predictions are only influenced by local neighbors, whereas many propagation steps allow predictions to be influenced by more distant nodes as well). Thus we can further enrich the feature set of our stacker models by concatenating together the predictions produced after different numbers of propagation steps. With this expanded feature set, our stacker models learn to aggregate not only the predictions of different models, but differently smoothed versions of these predictions as well. This allows the stacker model to adaptively decide how to best account for dependencies induced by the graph structure.

More precisely, if we let  $\mathbf{F}^{(t)}$  denote the predictions (concatenated across all base model types) for labeled and unlabeled nodes after t smoothing steps, then the feature input to each stacker model is given by the original node features X concatenated with  $[\mathbf{F}^{(0)}, ..., \mathbf{F}^{(T)}]$ . Here the predictions for labeled nodes are always OOF, obtained via bagging. Another fundamental difference between our approach and stack ensembling in the IID setting is *the use of unlabeled (test) nodes at each intermediate layer of the stack.* By including unlabeled nodes in the propagation, these nodes influence the features used to train subsequent stacker models at labeled nodes. This can even further reduce potential distribution shift in the stacker models' features between the labeled and unlabeled nodes, which ensures better generalization.

Graph machine learning models for non-IID data typically do not use bagging, seemingly because there has not been a rigorous study on the effect of bagging in relation to propagation models. Furthermore, bagging traditionally serves as a means of variance reduction which only brings limited performance benefits for large datasets (Breiman, 1996). In contrast, our stacking framework adopts bagging primarily as a means to mitigate the catastrophic effects of label leakage. While bagging can effectively mitigate label information from being directly encoded in stacker model features in the IID setting, it is not clear whether this property still holds with graph-structured dependence between nodes. A particular concern is the fact that the propagation of base model predictions across the graph implies label information is shared across the k-fold chunks used to hold-out some nodes from some models. In the next section, we theoretically study this issue and prove that bagging can still mitigate the effects of label leakage even in the non-IID graph setting. Our subsequent experiments (see Table 4) reveal that bagging produces substantial performance gains in practical applications of stack ensembling with graph propagation.

# 5 THEORETICAL ANALYSIS

Label utilization is a common technique in which the outputs of a model are concatenated with input features and then used to train a stacking layer. Unfortunately, layer-wise training with label utilization is susceptible to the label leakage problem. Although prior work (Sun & Wu, 2021; Zhang et al., 2021) has mentioned heuristic ways to address label leakage via graph propagation, it is unclear how generally applicable this strategy is in practice. Moreover, there is a natural trade-off between avoiding label leakage via graph propagation, and well-known oversmoothing effects in GNN models.

In this section we employ a powerful theoretical tool, Differential Privacy (Mironov, 2017), to showcase the advantage of bagging in our proposed BestowGNN. Our analysis will show that BestowGNN enjoys strong generalization under the Rényi Differential Privacy framework. In fact this is the first work that establishes that bagging in graph predictors is useful and mitigates label leakage. Specifically, BestowGNN can preserve the privacy (or information sharing) of labels between bags, that would otherwise be compromised by graph propagation.

To this end, we first introduce the definition of Rényi Differential Privacy, which is a relaxation of Differential Privacy based on the Rényi Divergence.

**Definition 1.** (*Rényi Differential Privacy (Mironov, 2017)*). Consider a randomized algorithm  $\mathcal{M}$  mapping from  $\mathcal{D}$  to a real-value  $\mathcal{R}$ . Such an algorithm is said to have  $\epsilon$ -Rényi Differential Privacy of order  $\alpha$  if for any  $D, D' \in \mathcal{D}$  with  $d_H(D, D') = 1$ , where  $d_H$  is the Hamming distance (D, D') are also referred to as adjacent datasets), we have that

$$D_{\alpha}(\mathcal{M}(D)||\mathcal{M}(D')) \triangleq \frac{1}{\alpha - 1} \log E_{x \sim \mathcal{M}(D')} \left(\frac{\mathcal{M}(D)}{\mathcal{M}(D')}\right)^{\alpha} \le \epsilon.$$
(8)

In plain words, this definition establishes that the output of an algorithm does not change significantly, as measured by the Rényi divergence  $D_{\alpha}(\mathcal{M}(D)||\mathcal{M}(D'))$ , when the data changes slightly. The idea behind this framework is that if each individual data sample has only a small effect on the resulting model, the model cannot be used to infer information about any single individual.

We then have the following result:

**Theorem 1.** Assume base model m is a multi-layer (two-layer) perceptron and that node features X are sampled from a multivariate Gaussian as in (Jia & Benson, 2021):

$$\boldsymbol{X} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Gamma}^{-1}), \qquad \boldsymbol{\Gamma} = c_1 \boldsymbol{I}_n + c_2 \boldsymbol{L},$$

where  $I_n$  is an identity matrix and L is the normalized graph Laplacian. Here  $c_1$  controls a noise level and  $c_2$  the smoothness over the whole graph.  $E(x_0; D_{\mathcal{L}})$  and  $F(x_0; D_{\mathcal{L}})$  are predictions produced by BestowGNN for a data point  $x_0$  with and without bagging mode, respectively. If E has sensitivity 1 and lower magnitude bound L, i.e., for any two adjacent  $D, D' \in D$  :  $|E(x_0; D) - E(x_0; D')| \leq 1$ and  $|E| \geq L$ , then E satisfies  $(\frac{1}{2}, \frac{1}{4\sigma^2 L^2} + \frac{1}{2L^2})$ -Rényi Differential Privacy, where  $\sigma^2$  depends on graph structure  $\mathcal{G}$ . Meanwhile, F has no privacy guarantee, i.e., the Rényi differential privacy loss (8) is unbounded.

The proof is deferred to the supplementary. Theorem 1 indicates that bagging with graph propagation can well preserve the privacy of  $D_{\mathcal{L}} = \{x_i, y_i\}_{i \in \mathcal{L}}$  between different chunks while non-bagging would have a high risk of leaking the information of  $D_{\mathcal{L}}$ . For layer-wise training with label utilization, the output of the model  $E(x_0; D_{\mathcal{L}})$  is concatenated with input features and then used to train next stacking layer, and bagging can effectively mitigate the label leakage issue since the information of true label is well preserved at the first layer, while no-bagging exposes the true labels and can lead to over-fitting issue for next stacking layer.

## 6 EXPERIMENTS

**Setup.** We study the effectiveness of our approach by comparing performance against state-ofthe-art baselines in node regression and classification tasks. For node regression with **tabular node features**, we consider four real-world graph datasets used for benchmarking by Ivanov & Prokhorenkova (2021): House, County, VK and Avazu. As node classification tasks, we adopt one datasets with **numerical features**: Reddit; and two datasets with **raw text features**: OGB-Arxiv and OGB-Products. More details about the datasets are provided in the supplementary.

We compare our method against various baselines, starting with purely tabular baseline models or language models where the graph structure is ignored. Our first baseline is **Autogluon** (Erickson et al., 2020), an AutoML system for IID tabular or text data that is completely unaware of the graph structure (here we simply treat nodes as IID). Next, we consider **AutoGluon + C&S**, which performs Correct and Smooth (Huang et al., 2020) as a posthoc processing step on top of AutoGluon's predictions, in order to at least account for the graph structure during inference. For node regression tasks we also consider some popular GNN models: **GCN** (Kipf & Welling, 2016), **GAT** (Veličković et al., 2017), and a hybrid strategy **BGNN** (Ivanov & Prokhorenkova, 2021), which combines Gradient Boosted Decision Trees (also a model intended for IID data) with GNNs via end-to-end training in a manner that is graph-aware.

For node classification, we firstly consider Reddit with original numerical features. We compare with **GraphSAGE** (Hamilton et al., 2017) and **PCAPass + Tree** (Sadowski et al., 2022), which combines PCA and message passing to generate node embeddings and leverages tree-based model for node classification.

We also consider OGB-Arxiv and OGB-Products with raw text as node features (as opposed to precomputed text embeddings as node features such as the low-dimensional homogeneous embeddings provided by OGB). We compare with **GIANT-XRT + MLP**, **GIANT-XRT + GRAPHSAGE** and **GIANT-XRT + GRAPHSAINT**, which extracts numerical embeddings from text features via a transformer trained through self-supervised learning and feed these high quality embeddings to a multi-layer perceptron or sampling based GNN model. For the smaller OGB-Arxiv dataset, we also consider standard GNN models: **GCN** (Kipf & Welling, 2016), **GAT** (Veličković et al., 2017) and **Ensemble GCN**, a natural baseline/competitor which divides all training nodes into k chunks, trains a GCN model for each chunk and then ensembles the results. Finally, we compare against SOTA model for OGB-Arxiv and OGB-Products with GIANT-XRT embedding and low-dimensional homogeneous embedding from OGB leaderboard. To our knowledge, there is not a consistent method with superior performance across each dataset. So we compare our *single* general framework with *different* SOTA models for each dataset to ensure we are competing against the best in each case; i.e., *there is no single existing model that is SOTA across them all*. We evaluate our method **BestowGNN**, which incorporates the graph information through propagation operations in each stacking layer.

**Results.** In Table 1 we present the results for the node regression task with tabular node features. The baseline GNN models are challenged by the tabular node features. AutoGluon is an ensemble of various base models (e.g., Gradient Boosted Decision Trees, fully-connected neural networks) intended for IID data without considering graph structure. We observe that **Autogluon + C&S** outperforms **Autogluon**, demonstrating that graph information can greatly boost the performance of models intended for IID data. Incorporating the graph structure at each stacking layer, our **BestowGNN** method performs better than **BGNN** on all datasets.

Tables 2 and 3 show the results for node classification with either raw text features or numerical embeddings. Our method BestowGNN outperforms all baselines regardless of whethor or not they leverage the raw text or OGB embeddings (or numerical Reddit embeddings). Note that OGB-Arxiv and OGB-Products have *different* SOTA models in the OGB leaderboard, for instance: **AGDN + BoT + self-KD + C&S** are architectural components from the best existing model for OGB-Arxiv, while **GAMLP + RLU + SCR + C&S** undergird the best existing model for OGB-Products. These SOTA models consisting of data-specific modules/components are manually composed to perform particularly well only for one specific dataset. In contrast, **BestowGNN** uses essentially the same architecture with minor/standard hyperparameter tuning to fit all datasets. Comparison of **BestowGNN** with AutoGluon demonstrates how incorporating graph information at each stacking layer can further improve the node classification performance of this AutoML system. More experiments details and computing cost are deferred to the supplementary.

**Ablation.** The key ingredients of our framework are bagging/ensembling and graph propagation. Table 4 shows an ablation study involving these components using OGB-Arxiv and OGB-Products with original OGB embeddings. From these results we observe that bagging modes can outperform

Data set	House	County	Vk	Avazu		
GCN	$0.63 \pm 0.01$	$1.48 \pm 0.08$	$7.25 \pm 0.19$	$0.1141 \pm 0.02$	Method	Reddit
BGNN	$0.34 \pm 0.01$ $0.50 \pm 0.01$	$1.43 \pm 0.00$ $1.26 \pm 0.08$	$7.22 \pm 0.19$ $6.95 \pm 0.21$	$0.1134 \pm 0.01$ $0.109 \pm 0.01$	PCAPass + XGBoost GraphSAGE	$96.26 \pm 0.02$ $95.40 \pm 0.22$
AutoGluon AutoGluon + C&S	$\begin{array}{c} 0.618 \pm 0.01 \\ 0.477 \pm 0.01 \end{array}$	$\begin{array}{c} 1.379 \pm 0.08 \\ 1.162 \pm 0.09 \end{array}$	$\begin{array}{c} 7.176 \pm 0.21 \\ 6.995 \pm 0.21 \end{array}$	$\begin{array}{c} 0.117 \pm 0.018 \\ 0.107 \pm 0.015 \end{array}$	AutoGluon BestowGNN	$95.83 \pm 0.00$ 96.44 ±0.00
BestowGNN	$\textbf{0.467} \pm \textbf{0.007}$	$\textbf{1.145} \pm \textbf{0.083}$	$\textbf{6.918} \pm \textbf{0.220}$	$\textbf{0.105} \pm \textbf{0.013}$		

Table 1:	Mean squared	l error results	for four node
regression	datasets with	tabular node	features.

**Table 3:** Node classification accuracy for OGB-Arxiv and OGB-Products achieved by various methods. Rows labeled TEXT contain methods including SOTA models trained on the **raw text** features at each node, while those labeled OGB indicate models trained on precomputed **numerical embeddings** provided by OGB as node features. *SOTA models vary from each dataset with different embeddings/architectures, but BestowGNN has consistently superior performance for each dataset; similarly for Table 1 results above.* 

### **OGB-Arxiv**

#### **OGB-Products**

 
 Table 2: Node classification accuracy for Reddit with numerical node features.

Feature	Method	Test Acc (Validation)	Feature	Method	Test Acc (Validation)
OGB	GCN GAT + C&S AGDN+BoT+self-KD+C&S Ensemble GCN	$\begin{array}{c} 73.06 \pm 0.24 \ (74.42 \pm 0.12) \\ 73.86 \pm 0.14 \ (74.84 \pm 0.07) \\ 74.31 \pm 0.14 \ (75.18 \pm 0.09) \\ 73.22 \pm 0.12 \ (74.64 \pm 0.01) \end{array}$	OGB	DeeperGCN + FLAG GAT + FLAG GAMLP+RLU+SCR+C&S Ensemble GAT	$\begin{array}{c} 81.93 \pm 0.31 \ (92.21 \pm 0.37) \\ 81.76 \pm 0.45 \ (92.51 \pm 0.06) \\ 85.20 \pm 0.08 \ (93.04 \pm 0.05) \\ 80.01 \pm 0.20 \ (93.24 \pm 0.05) \end{array}$
TEXT	GIANT-XRT+MLP GIANT-XRT+graphSAGE GIANT-XRT+GCN GIANT-XRT+RevGAT+KD	$\begin{array}{c} 73.06 \pm 0.11 \ (74.32 \pm 0.09) \\ 74.35 \pm 0.14 \ (75.95 \pm 0.11) \\ 75.28 \pm 0.17 \ (76.87 \pm 0.04) \\ 76.15 \pm 0.10 \ (77.16 \pm 0.09) \end{array}$	TEXT	GIANT-XRT+MLP GIANT-XRT+graphSAGE GIANT-XRT+graphSAINT GIANT-XRT+SAGN+SLE	$\begin{array}{c} 80.49 \pm 0.28 \ (92.10 \pm 0.09) \\ 81.99 \pm 0.45 \ (93.38 \pm 0.05) \\ 84.15 \pm 0.22 \ (93.18 \pm 0.04) \\ 85.47 \pm 0.29 \ (\text{-}) \end{array}$
TEXT	AutoGluon AutoGluon + C&S	$\begin{array}{c} 73.05\pm 0.00 \; (74.33\pm 0.00) \\ 75.34\pm 0.00 \; (76.67\pm 0.00) \end{array}$	TEXT	AutoGluon AutoGluon + C&S	$\begin{array}{c} 77.10 \pm 0.06 \; (91.78 \pm 0.03) \\ 79.03 \pm 0.12 \; (93.62 \pm 0.03) \end{array}$
TEXT	BestowGNN	$76.19 \pm 0.02~(77.25 \pm 0.05)$	TEXT	BestowGNN	$85.48 \pm 0.03~(93.93 \pm 0.02)$

**Table 4:** BestowGNN ablation study with  $(\checkmark)$  and without bagging  $(\aleph)$ . Here T is the number of graph propagation steps, thus T = 0 represents a baseline model that completely ignores graph structure.

STEP $T$	AR	XIV	PRODUCTS	
	1	×	1	X
0	$55.70\pm0.33$	$54.14\pm0.29$	$62.28 \pm 0.35$	$62.05\pm0.19$
1	$66.25\pm0.27$	$64.57\pm0.76$	$74.18 \pm 0.21$	$72.61\pm0.66$
2	$69.34 \pm 0.16$	$67.37 \pm 0.44$	$77.07 \pm 0.32$	$74.61\pm0.58$
3	$70.01\pm0.16$	$68.08 \pm 0.74$	$78.11 \pm 0.19$	$75.79\pm0.49$
4	$70.43\pm0.21$	$68.72\pm0.63$	$78.76 \pm 0.60$	$76.86\pm0.17$

no-bagging modes for each number propagation step, demonstrating that bagging can effectively mitigate label leakage and over-fitting issues even in a graph-aware propagation setting.

## 7 DISCUSSION

While real-world graph data come with heterogeneous feature types, existing GNN models are primarily suited for (adequately preprocessed) numerical features. For IID supervised learning, it is well-known that the best models for different feature types vary based on dataset and data-type, and that a learning system aiming to output good predictions across a variety of datasets should leverage a heterogeneous collection of different types of models (Erickson et al., 2020). There is little reason the situation should be different for graph data. In this paper, we demonstrate the first working system that can utilize arbitrary heterogeneous collections of models for arbitrary graph datasets with heterogeneous feature-types (numerical, categorical, text). This is achieved by means of a novel graph-aware stack ensembling technique that takes the graph structure into account without restricting how individual models are trained. Our graph-aware propagation techniques leverage specific properties of stack ensembling that allow our proposed methodology to outperform both many complex GNNs as well as existing approaches in which propagation is only applied to the predictions output by an IID base model (e.g., AutoGluon+C&S, etc.).

## REFERENCES

Owe Axelsson. Iterative Solution Methods. Cambridge University Press, 1996.

- Oleg Bezrukavnikov and Rhema Linder. A neophyte with automl: Evaluating the promises of automatic machine learning tools. *arXiv preprint arXiv:2101.05840*, 2021.
- Matthias Blohm, Marc Hanussek, and Maximilien Kintz. Leveraging automated machine learning for text classification: Evaluation of automl tools and comparison with human performance. *arXiv* preprint arXiv:2012.03575, 2020.

Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.

- Rich Caruana, Alexandru Niculescu-Mizil, Geoff Crew, and Alex Ksikes. Ensemble selection from libraries of models. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 18, 2004.
- Jiuhai Chen, Jonas Mueller, Vassilis N Ioannidis, Soji Adeshina, Yangkun Wang, Tom Goldstein, and David Wipf. Convergent boosted smoothing for modeling graph data with tabular node features. *arXiv preprint arXiv:2110.13413*, 2021.
- Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of* the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 257–266, 2019.
- Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic, and Inderjit S Dhillon. Node feature extraction by self-supervised multi-scale neighborhood prediction. *arXiv preprint arXiv:2111.00064*, 2021.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*, 2020.
- Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv:2003.06505*, 2020.
- Rasool Fakoor, Jonas Mueller, Nick Erickson, Pratik Chaudhari, and Alexander J Smola. Fast, accurate, and simple models for tabular data via augmented distillation. In *Advances in Neural Information Processing Systems*, 2020.
- Sergey Feldman. Which machine learning classifiers are best for small datasets? An empirical study. https://www.data-cowboys.com/blog/, 2021.
- Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. *Advances in neural information processing systems*, 28, 2015.
- William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 1025–1035, 2017.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *arXiv* preprint arXiv:2005.00687, 2020.
- Qian Huang, Horace He, Abhay Singh, Ser-Nam Lim, and Austin R Benson. Combining label propagation and simple models out-performs graph neural networks. *arXiv preprint arXiv:2010.13993*, 2020.
- Sergei Ivanov and Liudmila Prokhorenkova. Boost then convolve: Gradient boosting meets graph neural networks. *arXiv preprint arXiv:2101.08543*, 2021.

- Junteng Jia and Austin R Benson. A unifying generative model for graph learning algorithms: Label propagation, graph convolutions, and combinations. *arXiv preprint arXiv:2101.07730*, 2021.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. LightGBM: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, 2017.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. *arXiv preprint arXiv:1810.05997*, 2018.
- Yehuda Koren. The bellkor solution to the netflix grand prize, 2009. URL https://www2.seas.gwu.edu/~simhaweb/champalg/cf/papers/KorenBellKor2009.pdf.
- Erin LeDell and Sebastien Poirier. H2o automl: Scalable automatic machine learning. In *Proceedings* of the AutoML Workshop at ICML, volume 2020, 2020.
- Yuxiao Lin, Yuxian Meng, Xiaofei Sun, Qinghong Han, Kun Kuang, Jiwei Li, and Fei Wu. Bertgcn: Transductive text classification by combining gcn and bert, 2021.
- Yao Ma, Xiaorui Liu, Tong Zhao, Yozen Liu, Jiliang Tang, and Neil Shah. A unified view on graph neural networks as graph signal denoising. *arXiv preprint arXiv:2010.01777*, 2020.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Ilya Mironov. Rényi differential privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF), pp. 263–275. IEEE, 2017.
- Xuran Pan, Shiji Song, and Gao Huang. A unified framework for convolution-based graph neural networks, 2021. URL https://openreview.net/forum?id=zUMD--Fb9Bt.
- Bambang Parmanto, Paul W Munro, and Howard R Doyle. Reducing variance of committee prediction with resampling techniques. *Connection Science*, 8(3-4):405–426, 1996.
- Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. CatBoost: unbiased boosting with categorical features. In *Advances in Neural Information Processing Systems*, 2018.
- Krzysztof Sadowski, Michał Szarmach, and Eddie Mattia. Dimensionality reduction meets message passing for graph node embeddings. *arXiv preprint arXiv:2202.00408*, 2022.
- Xingjian Shi, Jonas Mueller, Nick Erickson, Mu Li, and Alexander J Smola. Benchmarking multimodal automl for tabular data with text fields. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, 2021.
- Chuxiong Sun and Guoshi Wu. Scalable and adaptive graph neural networks with self-label-enhanced training. *arXiv preprint arXiv:2104.09376*, 2021.
- Kai Ming Ting and Ian H Witten. Stacking bagged and dagged models. In *International Conference* on Machine Learning, 1997.
- Mark J Van der Laan, Eric C Polley, and Alan E Hubbard. Super learner. *Statistical applications in genetics and molecular biology*, 6(1), 2007.
- Tim Van Erven and Peter Harremos. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820, 2014.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.

- Yangkun Wang, Jiarui Jin, Weinan Zhang, Yong Yu, Zheng Zhang, and David Wipf. Bag of tricks for node classification with graph neural networks. *arXiv preprint arXiv:2103.13355*, 2(3), 2021.
- David H Wolpert. Stacked generalization. Neural networks, 5(2):241-259, 1992.
- Yongyi Yang, Tang Liu, Yangkun Wang, Jinjing Zhou, Quan Gan, Zhewei Wei, Zheng Zhang, Zengfeng Huang, and David Wipf. Graph neural networks inspired by classical iterative algorithms. *arXiv preprint arXiv:2103.06064*, 2021.
- Jason Yoo, Tony Joseph, Dylan Yung, S Ali Nasseri, and Frank Wood. Ensemble squared: A meta automl system. *arXiv preprint arXiv:2012.05390*, 2020.
- Yuning You, Tianlong Chen, Zhangyang Wang, and Yang Shen. L2-gcn: Layer-wise and learned efficient training of graph convolutional networks. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 2127–2135, 2020.
- Hongwei Zhang, Tijin Yan, Zenjun Xie, Yuanqing Xia, and Yuan Zhang. Revisiting graph convolutional network on semi-supervised node classification from an optimization perspective. arXiv preprint arXiv:2009.11469, 2020.
- Wentao Zhang, Ziqi Yin, Zeang Sheng, Wen Ouyang, Xiaosen Li, Yangyu Tao, Zhi Yang, and Bin Cui. Graph attention multi-layer perceptron. *arXiv preprint arXiv:2108.10097*, 2021.
- Dengyong Zhou, Olivier Bousquet, Thomas Navin Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. *Advances in Neural Information Processing Systems*, 2004.
- Meiqi Zhu, Xiao Wang, Chuan Shi, Houye Ji, and Peng Cui. Interpreting and unifying graph neural networks with an optimization framework. *arXiv preprint arXiv:2101.11859*, 2021.

# **Supplementary Materials**

# A PROOF OF THEOREM 1.

**Preliminary 1.** Firstly, we derive the format of  $E(x_0; D_{\mathcal{L}})$  and  $F(x_0; D_{\mathcal{L}})$ . Suppose BestowGNN randomly splits the labeled nodes  $D_{\mathcal{L}}$  into 2 disjoint chunks  $D_1 = \{X_1, Y_1\}, D_2 = \{X_2, Y_2\}$ . BestowGNN trains a model  $m \in \mathcal{M}$  with a different data chunk held-out. Model m is defined by a set of parameters collected in  $\theta$  namely, which is defined as  $m(X; \theta)$ . In the following, we will express the predicted labels from model m under the bagging and non-bagging settings. We compare the predicted labels under both settings and establish that our bagging solution is less amenable to label leakage.

The model m will learn different parameters for each chunk and those are denoted as  $\theta_1$  for the chunk I and  $\theta_2$  for the chunk II, namely  $\theta_1 = \theta(D_1)$  and  $\theta_2 = \theta(D_2)$ . Next, BestowGNN produces prediction  $\hat{Y}_1, \hat{Y}_2$  on out-of-fold data, i.e.,  $\hat{Y}_1 = m(X_1; \theta_2)$  and  $\hat{Y}_2 = m(X_2; \theta_1)$ . The prediction for unlabeled nodes is  $\hat{Y}_{\mathcal{U}} = \frac{1}{2}[m(X_{\mathcal{U}}; \theta_1) + m(X_{\mathcal{U}}; \theta_2)]$  as explained in (7). Consider one data point  $x_0$  from the unlabeled dataset  $D_{\mathcal{U}}$ , the prediction of  $x_0$  is given by  $\hat{y}_0 = \frac{1}{2}[m(x_0; \theta_1) + m(x_0; \theta_2)]$ . Next, we perform one step graph-aware propagation on  $\hat{y}_0$ .

$$\hat{\boldsymbol{y}}_{0}^{(1)} = \sum_{u \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{\mathcal{U}}} \hat{\boldsymbol{y}}_{u} + \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} \hat{\boldsymbol{y}}_{v} + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} \hat{\boldsymbol{y}}_{w}$$

$$= \sum_{u \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{\mathcal{U}}} \frac{1}{2} [m(\boldsymbol{x}_{u}; \boldsymbol{\theta}_{1}) + m(\boldsymbol{x}_{u}; \boldsymbol{\theta}_{2})] + \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}; \boldsymbol{\theta}_{2}) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}; \boldsymbol{\theta}_{1})$$
(9)

where  $\hat{y}_0^{(1)}$  is the aggregated results from one-hop neighbor  $\mathcal{N}(\boldsymbol{x}_0)$ , which may belongs to  $D_{\mathcal{U}}, D_1$  and  $D_2$ .

Next, we consider the no-bagging mode, where the predictions of  $X_1, X_2$  are changed into  $Y_1 = m(X_1; \theta_1)$  and  $\tilde{Y}_2 = m(X_2; \theta_2)$ . Notice that with bagging mode we use the parameters from a different bag, while without bagging we use the parameters from the same bag. The prediction of the test point  $x_0$  is once again  $\tilde{y}_0 = \frac{1}{2}[m(x_0; \theta_1) + m(x_0; \theta_2)]$ , which is identical to the bagging mode. We perform the same graph-aware propagation on  $\tilde{y}_0$ .

$$\widetilde{\boldsymbol{y}}_{0}^{(1)} = \sum_{u \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{\mathcal{U}}} \widetilde{\boldsymbol{y}}_{u} + \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} \widetilde{\boldsymbol{y}}_{v} + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} \widetilde{\boldsymbol{y}}_{w}$$

$$= \sum_{u \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{\mathcal{U}}} \frac{1}{2} [m(\boldsymbol{x}_{u}; \boldsymbol{\theta}_{1}) + m(\boldsymbol{x}_{u}; \boldsymbol{\theta}_{2})] + \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}; \boldsymbol{\theta}_{1}) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}; \boldsymbol{\theta}_{2}).$$
(10)

Next, we compare the terms among the predicted labels from the two settings, namely (9) and (10). The first term  $\sum_{u \in \mathcal{N}(\boldsymbol{x}_0) \cap D_{\mathcal{U}}} \frac{1}{2} [m(\boldsymbol{x}_u; \boldsymbol{\theta}_1) + m(\boldsymbol{x}_u; \boldsymbol{\theta}_2)]$  is the same for (9) and (10) and can be cancelled. In order to facilitate the exposition of the theoretical contributions we will define functions for the different terms in (9) and (10). We define  $\boldsymbol{E}(\boldsymbol{x}_0; D_{\mathcal{L}})$ , that is a function formulating the relation between training data  $D_{\mathcal{L}}$  and the prediction for test data  $\boldsymbol{x}_0$  under bagging mode.

$$\boldsymbol{E}(\boldsymbol{x}_0; D_{\mathcal{L}}) := \sum_{v \in \mathcal{N}(\boldsymbol{x}_0) \cap D_1} m(\boldsymbol{x}_v; \boldsymbol{\theta}(D_2)) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_0) \cap D_2} m(\boldsymbol{x}_w; \boldsymbol{\theta}(D_1)).$$
(11)

Similarly, we define the function  $F(x_0; D_L)$  formulating the relation between training data  $D_L$  and the prediction for test data  $x_0$  under the no-bagging mode:

$$\boldsymbol{F}(\boldsymbol{x}_{0}; D_{\mathcal{L}}) := \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}; \boldsymbol{\theta}(D_{1})) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}; \boldsymbol{\theta}(D_{2})).$$
(12)

Notice here  $\theta(D_1)$  is the model parameters of Chunk I involving information of true label  $Y_1$ . We aim to examine bagging and stacking strategies effectively preserve the information of label  $Y_1$  via

introducing randomness to the function  $E(x_0; D_L)$  while  $F(x_0; D_L)$  has high risk of leaking the information of true label  $Y_1$ .

We first reiterate the definition of Rényi Differential Privacy.

**Definition 1.** (*Rényi Differential Privacy (Mironov, 2017)*). Consider a randomized algorithm  $\mathcal{M}$  mapping from  $\mathcal{D}$  to real-value  $\mathcal{R}$ . Such an algorithm is said to have  $\epsilon$ -Rényi Differential Privacy of order  $\alpha$  ( $\alpha > 1$ ) if any  $D, D' \in \mathcal{D}$  with  $d_H(D, D') = 1$ , where  $d_H$  is the Hamming distance (D, D') are also referred to as adjacent datasets):

$$D_{\alpha}(\mathcal{M}(D)||\mathcal{M}(D')) = \frac{1}{\alpha - 1} \log E_{x \sim \mathcal{M}(D')} \left(\frac{\mathcal{M}(D)}{\mathcal{M}(D')}\right)^{\alpha} \le \epsilon.$$
(13)

To proceed in a quantifiable way, we rely on some preliminary results for Rényi Differential privacy and generative model for graph learning algorithms.

**Proposition 1.** Rényi differential privacy is preserved by post-processing (Mironov, 2017). If  $F(\cdot)$  has  $\epsilon$ -Rényi Differential Privacy, then for any randomized or deterministic function g,  $g(F(\cdot))$  satisfies  $\epsilon$ -Rényi Differential Privacy.

**Proposition 2.** The closed-form expression of the Rényi divergence between any two Gaussian distributions is given by  $D_{\alpha}(\mathcal{N}(\mu_0, \sigma_0^2)||\mathcal{N}(\mu_1, \sigma_1^2)) = \frac{\alpha(\mu_1 - \mu_0)^2}{2\sigma_{\alpha}^2} + \frac{1}{1 - \alpha} \ln \frac{\sigma_{\alpha}}{\sigma_0^{1 - \alpha} \sigma_1^{\alpha}}$ , provided that  $\sigma_{\alpha}^2 = (1 - \alpha)\sigma_0^2 + \alpha\sigma_1^2 > 0$  (Van Erven & Harremos, 2014).

**Proposition 3.** Assume f has sensitivity 1 and lower magnitude bound L, i.e., for any pair of adjacent datasets  $D, D' \in \mathcal{D}$ :  $|f(D) - f(D')| \le 1$  and  $|f| \ge L$ , and define the Gaussian multiplicative mechanism

$$GM_{\mu,\sigma}f(D) = f(D)\mathcal{N}(\mu,\sigma^2)$$

Then  $GM_{\mu,\sigma}f$  satisfies  $(\frac{1}{2}, \frac{1}{4\sigma^2L^2} + \frac{1}{2L^2})$ -Rényi Differential Privacy.

Proof. According to Proposition (2):

$$\begin{split} &D_{1/2} \left( \mathcal{N}(f(D) + \mu, f^2(D)\sigma^2) || \mathcal{N}(f(D') + \mu, f^2(D')\sigma^2) \right) \\ &= \frac{(f(D) - f(D'))^2}{2\sigma^2 \left(f^2(D) + f^2(D')\right)} + \ln[\frac{1}{2} \left(f^2(D) + f^2(D')\right)] - \ln|f(D)| - \ln|f(D')| \\ &= \frac{1}{2\sigma^2} - \frac{f(D)f(D')}{\sigma^2 \left(f^2(D) + f^2(D')\right)} + \ln[\frac{1}{2} \left(f^2(D) + f^2(D')\right)] - \ln|f(D)f(D')| \\ &= \frac{1}{2\sigma^2} - \frac{f(D)f(D')}{\sigma^2 \left(f^2(D) + f^2(D')\right)} + \ln|\frac{f^2(D) + f^2(D')}{2f(D)f(D')}| \\ &\leq \frac{1}{2\sigma^2} \frac{1}{f^2(D) + f^2(D')} + \ln(\frac{1}{2|f(D)f(D')|} + 1) \\ &\leq \frac{1}{4\sigma^2 L^2} + \ln(\frac{1}{2L^2} + 1) \\ &\leq \frac{1}{4\sigma^2 L^2} + \frac{1}{2L^2}. \end{split}$$

The first inequality follows from  $|f(D) - f(D')| \leq 1$ , take square for both side  $f^2(D) + f^2(D') \leq 1 + 2f(D)f(D')$ . Then we have  $\frac{1}{2\sigma^2} - \frac{f(D)f(D')}{\sigma^2(f^2(D) + f^2(D'))} \leq \frac{1}{2\sigma^2} \frac{1}{f^2(D) + f^2(D')}$  and  $\frac{f^2(D) + f^2(D')}{2|f(D)f(D')|} \leq \frac{1}{2|f(D)f(D')|} + 1$ , the first inequality holds.

**Proposition 4.** If f has sensitivity 1, i.e., for any pair of adjacent datasets  $D, D' \in \mathcal{D}$ :  $|f(D) - f(D')| \leq 1$ . Define the Gaussian additive mechanism

$$\boldsymbol{G}\boldsymbol{A}_{\sigma}f(D) = f(D) + \mathcal{N}(0,\sigma^2),$$

then Gaussian additive mechanism  $GA_{\sigma}f$  satisfies  $(\alpha, \frac{\alpha}{2\sigma^2})$ -Rényi Differential Privacy (Mironov, 2017).

**Proposition 5.** Consider a multivariate Gaussian distribution, and the random variables are partitioned into two groups  $(z_P, z_Q)$ , the distribution is block matrix format

$$egin{pmatrix} oldsymbol{z}_P \ oldsymbol{z}_Q \end{pmatrix} \sim \mathcal{N}\left( egin{bmatrix} oldsymbol{ar{z}}_P \ oldsymbol{ar{z}}_Q \end{bmatrix}, \quad egin{bmatrix} oldsymbol{\Gamma}_{PP} & oldsymbol{\Gamma}_{PQ} \ oldsymbol{\Gamma}_{QQ} \end{bmatrix}^{-1} \end{pmatrix},$$

where  $\begin{bmatrix} \Gamma_{PP} & \Gamma_{PQ} \\ \Gamma_{QP} & \Gamma_{QQ} \end{bmatrix}$  is precision (inverse covariance) matrix. Then the marginal and conditional distribution can be written as

$$\boldsymbol{z}_{P} \sim \mathcal{N}\left(\bar{\boldsymbol{z}}_{P}, (\boldsymbol{\Gamma}_{PP} - \boldsymbol{\Gamma}_{PQ}\boldsymbol{\Gamma}_{QQ}^{-1}\boldsymbol{\Gamma}_{QP})^{-1}\right),\tag{14}$$

$$\boldsymbol{z}_{P}|\boldsymbol{z}_{Q}=\boldsymbol{z}_{Q}\sim\mathcal{N}\left(\bar{\boldsymbol{z}}_{P}-\boldsymbol{\Gamma}^{-1}\boldsymbol{\Gamma}_{PQ}(\boldsymbol{z}_{Q}-\bar{\boldsymbol{z}}_{Q})\right). \tag{15}$$

Before proceeding to our specific results in the main paper, we also need to describe the graph setting.

**Preliminary 2.** Let  $\mathcal{G} = (V, E)$  be an undirected graph, where V is the set of n nodes and E is the set of edges. The adjacency matrix of  $\mathcal{G}$  is  $W \in \mathcal{R}^{n \times n}$ , the diagonal degree matrix is  $D \in \mathcal{R}^{n \times n}$ . The normalize graph Laplacian can be written as  $N = I - D^{-1/2}WD^{-1/2} = I - S$ . We use  $X \in \mathcal{R}^{n \times p}$  for the feature matrix, where p is the dimension of features. We assume all vertex features X are jointly sampled from a multivariate Gaussian distribution (Jia & Benson, 2021), namely

$$\boldsymbol{X} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Gamma}^{-1}), \qquad \boldsymbol{\Gamma} = c_1 \boldsymbol{I}_n + c_2 \boldsymbol{N},$$
 (16)

where  $I_n$  is identical matrix, N is normalized graph Laplacian. Here  $c_1$  controls noise level and  $c_2$  controls the smoothness over the whole graph.

We now proceed to our specific results in the main paper.

**Theorem 1.** Assume base model m to be a multi-layer (two-layer) perceptron and node features X is sampled from a multivariate Gaussian as in Jia & Benson (2021):

$$X \sim \mathcal{N}(\mathbf{0}, \Gamma^{-1}), \qquad \Gamma = c_1 \mathbf{I}_n + c_2 \mathbf{L},$$

where  $I_n$  is an identity matrix and L is the normalized graph Laplacian. Here  $c_1$  controls noise level and  $c_2$  controls the smoothness over the whole graph.  $E(x_0; D_L)$  and  $F(x_0; D_L)$  are predictions produced by BestowGNN for a data point  $x_0$  with and without bagging mode, respectively. If E has sensitivity 1 and lower magnitude bound L, i.e., for any two adjacent  $D, D' \in D$ :  $|E(x_0; D) - E(x_0; D')| \leq 1$  and  $|E| \geq L$ , then E satisfies  $(\frac{1}{2}, \frac{1}{4\sigma^2 L^2} + \frac{1}{2L^2})$ -Rényi Differential Privacy, where  $\sigma^2$  depends on graph structure G. Meanwhile, F has no privacy guarantee, i.e., the Rényi differential privacy loss (8) is unbounded.

*Proof.* Given the definition of function *E* from above, we have that

$$\boldsymbol{E}(\boldsymbol{x}_{0}; D_{\mathcal{L}}) = \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}; \boldsymbol{\theta}(D_{2})) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}; \boldsymbol{\theta}(D_{1}))$$

$$= \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}\boldsymbol{\theta}(D_{2})) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}\boldsymbol{\theta}(D_{1})),$$
(17)

where the second equality follows from the MLP assumption. Similarly for F we have

$$F(\boldsymbol{x}_{0}; D_{\mathcal{L}}) = \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}; \boldsymbol{\theta}(D_{1})) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}; \boldsymbol{\theta}(D_{2}))$$

$$= \sum_{v \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{1}} m(\boldsymbol{x}_{v}\boldsymbol{\theta}(D_{1})) + \sum_{w \in \mathcal{N}(\boldsymbol{x}_{0}) \cap D_{2}} m(\boldsymbol{x}_{w}\boldsymbol{\theta}(D_{2})).$$
(18)

We now define the adjacent datasets D and D' as follows. Assume  $D = D_1$ ; one data point  $\{x', y'\}$  is then randomly selected from Chunk I and removed  $\{x', y'\}$  from  $D_1$  forming  $D' = D_1 \setminus \{x', y'\}$ . Meanwhile, the unlabeled set  $D_{\mathcal{U}}$  and  $D_2$  remain the same. Our goal is to examine the extent to which E and F may leak information pertaining to  $\{x', y'\}$  when  $\{x', y'\}$  is removed from  $D_1$  as described above.

Denote  $x_v, x_w$  as training data in chunk I and chunk II. Assume  $\begin{pmatrix} x_v \\ x_w \end{pmatrix}$  is drawn from a multivariate Gaussian distribution:

$$\begin{pmatrix} \boldsymbol{x}_{v} \\ \boldsymbol{x}_{w} \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Gamma}_{vv} & \boldsymbol{\Gamma}_{vw} \\ \boldsymbol{\Gamma}_{wv} & \boldsymbol{\Gamma}_{ww} \end{bmatrix}^{-1} \right),$$
(19)

where  $\begin{bmatrix} \Gamma_{vv} & \Gamma_{vw} \\ \Gamma_{wv} & \Gamma_{ww} \end{bmatrix} = aI + bN$ , *I* is identical matrix, *N* is normalized graph Laplacian, *a* controls noise level and *b* controls the smoothness over the whole graph.

From Proposition 5, the condition distribution of  $x_w$  given  $x_v = x_v$  can be written as

$$\boldsymbol{x}_w | \boldsymbol{x}_v = \boldsymbol{x}_v \sim \mathcal{N}(-\boldsymbol{\Gamma}_{ww}^{-1}\boldsymbol{\Gamma}_{wv}\boldsymbol{x}_v, \boldsymbol{\Gamma}_{ww}^{-1}).$$

Condition on the data  $D_1$ , the distribution of  $D_2$  is a conditional multivariate Gaussian distribution with mean  $-\Gamma_{ww}^{-1}\Gamma_{wv}x_v$  and variance  $\Gamma_{ww}^{-1}$ . Furthermore, multiplicative Gaussian distribution  $x_w\theta(D_1)$  introduces a Gaussian random noise into (17). According to Proposition (1) and (3), Esatisfies  $(\frac{1}{2}, \frac{1}{4\sigma^2 L^2} + \frac{1}{2L^2})$ -Rényi Differential Privacy, where  $\sigma^2$  depends on  $\Gamma_{ww}^{-1}$  decided by graph structure.

Meanwhile, although (18) is deterministic, we can manually add Gaussian noise  $\mathcal{N}(0, \sigma^2)$  such that  $\mathbf{F}$  satisfies  $\frac{\alpha}{2\sigma^2}$ -Rényi Differential Privacy via Proposition (4). However, if we then let  $\sigma \to 0$  to reproduce  $\mathbf{F}$ , we have that  $\frac{\alpha}{2\sigma^2} \to \infty$ , indicating that in fact  $\mathbf{F}$  has no privacy guarantee.

# **B** EXPERIMENT DETAILS

## **B.1** DATA DESCRIPTIONS

**House**: node features are the property of house, edges connect the neighbors, the task is to predict the price of the house. **County**: each node is a county and edges connect two counties sharing a border, the task is to predict the unemployment rate for a county. **VK**: each node is a person and edges connect two people based on the friendships, the task is to predict the age of each person. **Avazu**: each node is a device and edges connect two devices if they appear on the same site with the same application, the target is the click-through-rate of a node. For **House**, **County**, **VK** and **Avazu** datasets, Training/validation/testing are randomly split with 6/2/2 ratio and all experiments results are averaged over 5 trails.

OGB-Arxiv, OGB-Products are standard datasets from OGB-leaderboards and all training/validation/testing splits follow the standard data splitting from OGB-leaderboards. Reddit is standard datset from Deep Graph Library (DGL).

### B.2 BASE MODELS

Specifically, we consider LightGBM boosted Tress (GBM) (Ke et al., 2017), CatBoost boosted trees (CAT) (Prokhorenkova et al., 2018), fully-connected neural networks (NN), Extremely Randomized Trees (RT), Random Forests (RF), K Nearest Neighbors (KNN), Label Propagation (LP) (Huang et al., 2020) and Transformer with electra pretrained model (Text) (Training epoch is 12) (Clark et al., 2020). For the first layer, we keep the typical models, for example, Gradient Boosted Decision Trees for Tabular data, Transformer models for text data. For second stacking layer, we use all of models except extremely low-efficient models for large dataset, for example, KNN and Catboost slow down the training procedure for OGB-products dataset. All details about the base models can be found in table 5. The parameters about all models can be referred to AutoGluon (Erickson et al., 2020).

## **B.3** PARAMETERS FOR GRAPH-AWARE PROPAGATION

We do graph-aware propagation for the prediction to incorporate the graph structure. Table 6 shows two hyperparameters considered in the propagation part: weight  $\lambda$  and number of propagation step T. We also present the hyperparameters for Correct and Smooth in Table 7.

# Table 5: Base models

DATASET	FIRST LAYER	SECOND LAYER	
HOUSE/COUNTY/VK/AVAZU	CAT, GBM, NN	KNN, GBM, RF, RT, CAT, NN	
OGB-ARXIV	TEXT, GBM, NN	GBM, RF, RT, NN	
OGB-PRODUCTS	Text, LP	GBM, RF, RT, NN	
Reddit	NN, RF, RT, GBM	NN, RF, RT, KNN, GBM	

## Table 6: Hyperparameters

DATASET	$\lambda$	INPUT FOR STACKING LAYER
HOUSE/COUNTY/VK/AVAZU	0.9	$(m{X}, \{m{F}_m^{(0)}, m{F}_m^{(1)}, m{F}_m^{(2)}, m{F}_m^{(3)}, m{F}_m^{(4)}, m{F}_m^{(5)}\})$
OGB-ARXIV	0.95	$(m{X}, \{m{F}_m^{(0)}, m{F}_m^{(1)}, m{F}_m^{(3)}, m{F}_m^{(5)}, m{F}_m^{(7)}, m{F}_m^{(9)}\})$
OGB-PRODUCTS REDDIT	0.97 0.95	$(\boldsymbol{X}, \{\boldsymbol{F}_m^{(0)}, \boldsymbol{F}_m^{(1)}, \boldsymbol{F}_m^{(3)}, \boldsymbol{F}_m^{(5)}, \boldsymbol{F}_m^{(7)}, \boldsymbol{F}_m^{(9)}\}) \\ (\boldsymbol{X}, \{\boldsymbol{F}_m^{(0)}, \boldsymbol{F}_m^{(1)}, \boldsymbol{F}_m^{(2)}, \boldsymbol{F}_m^{(3)}, \boldsymbol{F}_m^{(4)}, \boldsymbol{F}_m^{(5)}\})$

## Table 7: Hyperparameters for C&S

DATASET	$\lambda_1$	KERNEL TYPE	$\lambda_2$	KERNEL TYPE	NUM_PROPAGATION
HOUSE/COUNTY/AVAZU	0.8	DA	0.5	DA	5
VK	0.8	DA	-	-	5
OGB-ARXIV	0.9	DA	0.1	AD	50
OGB-PRODUCTS	0.3	DAD	0.3	AD	50

 Table 8: Training time tested on AWS g4dn.12xlarge machine.

DATASET	BASE MODEL	TIME(S)
HOUSE	GBM, NN	52
County	GBM, NN	18
VK	GBM, NN	119
Avazu	GBM, NN	15
OGB-Arxiv	NN	199
OGB-PRODUCTS	NN	837

## B.4 COMPUTING COST

The computing cost depends on the ensemble models we select (e.g., transformer models can take more computing resources relying on the implementation, including more emsemble models leads to more computing cost). So it's hard to consistently measure the training/inference time or memory consumption. But the computing cost is in a competitive range since the integration of the bagging and ensembling parts key to our model can be efficiently implemented, e.g., via open source packages like AutoGluon that we used. In Table 8, we present the training time of different datasets with basic ensemble models. For instance, the training time for OGB-products with OGB embeddings is around 800s, while for GraphSage it is about 1000s for 100 epochs.

# ADDITIONAL REFERENCES FOR THE SUPPLEMENTARY

- K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*, 2020.
- N. Erickson, J. Mueller, A. Shirkov, H. Zhang, P. Larroy, M. Li, and A. Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv:2003.06505*, 2020.
- Q. Huang, H. He, A. Singh, S.-N. Lim, and A. R. Benson. Combining label propagation and simple models out-performs graph neural networks. arXiv preprint arXiv:2010.13993, 2020.
- J. Jia and A. R. Benson. A unifying generative model for graph learning algorithms: Label propagation, graph convolutions, and combinations. *arXiv preprint arXiv:2101.07730*, 2021.
- G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu. LightGBM: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, 2017.
- I. Mironov. Rényi differential privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF), pages 263–275. IEEE, 2017.
- L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin. CatBoost: unbiased boosting with categorical features. In *Advances in Neural Information Processing Systems*, 2018.