# Stable Prediction on Graphs with Agnostic Distribution Shift

Shengyu Zhang[1], Kun Kuang[1], Jiezhong Qiu[2], Jin Yu[3], Zhou Zhao[1], Hongxia Yang[3], Zhongfei Zhang[1] and Fei Wu[1]

[1] College of Computer Science and Technology, Zhejiang University, China
[2] Department of Computer Science and Technology, Tsinghua University, China
[3] Alibaba Group, China

{sy_zhang,kunkuang,zhaozhou,wufei}@zju.edu.cn
{kola.yu,yang.yhx}@alibaba-inc.com
qiujz16@mails.tsinghua.edu.cn

## ABSTRACT

Graph is a flexible and effective tool to represent complex structures in practice and graph neural networks (GNNs) have been shown to be effective on various graph tasks with randomly separated training and testing data. In real applications, however, the distribution of training graph might be different from that of the test one (*e.g.*, users' interactions on the user-item training graph and their actual preference on items, *i.e.*, testing environment, are known to have inconsistencies in recommender systems). Moreover, the distribution of test data is always agnostic when GNNs are trained. Hence, we are facing the agnostic distribution shift between training and testing on graph learning, which would lead to unstable inference of traditional GNNs across different test environments. To address this problem, we propose a novel stable prediction framework for GNNs, which permits both locally and globally stable learning and prediction on graphs. In particular, since each node is partially represented by its neighbors in GNNs, we propose to capture the stable properties for each node (locally stable) by re-weighting the information propagation/aggregation processes. For global stability, we propose a stable regularizer that reduces the training losses on heterogeneous environments and thus warping the GNNs to generalize well. We conduct extensive experiments on several graph benchmarks and a noisy industrial recommendation dataset that is collected from 5 consecutive days during a product promotion festival. The results demonstrate that our method outperforms various SOTA GNNs for stable prediction on graphs with agnostic distribution shift, including shift caused by node labels and attributes.

## KEYWORDS

Stability; Stable Prediction; Graph Neural Network; Distribution Shift

## 1 INTRODUCTION

Graphs are commonly used to represent the complex structures of interacting entities in a flexible and effective manner, with applications and domains varying from social networks, knowledge graphs, and recommender systems. The challenging and open-ended nature of graph representation learning lends itself to a variety of diverse models [13]. Recent advances in the literature have convincingly demonstrated the high capability of Graph Neural Networks (GNNs) [18, 31, 35] on graph representation learning and inference. GNNs typically follow a neighborhood aggregation scheme, *i.e.*, a node is represented by recursively aggregating other nodes' information within the neighborhood [11, 43]. Since messages pass along the edges, GNNs can effectively capture the high-order structural correlations between nodes.

Most of existing GNNs are proposed and evaluated with an underlying assumption, *i.e.*, the distribution of collected training data is consistent with the distribution of testing data. Under this assumption, many studies follow the convention of performing random splits on graph benchmarks [18, 42] to train and test GNNs. Unfortunately, this assumption can hardly be satisfied in real-world application-specific systems [15, 24] due to the ubiquitous selection bias, leading to distribution shift between training and test data. Moreover, the test distribution is always agnostic during optimizing GNNs on training graph data. The agnostic distribution shift includes shift caused by both node labels and node attributes. For example, in recommender systems, a notorious problem is that the logged interactions collected from current undergoing systems may warp (graph) recommenders biased towards popular items, *i.e.*, label-related distribution shift. Such a shift may degrade the recommendation effectiveness on deployed environment, *e.g.*, causing unfairness towards less exposed items. In addition, male users often have less logged interactions on e-commerce platforms since they often online-shop with specific purposes. Therefore, the (graph) recommenders may be trained with bias towards female users and thus deteriorates the experience of male users, *i.e.*, node attribute related distribution shift. Generally, the agnostic distribution shift between training and testing would render traditional GNNs over-optimized on the labeled training samples and render their predictions error-prone on test samples, resulting in unstable predictions. Therefore, learning GNNs that are resilient to distribution shifts and able to

make stable predictions on graphs will be important for real-world applications.

Many techniques [4, 9, 16, 23, 33] designed for general machine learning problems are proposed in the literature of stable learning. These methods typically propose to re-weight training samples with a density ratio, thus driving the training distribution closer to the testing distribution. However, how to learn structural node representations based on which we can make stable prediction on graphs across agnostic environments remains largely unexplored in the literature of graph neural networks. More recently, [14] proposes to learn a stable graph that captures general relational patterns, and directly optimizes the adjacency matrix with a pre-defined set generation task. However, they mainly aim to obtain an optimized and static adjacency matrix for a given graph and cannot dynamically process unseen neighborhood structures. Their framework also has the disadvantage of being tightly coupled with the set generation problem and cannot be easily adapted to other graph tasks.

In this paper, we aim to learn task-agnostic GNNs that can make stable predictions on unseen neighborhood structures. Stable prediction poses a quantitative goal that, given *one* observational graph environment, a stably trained GNN should achieve a high average score and a low variance of scores on *multiple* agnostic testing environments. We propose a novel stable prediction framework for GNNs, which permits both locally stable learning at the node-level, and globally stable learning at the environment-level. The proposed framework first performs biased selection on the observational environment and constructs multiple training environments. For *locally stable learning*, we start from the perspective that each node is partially represented by other nodes in the neighborhood in GNNs, and we propose to capture stable properties by re-weighting the neighborhood aggregation process. Under the invariant pre-diction assumption for causal inference [30], we can improve the stability of predictions when the model relies more on stable prop-erties. For *globally stable learning*, we inspect the training of GNNs at the environment-level and empirically find that the losses for different environments progressively diverge in biased training, which eventually leads to unstable performance across environ-ments. Therefore, we propose to reduce the training gap between environments to explicitly warp the GNNs to generalize well across environments.

We conduct experiments on various public graph benchmarks and a real-world recommendation dataset that is collected from a world-leading e-commerce platform during an annual product-promotion festival where selection biases naturally exist [1]. We evaluate a bunch of generic SOTA GNNs and methods that are specifically designed for mitigating selection biases. We concern both traditional task-specific evaluation metrics and protocols that are especially designed for stable learning [20, 21]. Extensive results demonstrate the capability of our framework on learning GNNs that make stable predictions on graphs. In summary, the contribution of this paper are:

- We propose to achieve stable prediction on graphs for explic-itly reducing the performance variances across environments with distribution shifts. This is less explored in the literature.

---

[1]We will release the desensitized version to promote further investigations.

- We devise a novel stable prediction framework for GNNs, which captures stable properties for each node, based on which we learn node representations and make predictions (*locally stable*), and regularizes the training of GNNs on het-erogeneous environments (*globally stable*).
- We conduct comprehensive experiments on different public graph benchmarks and a noisy industrial recommendation dataset, which jointly demonstrate the effectiveness of the proposed framework.

## 2 RELATED WORKS

The recent advances in graph neural networks [7, 18, 19, 27, 35, 38, 39, 47] have convincingly demonstrated high capability in cap-turing structural and relational patterns within graphs. Typically, GNNs follow a message passing schema for representation learn-ing by transforming and aggregating the information from other nodes within the neighborhood. Different neighborhood aggrega-tion variants [6, 8, 10, 17, 22, 25, 31, 34, 36, 40, 41, 44–46, 48, 50] have been proposed to reach state-of-the-art performances in vari-ous tasks. Typically, GAT [35] introduces the attention mechanism into the information aggregation process. SGC [39] simplifies the original Graph Convolutional Network [18] by linearly propagating information and collapsing weights among graph layers. APPNP [19] extends the utilized neighborhood for node representation and achieves an adjustable neighborhood for classification.

Many existing GNNs are evaluated on graph benchmarks that are randomly split [18, 42], *i.e.*, the training and testing data share similar data distribution. However, in real-world applications, *e.g.*, recommender systems, training samples are observed and collected with selection bias, leading to inconsistencies between the training and testing distribution. Few works in the literature investigate such a real-world problem. Recently, GNM [49] confronts a related problem named non-ignorable non-response, which indicates that the unlabeled nodes are missing not at random (MNAR). However, they only consider distribution caused by node labels and neglect distribution shift related to node attributes, and they solely dis-cuss binary-class datasets in the experiments, which can be less plausible for real-world scenarios. In this paper, we propose a frame-work that can alleviate the negative effects from shift related to both labels and attributes, and obtains stable predictions on various graph benchmarks and real-world recommendation datasets. [14] proposes to learn a static adjacency matrix for a given graph and expects that the learned adjacency matrix captures general rela-tional patterns that are free from selection biases. We largely differ from this work as illustrated in the Introduction.

Recently, many methods [4, 9, 16, 23, 33] are proposed to address selection bias for general machine learning problems. They typi-cally align the training distribution with the testing distribution by re-weighting the training samples. Following these works, [20, 21] first propose the stable learning framework as well as the evalu-ation protocols. They mainly decorrelate different dimensions of the hidden variable, and they prove that the prediction based on decorrelated variables should be stable. More recently, GAT-DVD [3] directly borrows such a decorrelation idea to the literature of GNNs. In this paper, instead of borrowing generic off-the-shelf

stable learning tools, we propose to inspect the neighborhood aggregating process and design graph-specific architectures. Besides, we show a large improvement over GAT-DVD in the experiments.

## 3 PROBLEM FORMULATION

Let $\mathcal{X}$ denote the space of observed features, $\mathcal{A}$ denote the space of adjacency matrix, and $\mathcal{Y}$ denote the outcome space. Following [30], we define a graph **environment** as the joint distribution $P_{XAY}$ on $\mathcal{X} \times \mathcal{A} \times \mathcal{Y}$ and use $\mathcal{E}$ to denote the set of all environments. For each environment, we have a graph dataset $G^e = (\mathbf{X}^e, \mathbf{A}^e, Y^e)$, where $\mathbf{X}^e \in \mathcal{X}$ are node features, $\mathbf{A}^e \in \mathcal{A}$ is the adjacency matrix, and $Y^e \in \mathcal{Y}$ is the response variable (*e.g.*, node labels in the node classification problem). The joint distribution of features and outcomes on $(\mathbf{X}, \mathbf{A}, Y)$ can vary across environments, *i.e.*, $P^e_{XAY} \neq P^{e'}_{XAY}$ for $e, e' \in \mathcal{E}$, and $e \neq e'$. In this paper, we aim to learn node representations based on which we can make stable predictions across environments with various degrees of selection biases. Before giving the specific definition of stable prediction, we first define *Average_Score* and *Stability_Score* similar to [20] of a GNN model as:

$$Average\_Score = \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathrm{S}\left(G^e\right), \tag{1}$$

$$Stability\_Error = \sqrt{\frac{1}{|\mathcal{E}| - 1} \sum_{e \in \mathcal{E}} \left(\mathrm{S}\left(G^e\right) - Average\_Score\right)^2}. \tag{2}$$

where $|\mathcal{E}|$ denotes the number of environments evaluated, and $\mathrm{S}(G^e)$ refers to the predictive score computed with a specific numerical assessment on dataset $G^e$. Therefore, *Average_Score* indicates the overall performance of the learned GNN on heterogeneous testing environments, and *Stability_Error* indicates how differently the learned GNN performs on these environments. Based on these two metrics, we define the problem of stable prediction on graphs:

PROBLEM 1 (STABLE PREDICTION ON GRAPHS). ***Given*** *one training environment* $e \in \mathcal{E}$ *with dataset* $G^e = (\mathbf{X}^e, \mathbf{A}^e, Y^e)$*, the task is to* ***learn*** *node representations based on which the predictions yield high Average_Score but small Stability_Error across environments* $\mathcal{E}$*.*

## 4 METHODS

### 4.1 A Retrospect of Modern GNNs

Here we briefly summarize typical GNNs and then give an illustration on why they suffer from the distribution shift. Modern GNNs follow a neighborhood aggregation schema by iteratively aggregating representations of its neighbors to update the representation of the target node. One iteration can be generally formulated as:

$$\mathbf{x}_i^* = \text{AGGREGATE}\left(\left\{\mathbf{x}_j : j \in \mathcal{N}_i\right\}\right), \tag{3}$$

$$\mathbf{x}_i' = \text{COMBINE}\left(\mathbf{x}_i^*, \mathbf{x}_i\right). \tag{4}$$

where $\mathcal{N}_i$ denotes the indices of nodes in the neighborhood of node $i$, $\mathbf{x}_i'$ denotes the updated representation of node $i$. The modeling of AGGREGATE and COMBINE can be various and essential for different GNN architectures. Most GNNs are optimized to obtain an overall best performance on the observational environment, *i.e.*, $G^0 = (\mathbf{X}^0, \mathbf{A}^0, Y^0)$. In this way, we may expect an overly-optimized solution after training and some distribution-specific patterns may warp the GNN biased towards a globally sub-optimal solution since

the distribution of real-world testing data can be agnostic and mostly yields shift from the training graph distribution in real-world applications (*e.g.*, recommender systems). We denote this phenomenon as *global instability*.

In the local view, *i.e.*, neighborhood aggregation process, each neighbor node contributing to the final aggregated representation can be viewed as a property of the root node. Most GNNs are proposed disregarding a fundamental problem, *i.e.*, "what are the **stable** properties of the root node?" The connection between capturing stable properties and stable prediction can be interpreted as a plausible hypothesis:

HYPOTHESIS 1. *If we consider all stable properties ("direct causes") of the target node of interest, then the conditional distribution of the outcome given the stable properties will not change when interventions (e.g., biased selection) are performed without affecting the target and stable properties themselves.*

This stability idea is closely linked to causality and has been discussed or empirically demonstrated to be effective in the literature [2, 12, 28]. In other words, we can achieve stable prediction on graphs across environments with various selection biases by capturing the stable properties for each node. Indeed, some existing GNNs (*e.g.*, GAT [35]) propose to specify different weights to different nodes in a neighborhood and somewhat "stabilize" the weights learning process by leveraging multiple heads. However, we argue that the weights they learned explore the subtle statistical relationship in training data and can be far from being stable due to potential confounders [29]. A confounder denotes the common cause of the predictor and the outcome, which may lead to spurious correlations between the root node (outcome) and some "important" neighbors (predictor). We denote this phenomenon as *local instability*.

### 4.2 Stable Learning on Graphs

The essence of our methodology for stable learning on graphs is depicted in Figure 1. Given the environments priorly generated by various selection biases, we propose the *locally stable learning* strategy, which explicitly captures the stable properties in learning the representation of each node by re-weighting the information propagation/aggregation processes, and the *globally stable learning*, which regularizes the errors across environments to be closer and thus improving the stability. We borrow the idea from the literature of causality [30] and propose to first perform several interventions on the observational data based on selection biases, and create some environments $\mathcal{E}$ accordingly. We note that we keep an observational environment in $\mathcal{E}$.

*4.2.1 Locally Stable Learning.* In the neighborhood aggregation process, we regard nodes that are consistently important across environments as stable properties, as depicted in Figure 1. Therefore, the final representation of the target node of interest is primarily based on such stable properties.

Given one environment $e$, we obey a weighted schema of the commonly used neighborhood aggregation in most GNNs:

$$\mathbf{x}_i' = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij}^e \mathbf{W} \mathbf{x}_j\right), \tag{5}$$
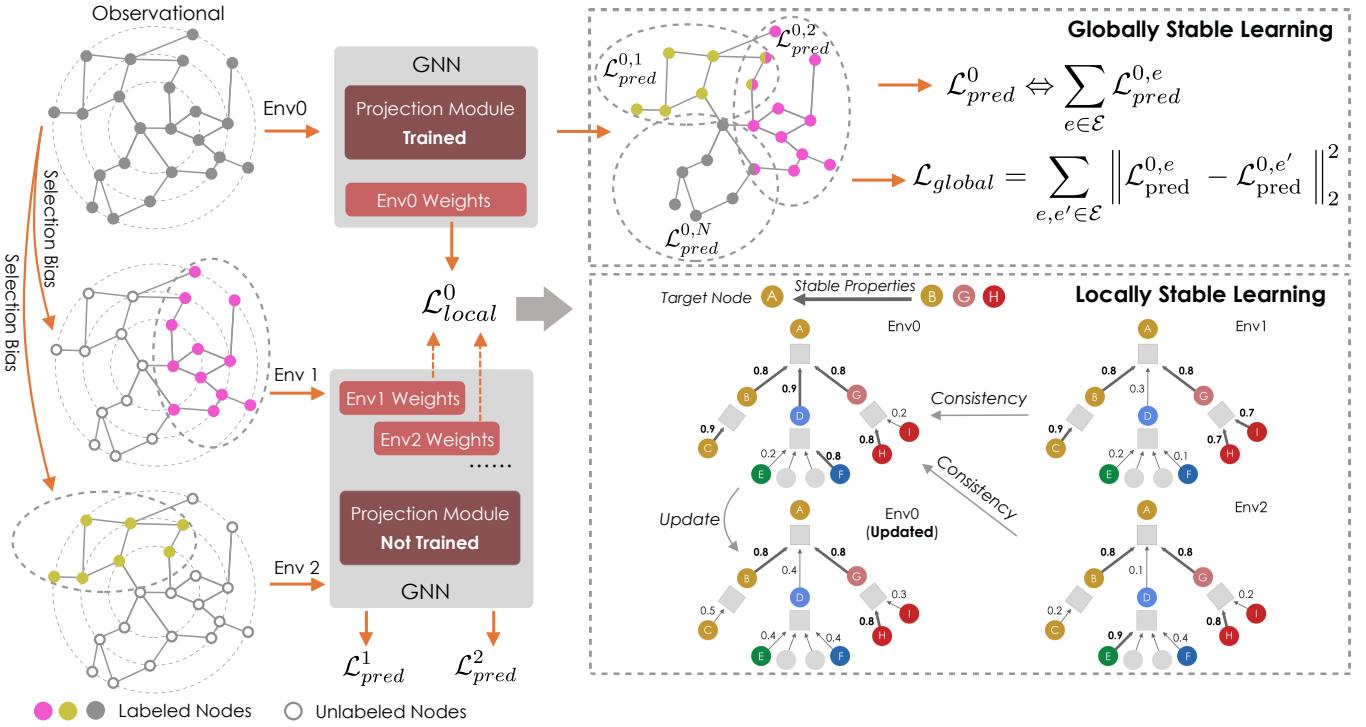
**Figure 1: The overall schema of the proposed framework, which consists of two essential components, *i.e.*, the *locally stable learning* that captures properties that are stable across environments in the representation learning of each target node, and the *globally stable learning* that explicitly balances the training of different environments.**

where $\mathbf{x}'_i$ denotes the updated representation of the target node $x_i$, $\mathbf{W}$ denotes the weight matrix of the linear transformation, $\sigma$ denotes the nonlinear activation function, $\mathcal{N}_i$ denotes the indices set of node $x_i$'s neighbors, and $\alpha^e_{ij}$ is the weight that indicates the importance of one property $x_j$ to the target node $x_i$ learned in environment $e$. The weight $\alpha^e_{ij}$ can be learned in various ways and we adopt the formulation of graph attention as the weight predictor $\varphi^e$ for its simplicity:

$$\alpha^e_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\mathbf{a}^e\left[\mathbf{W}\mathbf{x}_i\|\mathbf{W}\mathbf{x}_j\right]\right)\right)}{\sum_{k\in\mathcal{N}_i}\exp\left(\text{LeakyReLU}\left(\mathbf{a}^e\left[\mathbf{W}\mathbf{x}_i\|\mathbf{W}\mathbf{x}_k\right]\right)\right)} \quad (6)$$

where LeakyReLU denotes the nonlinearity (with negative input slope 0.2), and $\mathbf{a}^e$ denotes the parameter vector to determine the weight $\alpha^e_{ij}$ for environment $e$. The weights in graph attention indicate relative importance in the neighborhood, and it is also acceptable to use the following formulation, which has a sense of absolute importance:

$$\alpha^e_{ij} = \text{sigmoid}\left(\mathbf{a}^e\left[\mathbf{W}\mathbf{x}_i\|\mathbf{W}\mathbf{x}_j\right]\right), \quad (7)$$

We note that each environment $e$ has its parameter vector $\mathbf{a}^e$, which helps find important properties especially for this environment (or data distribution). To identify and leverage properties that are consistently important across all environments $\mathcal{E}$, a straightforward solution is to treat properties with weights higher than a certain threshold as stable properties and manually encourage the model to

rely on them for prediction while suppressing the others. However, this solution has at least two drawbacks: 1) manually modifying the weights of properties will lead to inconsistencies between the current weight prediction module and other modules in GNN since they will be trained based on the modified weights rather than the predicted weights; 2) some properties that are not important may happen to obtain high weights during training and model regarding them as stable properties may draw false conclusions.

In this regard, we propose to identify and leverage stable properties in a soft way: 1) we propose to pull the weight predicted across environments by using a distance loss; 2) we will regard properties with weights that are both high in value and can be easily pulled together as stable properties and suppress the others. In this way, we can confront the randomness of gradient descent based training (*i.e.*, the second problem) by pulling sensitive weights towards the corresponding insensitive ones, and confront the inconsistent training problem (*i.e.*, the first problem) since other parts in GNNs will always reply on the currently predicted weights. To further simplify the training process, we note that the first pulling process can also help to suppress properties that yield inconsistent weights since they will be driven towards the average. Therefore, the locally stable regularizer can be written as the following:

$$\mathcal{L}^e_{local} = \sum_{i\in\mathcal{V}}\sum_{j\in\mathcal{N}_i}\sum_{e'\in\mathcal{E}}\text{Dist}(\alpha^e_{ij},\alpha^{e'}_{ij}), \quad (8)$$

where $\mathcal{V}$ denotes the indices set of all nodes, and Dist denotes a distance function which is set as the L2 distance in our experiment for its empirical effectiveness. When there are multiple GNN layers equipped with the re-weighting module, the regularizer can be written as:

$$\mathcal{L}_{local}^e = \sum_{l=1}^{N_l} \mathcal{L}_{local}^{e,l}. \tag{9}$$

where $N_l$ denotes the number of GNN layers equipped with the re-weighting module.

*4.2.2 Globally Stable Learning.* Besides investigating stable prediction from a local view, *i.e.*, learning node representations that capture stable properties, we further investigate stable prediction from a global and environment-level view. Recall that one of the fundamental goals of stable prediction is to reduce *Stability_Error*, *i.e.*, the standard deviation of scores across environments. However, although we can individually compute scores for all environments, it is not directly applicable to optimize objectives related to non-differentiable metrics. We start with a perspective that model performances correlate well with training losses and empirically find that, in the observational environment, the sub-losses that belong to different environments progressively diverge during training with selection bias. Such divergence in training losses will eventually lead to gaps in testing performance. In this regard, we propose to explicitly reduce the gap between losses across environments to rectify the unstable training. Mathematically, the proposed stable regularizer for globally stable learning can be formularized as the following:

$$\mathcal{L}_{global} = \sum_{e,e' \in \mathcal{E}} \left( \mathcal{L}_{pred}^{0,e} - \mathcal{L}_{pred}^{0,e'} \right)^2, \tag{10}$$

where $\mathcal{L}_{pred}^0$ denotes the task-specific loss computed in the observational environment, *i.e.*, environment 0, and $\mathcal{L}_{pred}^{0,e}$ denotes the sub-loss that belongs to environment $e$, *i.e.*,

$$\mathcal{L}_{pred}^{0,e} = \sum_{l \in \mathcal{Y}_L^0} \mathbb{1}(l \in \mathcal{Y}_L^e) \mathcal{L}_{pred,l}^0. \tag{11}$$

where $\mathcal{Y}_L^e$ denotes the indices set of nodes that have labels for environment $e$. We note that minimizing the pair-wise distance of losses as defined in Equation 10 is equivalent to minimizing the variance of losses. Since the negative of loss function can be interpreted as a certain scoring function, it is equivalent to minimizing the stability error defined in Equation 2.

*4.2.3 Training.* With the proposed two building blocks, *i.e.*, locally and globally stable learning, we give a detailed illustration of how we train the entire framework. The training procedure is summarized in Algorithm 1. Given one observational environment 0 for training with dataset $G^0 = (\mathbf{X}^0, \mathbf{A}^0, Y^0)$, we perform biased selection with one or more factors (*e.g.*, node label or semantic node attributes) on $G^0$ at the beginning of training or each training epoch. After selection, we obtain several environments $\mathcal{E}$ with the corresponding graph datasets $\{G^e = (\mathbf{X}^e, \mathbf{A}^e, Y^e)\}_{e=0,...,|\mathcal{E}|}$. Each environment keeps a weight predictor that is individually trained, and all the environments share the same GNN backbone, which

---

**Algorithm 1:** Stable Learning on Graphs

**Input:** Observational graph datasets $G^0 = (\mathbf{X}^0, \mathbf{A}^0, Y^0)$
**Output:** Parameters of stably learned GNN backbone $\theta$ and weight predictors $\varphi^0$ for the observational environment
Perform biased selection on $G^0$ and obtain $\{G^e = (\mathbf{X}^e, \mathbf{A}^e, Y^e)\}_{e=1}^{|\mathcal{E}|}$. Initialize $\theta$ and $\{\varphi^e\}_{e=0}^{|\mathcal{E}|}$
**while** *not converged* **do**
    Freeze $\theta$                  // Disable gradients
    **for** *e = 1 to $|\mathcal{E}|$* **do**
        Optimize $\varphi^e$ to minimize $\mathcal{L}_{pred}^e$, and cache $\alpha^e$
    **end**
    Unfreeze $\theta$             // Enable gradients
    Compute $\mathcal{L}_{pred}^0$ and $\mathcal{L}_{global}$ as in Eq. 10-11, and
    compute $\mathcal{L}_{local}^0$ with cached $\{\alpha^e\}_{e=1}^{|\mathcal{E}|}$ as in Eq. 8.
    Optimize $\theta, \varphi^0$ to minimize $\mathcal{L}$ as in Eq. 12.
**end**

---

is solely trained on environment 0, *i.e.*, the observational environment. We train the entire framework following an environment-by-environment procedure and update the corresponding parameters with the task-specific objective $\mathcal{L}_{pred}^e$ (*e.g.*, node classification, recommendation). For the observational environment, we train the weight predictor for this environment and the GNN backbone with combined loss function:

$$\mathcal{L} = \mathcal{L}_{pred}^0 + \lambda_0 \mathcal{L}_{local}^0 + \lambda_1 \mathcal{L}_{global}. \tag{12}$$

We do not train other environments with the corresponding $\mathcal{L}_{local}^e$ since we mainly aim to rectify the importance weights in the observational environment. We note that the GNN backbone and the weight predictor in the observational environment will be used for testing/inference.

## 5 EXPERIMENTS

We note that instability in prediction arises due to the inconsistencies of the joint distribution of features, adjacency matrix, and outcomes $(\mathbf{X}, \mathbf{A}, Y)$ across environments (*e.g.*, training/testing environments) according to Section 3. To test the stability of all methods, for each dataset, we construct several training and testing environments by varying the joint distribution of $(\mathbf{X}, Y)$. We consider both shift related to node labels by directly varying $Y$ and shift related to node attributes by varying additional factors that well correlate with $\mathbf{X}$. We disregard the shift related to $\mathbf{A}$ in the experiment.

Following previous GNNs [18, 35] , we evaluate the proposed stable learning framework on public graph benchmarks, including the recently proposed Open Graph Benchmark [15], in Section 5.1. Since these datasets seldom accompany meaningful attributes, we are mainly concerned with *label*-related shift. We conduct hyper-parameter analysis to obtain a better understanding of the designs. We further demonstrate its efficacy on real-world recommendation scenarios that are known to be full of sample selection biases [5] in Section 5.2. We mainly report results with shift caused by node attributes as a complement to label shift on graph benchmarks. We
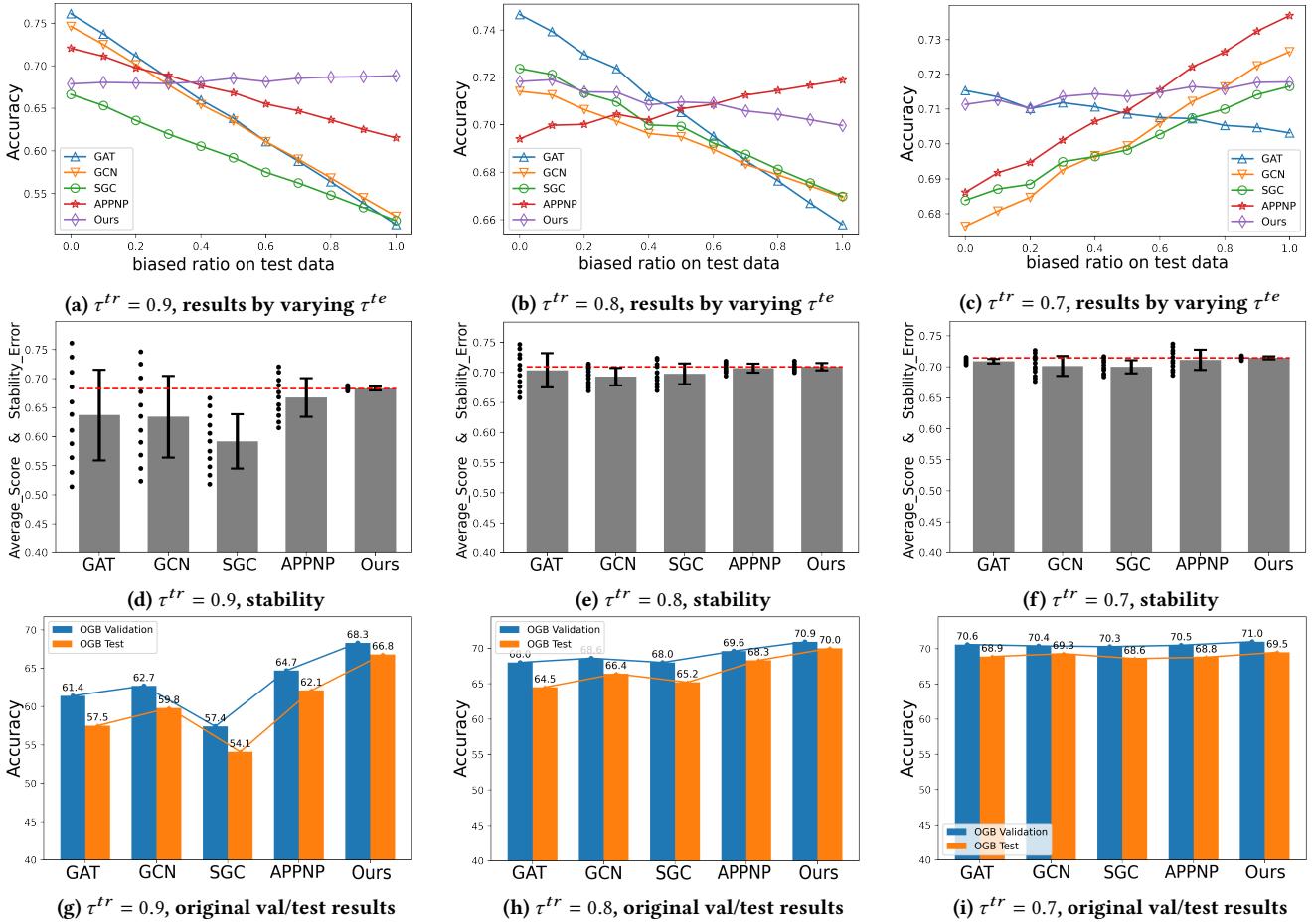
(a) $\tau^{tr} = 0.9$, results by varying $\tau^{te}$     (b) $\tau^{tr} = 0.8$, results by varying $\tau^{te}$     (c) $\tau^{tr} = 0.7$, results by varying $\tau^{te}$

(d) $\tau^{tr} = 0.9$, stability     (e) $\tau^{tr} = 0.8$, stability     (f) $\tau^{tr} = 0.7$, stability

(g) $\tau^{tr} = 0.9$, original val/test results     (h) $\tau^{tr} = 0.8$, original val/test results     (i) $\tau^{tr} = 0.7$, original val/test results

Figure 2: Testing results on OGB-Arxiv dataset by varying the training bias ratio $\tau^{tr}$ among $\{0.9, 0.8, 0.7\}$. Subfigures 2a - 2c show the metrics by varying the testing bias ratio. Subfigures 2d - 2f explictly show the stability of prediction by plotting the *Average_Score* and *Stability_Error*. Subfigures 2g - 2i show the results on the original validation and testing datasets.

also consider the settings with agnostic real-world selection biases in Section 5.2.3.

## 5.1 Experiments on Graph Benchmarks

*5.1.1 Experimental Setup.* We conduct experiments on the following datasets, of which the statistics are listed in Table 1:

**OGB Arxiv.** The OGB datasets are proposed with real-world (non-random) training/testing split, which is suitable for stability evaluation. We perform biased selection on the half nodes from the original training set based on their corresponding labels (*i.e.*, label shift) to construct the observational environment. We construct multiple testing environments from the remaining half nodes. For each environment, the probability of node $i$ with label $y$ to be selected can be $P(s_i = 1) = \tau$ if $y \geq 24$ and $1 - \tau$ otherwise. We note that other selection choices are acceptable and we here mainly aim to get originally equal-size environments by splitting in the middle, *i.e.*, 24. $\tau$ denotes the bias ratio indicating the severity of sample selection bias. We use $\tau^{tr}, \tau^{te}$ to denote the training and testing bias ratio, respectively. In the experiments, we vary the training

### Table 1: Statistics of datasets.

**(a) Dataset statistics.**

| Dataset | Citeseer | OGB-Arxiv | Recommendation |
|---|---|---|---|
| Nodes | 3,327 | 169,343 | 29,444 |
| Edges | 4,732 | 1,166,243 | 180,792 |

**(b) Detailed statistics of the recommendation dataset, which contains 5 consecutive days collected during a product promotion festival. (*IteracN* stands for the number of Interactions at Day N.)**

| Users | Items | Iterac1 | Iterac2 | Iterac3 | Iterac4 | Iterac5 |
|---|---|---|---|---|---|---|
| 9,052 | 20,392 | 180,792 | 161,035 | 174,511 | 181,166 | 233,373 |

bias ratio $\tau^{tr}$ among $\{0.9, 0.8, 0.7\}$ indicating heavy, medium, and light selection biases, following [20].

We consider several SOTA GNNs as comparison methods, including GCN [18], GAT [35], SGC [39], and APPNP [? ].
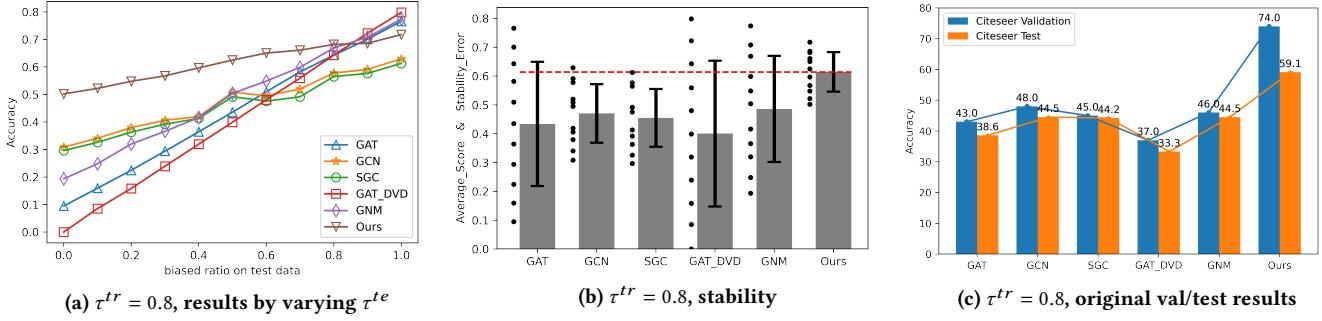
(a) $\tau^{tr} = 0.8$, **results by varying $\tau^{te}$**

(b) $\tau^{tr} = 0.8$, **stability**

(c) $\tau^{tr} = 0.8$, **original val/test results**

**Figure 3: Testing results on Citeseer dataset.**

**Citeseer.** We also report the results on a traditional benchmark, Citeseer [32]. The probability of node $i$ with label $y$ to be selected can be $P(s_i = 1) = \tau$ if $y \geq 3$ and $1 - \tau$ otherwise. We further add two comparison methods that are designed for alleviating selection biases, *i.e.*, GNM [49], and GAT-DVD [3]. We do not test them on OGB-Arxiv since they are not easily adaptable for large-scale datasets. We consider training bias ratio $\tau^{tr} = 0.8$, which is the same as the medium-level selection bias in the Arxiv dataset.

We follow [18, 35] to construct two-layer GCN and GAT. All other methods (including ours) also contain two graph layers for a fair comparison. All methods are with hidden size 250 for OGB-Arxiv and 64 for Citeseer, and learning rate 0.002.

*5.1.2 Stability comparison with SOTA GNNs.* The testing results on the OGB-Arxiv dataset and Citeseer dataset are shown in Figure 2 and Figure 3. Specifically, each figure in Figure 2a - 2c plot the evaluation results across different testing environments with bias ratio $\tau^{te}$ varying among $\{0.0, 0.1, 0.2, \ldots, 1.0\}$. Figure 2a - 2c differ from each by the training bias ratio $\tau^{tr}$. Since stable prediction expects both low performance variances across multiple testing environments and an overall high performance, we explicitly plot the corresponding Stability_Error and Average_Score in Figure 2d - 2f. In a nutshell, our framework achieves more stable results than SOTA GNNs, including both generic ones (GAT, GCN, SGC, and APPNP) and those designed for reducing selection biases (GAT-DVD and GNM). In different testing environments constructed by varying the testing bias ratio, we observe that most GNNs suffer from the distributional shifts and yield poorer performances when the testing distribution is more different from the training distribution (*e.g.*, the right of Figure 2a). Although our framework sacrifices some performance in testing environments with distribution closer to the training (*e.g.*, the left of Figure 2a), our framework obtains significantly higher Average_Score and lower Stability_Error, which are important metrics (as illustrated in Section 3) for stable prediction [20], across heterogeneous environments as indicated in Figure 2d - 2f, 3b.

By varying the training bias ratio from *light* to *heavy* (from Figure 2f to 2d), we can observe that the stability errors are increasing significantly for most conventional GNNs. This means that they are sensitive to the selection biases and the resulted distribution shifts. When the distributional shifts are heavier, the performance across multiple testing environments yields larger differences, *i.e.*,

unstable predictions. Our framework yields the least stability error that can be almost negligible, which demonstrates the superiority. On the other hand, when the distribution shifts are increasing from *light* to *heavy*, the least Average_Score drop of our framework demonstrates that we are not achieving some trivial solutions. For example, one trivial solution is to deteriorate the performance on data distributions similar to the training environments, which can somehow reduce the Stability_Error but have the cost of deteriorating the overall effectiveness. In a nutshell, our framework achieves stable prediction on graphs without sacrificing effectiveness. Although APPNP yields high Average_Score compared to other SOTA GNNs, their performance scores across different testing environments mostly vary significantly (by taking an in-depth analysis on Figure 2c or 2a), *i.e.*, unstable predictions.

Surprisingly, with the same hyper-parameters, GAT-DVD and GNM achieve more unstable results than other generic GNNs. As shown in Figure 3b, GAT-DVD and GNM yield similar or better Average_Score but significantly larger Stability_Error compared to other SOTA GNNs. We attribute these results to that GNM is designed for binary-class datasets and may perform poorly or need further improvements when extending it to multi-classes datasets, and that the work of GAT-DVD [3] is still under development.

*5.1.3 Performance on the original validation/testing datasets.* Since the Open Graph Benchmark splits the datasets under real-world settings, we also report the performance on the original OGB validation and test datasets, which can be viewed as real-world environments. As shown in Figure 2g-2i, our framework consistently outperforms the other methods. The improvement can be larger when the bias is more severe. These results again verifies the importance of alleviating selection biases and the effectiveness of our framework.

*5.1.4 Analysis on Hyper-parameters.* To investigate how the proposed two regularizers affect the stability of prediction, we vary the corresponding hyper-parameters, *i.e.*, $\lambda_0$ and $\lambda_1$, on the OGB-Arxiv dataset with training bias ratio $\tau^{tr} = 0.8$. Figure 4 summarizes the results. When increasing the $\lambda_0$ and $\lambda_1$ from a lower rate, we observe a significant stability improvement, *i.e.*, improved Average_Score and reduced Stability_Error. This demonstrates the effectiveness of our framework in the sense of ablation study. In other words, progressively reducing the impact of the two regularizers to a lower
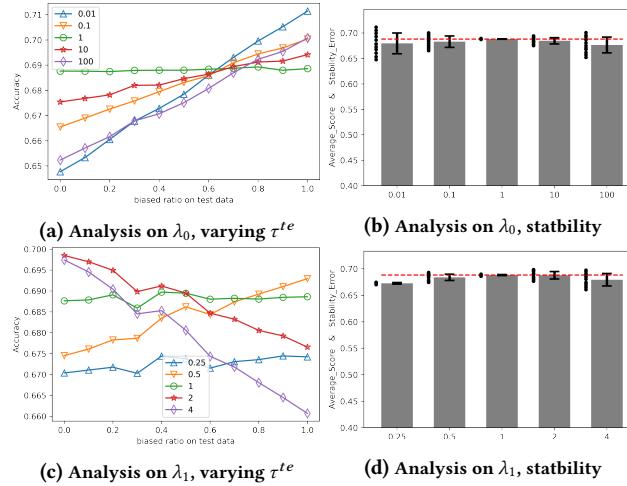
(a) Analysis on $\lambda_0$, varying $\tau^{te}$

(b) Analysis on $\lambda_0$, statbility

(c) Analysis on $\lambda_1$, varying $\tau^{te}$

(d) Analysis on $\lambda_1$, statbility

**Figure 4: Hyper-parameter analysis on the OGB-Arxiv dataset with training bias ratio $\tau^{tr} = 0.8$.**

rate will lead to a performance drop. We also observe that when the impact of the two regularizers becomes increasingly larger, *e.g.*, $\lambda_0 = 10$ or $100, \lambda_1 = 2$ or $4$, there is a performance drop in the model's stability. We attribute this phenomenon to that a large $\lambda$ of one regularizer (global or local) may deteriorate the capability of the other regularizer (local or global) as well as the task-specific capacity, thus leading to unstable predictions. When comparing the global and local regularizers, the effectiveness of local stable learning is less sensitive to the corresponding hyper-parameter $\lambda_0$. In other words, we can easily improve the model stability using the local regularizer without cautiously tuning $\lambda_0$. As for the global regularizer, relatively small or large $\lambda_1$ will harm the stability. This suggests that universally pulling the losses of each node across heterogeneous environments at the *same* rate might not be an optimal solution. Therefore, we leave finding adaptive pulling strategies for different nodes and environments as a promising future work.

## 5.2 Experiments on a Recommendation Dataset

*5.2.1 Experimental Setup.* To demonstrate the efficacy of the proposed framework on real-world applications, we conduct experiments on the user-item bipartite graph of recommender systems, where sample selection bias are arguably ubiquitous [5].

**Collecting a Real-world *Noisy* Dataset.** Specifically, we collect an industrial dataset from one of the world-leading e-commerce platforms from June 11th, 2020, to June 15th, 2020, when an annual product promotion festival is being celebrated. In such a period, the promotion strategies from online shop owners can be various and time-evolving. Therefore, inconsistencies between the users' clicks and their satisfactions naturally exist, leading to a distribution shift from the collected data and the real-world testing environment. We select users and items that have interactions in all five days, which means we do not consider the cold-start setting. We mainly consider click interactions, which is a common setting for the deep candidate generation phase in recommendation. We further select

users that have 200-300 interactions to ensure the quality of data. We split the interactions into several environments according to the day they happen, and the statistics are listed in Table 1. We use data samples from the first day for training and use the remaining for evaluation. We keep the gender and age attributes for users. For users that do not provide these attributes, we use the value predicted by the e-commerce platform. There are 8 age sections provided by the platform, including 1-18, 19-25, 26-30, 31-35, 36-40, 41-50, 51-60, and >=61. We group 1-18 and 19-25 into a holistic section and the remaining into another section. This split results in approximately equally-sized sections.

**Evaluation** We mainly focus on the deep candidate generation stage of recommendation. The user-item bipartite graph consists of user nodes and item nodes. A user is connected to an item when the user interacts with the item. Solely id feature is considered and transformed to dense vectors using a learnable embedding matrix, which is a common practice for deep candidate generation models [1, 37]. The embedding matrix is learned along with the graph recommenders. The collected dataset contains rich semantic attributes, and we discuss user gender/age in experiments, respectively. The probability of a user to be selected can be $P(s_i = 1) = \tau$ if $age \leq 25$ *or gender* $= M$ and $1 - \tau$ if $age > 25$ *or gender* $= F$. $M$ denotes male and $F$ denotes female. We also discuss agnostic selection bias by viewing each of the following days as an individual environment. We incorporate a widely used metric NDCG [1, 37]. Compared to other widely used metrics such as Recall and Hit Ratio, NDCG considers the positions of recommended items. NDCG can be formally written as:

$$\text{DCG@N} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sum_{r \in \hat{\mathcal{I}}_{u,N}} \frac{\mathbb{1}(r \in \mathcal{I}_u)}{\log_2(i_r + 1)} \quad (13)$$

$$\text{NDCG@N} = \frac{\text{DCG@N}}{\text{IDCG@N}} \quad (14)$$

where $\mathcal{U}$ is the set of users, N is the number of recommended items, and $\hat{i}_{u,k}$ indicates the $k$th item recommended for user $u$. $\mathbb{1}$ denotes the indicator function. IDCG@N denotes the ideal discounted cumulative gain and is the maximum possible value of DCG@N. consider the top 100 generated candidates of each model for evaluation, *i.e.*, NDCG@100. We report the percentage score in the results.

**Baselines** We consider the following models as baselines:

- *NGCG* leverages GCN to represent users and items based on the user-item bipartite graph.
- *LightGCN* linearly propagates user/item embeddings on the user-item graph and thus simplifying the NGCF.
- *Stable Graph Recommender*. We build the proposed stable graph recommender based on NGCF. NGCF largely follows the standard GCN model. The proposed stable graph recommender propagates
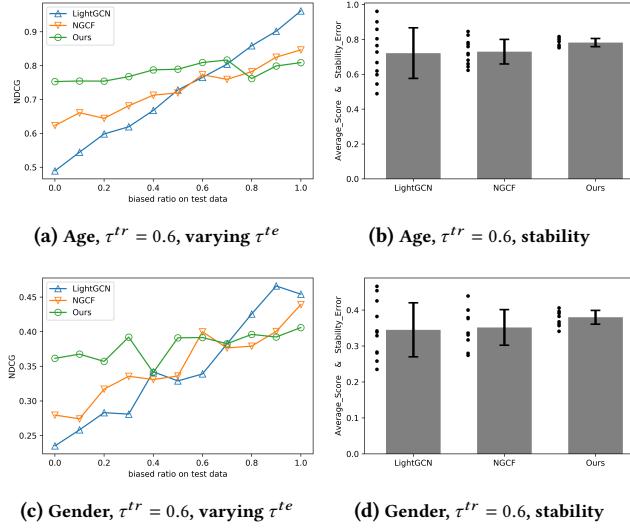
**(a) Age, $\tau^{tr} = 0.6$, varying $\tau^{te}$**

**(b) Age, $\tau^{tr} = 0.6$, stability**



**(c) Gender, $\tau^{tr} = 0.6$, varying $\tau^{te}$**

**(d) Gender, $\tau^{tr} = 0.6$, stability**

**Figure 5: Results on real-world recommendation dataset with distribution shift caused by node attributes.**



**(a) Age, $\tau^{tr} = 0.6$**

**(b) Gender, $\tau^{tr} = 0.6$**

**Figure 6: Results on recommendation dataset with real-world environments (each day as an individual environment).**

embeddings on the user-item bipartite graph as the following:

$$\alpha^e_{ui} = \text{sigmoid}\left(\mathbf{a}^e\left[\mathbf{e}^0_u\|\mathbf{e}^0_i\right]\right) \tag{15}$$

$$\mathbf{e}^{(k+1)}_u = \sigma\left(\mathbf{W}_1\mathbf{e}^{(k)}_u + \sum_{i\in\mathcal{N}_u}\frac{\alpha^e_{ui}}{\sqrt{|\mathcal{N}_u|\,|\mathcal{N}_i|}}\left(\mathbf{W}_1\mathbf{e}^{(k)}_i + \mathbf{W}_2\left(\mathbf{e}^{(k)}_i\odot\mathbf{e}^{(k)}_u\right)\right)\right) \tag{16}$$

$$\mathbf{e}^{(k+1)}_i = \sigma\left(\mathbf{W}_1\mathbf{e}^{(k)}_i + \sum_{u\in\mathcal{N}_i}\frac{\alpha^e_{ui}}{\sqrt{|\mathcal{N}_u|\,|\mathcal{N}_i|}}\left(\mathbf{W}_1\mathbf{e}^{(k)}_u + \mathbf{W}_2\left(\mathbf{e}^{(k)}_u\odot\mathbf{e}^{(k)}_i\right)\right)\right) \tag{17}$$

where $\mathbf{e}^{(k+1)}_u$ and $\mathbf{e}^{(k+1)}_i$ denote the updated user and item embedding after $k$ layers propagation. $\sigma$ denotes a certain nonlinear activation function and is LeakyReLU as NGCF. $\mathcal{N}_u$ denotes the set of interacted items for user $u$ and $\mathcal{N}_i$ denotes the set of interacted users for item $i$. $\mathbf{W}_1$ and $\mathbf{W}_2$ are learnable transformation matrices. $\alpha^e_{ui}$ denotes the importance of interaction $< u, i >$ for representing user $u$ and item $i$. We note that interactions that are consistently important across environments $\mathcal{E}$ are stable properties, as illustrated in Section 4.2.1. We keep a global $\alpha^e_{ui}$ for layers due to its efficiency and do not compute weights per layer. We note that deep candidate generation models, which recall Top K items from a billion-scale item gallery are largely sensitive to model efficiency. For example, one of the contribution of LightGCN is to remove the nonlinearity $\sigma$ in NGCF and thus improving efficiency. We train the stable graph recommender as illustrated in Section 4.2.3.

*5.2.2 Analysis on Stable Prediction.* The testing results with distribution shift caused by node attributes are shown in Figure 5. Specifically, we set training bias rate $\tau^{tr} = 0.6$ and test the model on testing environments with bias rate varying among $\{0.0, 0.1, 0.2, \ldots, 1.0\}$.
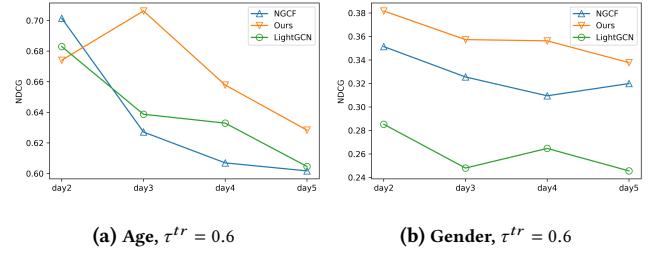
Overall, we observe that our framework achieves significantly more stable results than the other two SOTA graph recommenders. Specifically, the Average_Score is significantly improved with Stability_Error largely reduced. Noteworthy, similar to the findings to those on the graph benchmarks, the stable graph recommender improves the stability mostly by boosting the performance on environments where NGCG/LightGCN performs poorly. For example, as shown in Figure 6b, NGCF/LightGCN achieves about 0.6/0.5 NDCG score when $\tau^{tr} = 0.6$ and $\tau_{te} = 0.0$ while the proposed stable graph recommender significantly boost the NDCG to nearly 0.77 with the same distribution shift. These results further demonstrate the effectiveness of the proposed method on broader real-world applications.

*5.2.3 Evaluation with Real-world Environments.* To evaluate our framework on real-world environments with agnostic selection bias, we construct testing environments from days that follow the day where graph recommenders are trained. To be specific, the model is trained on the recommendation data logged on June 11th, 2020, and we test the trained model on the data logged on June 12-15th, 2020, individually. Therefore, there are four testing environments in total, and we report the results in Figure 6. We observe that there is a consistent performance improvement across heterogeneous real-world environments. These results suggest that our framework is capable of capturing stable properties that should persist across time in node representation learning. For example, a user might click some items due to external distractions, which can hardly represent the user's inherent interest. This phenomenon is known as the *natural noise* [26], and such clicks are unstable properties when representing the user node on the user-item bipartite graph. The proposed stable graph recommender captures non-noisy interactions for users on the user-item bipartite graph and obtains stable/inherent interest representations that are effective across time.

## 6 CONCLUSION AND FUTURE WORK

In this paper, we argue that real-world applications (*e.g.*, recommender systems) are full of selection bias, leading to a distribution shift from the collected graph training data to testing environments. We propose a novel stable prediction framework that permits *locally* stable learning by capturing properties that are stable across environments in node representation learning, and *globally* stable

learning to balance the training of GNNs on different environments. We conduct extensive experiments on newly proposed and traditional graph benchmarks as well as datasets collected on real-world recommender systems concerning shift related to both node labels and attributes. Results demonstrate the capability of the proposed framework on stable prediction on graphs.

We believe that the insights of the proposed framework are inspirational to future developments of stable learning techniques on graphs. Mitigating selection biases can be essential for deploying GNNs on real-world application systems (*e.g.*, confronting the inconsistencies between the training and testing distributions, and achieving fairness for different groups of users), while few works in the literature are proposed for learning stable GNNs. We plan to further investigate such a problem from the perspective of causal discovery/inference, *e.g.*, disentangling the stable/unstable properties for a given neighborhood structure or a hidden variable. Another future direction is to design techniques for application-specific selection biases, *e.g.*, exposure bias in recommender systems. We plan to make further investigations.

# REFERENCES

[1] 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation.. In *SIGIR*.
[2] John Aldrich. 1989. Autonomy. *Oxford Economic Papers* (1989).
[3] Anonymous Authors. 2020. Debiased Graph Neural Networks with Agnostic Label Selection Bias. *OpenReview* (2020).
[4] Steffen Bickel, Michael Brückner, and Tobias Scheffer. 2009. Discriminative Learning Under Covariate Shift. *J. Mach. Learn. Res.* (2009).
[5] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. 2020. Bias and Debias in Recommender System: A Survey and Future Directions. *CoRR* (2020).
[6] Ming Ding, Jie Tang, and Jie Zhang. 2018. Semi-supervised Learning on Graphs with Generative Adversarial Nets.. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018*.
[7] Yingtong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S. Yu. 2020. Enhancing Graph Neural Network-based Fraud Detectors against Camouflaged Fraudsters.. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*.
[8] Rong Duan and Yanghua Xiao. 2019. Enterprise Knowledge Graph From Specific Business Task to Enterprise Knowledge Management.. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*.
[9] Miroslav Dudík, Robert E. Schapire, and Steven J. Phillips. 2005. Correcting sample selection bias in maximum entropy density estimation.. In *Advances in Neural Information Processing Systems 18 [Neural Information Processing Systems, NIPS 2005, December 5-8, 2005, Vancouver, British Columbia, Canada]*.
[10] Changjun Fan, Li Zeng, Yuhui Ding, Muhao Chen, Yizhou Sun, and Zhong Liu. 2019. Learning to Identify High Betweenness Centrality Nodes from Scratch: A Novel Graph Neural Network Approach.. In *CIKM*.
[11] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. 2017. Neural Message Passing for Quantum Chemistry.. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*.
[12] Trygve Haavelmo. 1944. The probability approach in econometrics. *Econometrica: Journal of the Econometric Society* (1944).
[13] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Representation Learning on Graphs: Methods and Applications. *IEEE Data Eng. Bull.* 40, 3 (2017).
[14] Yue He, Peng Cui, Jianxin Ma, Hao Zou, Xiaowei Wang, Hongxia Yang, and Philip S. Yu. 2020. Learning Stable Graphs from Multiple Environments with Selection Bias.. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*.
[15] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open Graph Benchmark: Datasets for Machine Learning on Graphs.. In *NeurIPS*.
[16] Jiayuan Huang, Alexander J. Smola, Arthur Gretton, Karsten M. Borgwardt, and Bernhard Schölkopf. 2006. Correcting Sample Selection Bias by Unlabeled Data.. In *NIPS*.
[17] Yizhu Jiao, Yun Xiong, Jiawei Zhang, and Yangyong Zhu. 2019. Collective Link Prediction Oriented Network Embedding with Hierarchical Graph Attention.. In *CIKM*.

[18] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks.. In *ICLR*.
[19] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. 2019. Predict then Propagate: Graph Neural Networks meet Personalized PageRank.. In *ICLR*.
[20] Kun Kuang, Peng Cui, Susan Athey, Ruoxuan Xiong, and Bo Li. 2018. Stable Prediction across Unknown Environments.. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*.
[21] Kun Kuang, Ruoxuan Xiong, Peng Cui, Susan Athey, and Bo Li. 2020. Stable Prediction with Model Misspecification and Agnostic Distribution Shift.. In *AAAI*.
[22] Xiangsheng Li, Maarten de Rijke, Yiqun Liu, Jiaxin Mao, Weizhi Ma, Min Zhang, and Shaoping Ma. 2020. Learning Better Representations for Neural Information Retrieval with Graph Information.. In *CIKM*.
[23] Anqi Liu and Brian D. Ziebart. 2014. Robust Classification Under Sample Selection Bias.. In *NeurIPS*.
[24] Sharon L Lohr. 2009. *Sampling: design and analysis*. Nelson Education.
[25] Guixiang Ma, Nesreen K. Ahmed, Theodore L. Willke, Dipanjan Sengupta, Michael W. Cole, Nicholas B. Turk-Browne, and Philip S. Yu. 2019. Deep Graph Similarity Learning for Brain Data Analysis.. In *CIKM*.
[26] Michael P. O'Mahony, Neil J. Hurley, and Guenole C. M. Silvestre. 2006. Detecting noise in recommender system databases.. In *Proceedings of the 11th International Conference on Intelligent User Interfaces, IUI 2006, Sydney, Australia, January 29 - February 1, 2006*.
[27] Zhiqiang Pan, Fei Cai, Wanyu Chen, Honghui Chen, and Maarten de Rijke. 2020. Star Graph Neural Networks for Session-based Recommendation.. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*.
[28] Judea Pearl. 2009. *Causality*. Cambridge university press.
[29] Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. 2016. *Causal inference in statistics: A primer*. John Wiley & Sons.
[30] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. 2016. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* (2016).
[31] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2009. The Graph Neural Network Model. *IEEE Trans. Neural Networks* (2009).
[32] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. 2008. Collective Classification in Network Data. *AI Mag.* (2008).
[33] Hidetoshi Shimodaira. 2000. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference* (2000).
[34] Xianfeng Tang, Huaxiu Yao, Yiwei Sun, Yiqi Wang, Jiliang Tang, Charu C. Aggarwal, Prasenjit Mitra, and Suhang Wang. 2020. Investigating and Mitigating Degree-Related Biases in Graph Convoltuional Networks.. In *CIKM*.
[35] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks.. In *ICLR*.
[36] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jing Jiang. 2017. MGAE: Marginalized Graph Autoencoder for Graph Clustering.. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, Singapore, November 06 - 10, 2017*.
[37] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering.. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019*.
[38] Xiao Wang, Meiqi Zhu, Deyu Bo, Peng Cui, Chuan Shi, and Jian Pei. 2020. AM-GCN: Adaptive Multi-channel Graph Convolutional Networks.. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*.
[39] Felix Wu, Jr. Souza, Amauri H., Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Q. Weinberger. 2019. Simplifying Graph Convolutional Networks.. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*.
[40] Man Wu, Shirui Pan, Lan Du, Ivor W. Tsang, Xingquan Zhu, and Bo Du. 2019. Long-short Distance Aggregation Networks for Positive Unlabeled Graph Learning.. In *CIKM*.
[41] Xiancheng Xie, Yun Xiong, Philip S. Yu, and Yangyong Zhu. 2019. EHR Coding with Multi-scale Feature Attention and Structured Knowledge Graph Propagation.. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*.
[42] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In *ICLR*.
[43] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation Learning on Graphs with Jumping Knowledge Networks.. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*.

[44] Deqing Yang, Jingrui He, Huazheng Qin, Yanghua Xiao, and Wei Wang. 2015. A Graph-based Recommendation across Heterogeneous Domains.. In *Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM 2015, Melbourne, VIC, Australia, October 19 - 23, 2015.*

[45] Yuting Ye, Xuwu Wang, Jiangchao Yao, Kunyang Jia, Jingren Zhou, Yanghua Xiao, and Hongxia Yang. 2019. Bayes EMbedding (BEM): Refining Representation by Integrating Knowledge Graphs and Behavior-specific Networks.. In *CIKM.*

[46] Tianqi Zhang, Yun Xiong, Jiawei Zhang, Yao Zhang, Yizhu Jiao, and Yangyong Zhu. 2020. CommDGI: Community Detection Oriented Deep Graph Infomax.. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020.*

[47] Huan Zhao, Lanning Wei, and Quanming Yao. 2020. Simplifying Architecture Search for Graph Neural Network.. In *Proceedings of the CIKM 2020 Workshops co-located with 29th ACM International Conference on Information and Knowledge*

*Management (CIKM 2020), Galway, Ireland, October 19-23, 2020.*

[48] Tianxiang Zhao, Xianfeng Tang, Xiang Zhang, and Suhang Wang. 2020. Semi-Supervised Graph-to-Graph Translation.. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020.*

[49] Fan Zhou, Tengfei Li, Haibo Zhou, Hongtu Zhu, and Ye Jieping. 2019. Graph-Based Semi-Supervised Learning with Non-ignorable Non-response. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc.

[50] Xuliang Zhu, Xin Huang, Byron Choi, and Jianliang Xu. 2020. Top-k Graph Summarization on Hierarchical DAGs.. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020.*