# **Neural Bandit with Arm Group Graph**

Yunzhe Qi University of Illinois at Urbana-Champaign yunzheq2@illinois.edu Yikun Ban University of Illinois at Urbana-Champaign yikunb2@illinois.edu Jingrui He University of Illinois at Urbana-Champaign jingrui@illinois.edu

#### **ABSTRACT**

Contextual bandits aim to identify among a set of arms the optimal one with the highest reward based on their contextual information. Motivated by the fact that the arms usually exhibit group behaviors and the mutual impacts exist among groups, we introduce a new model, Arm Group Graph (AGG), where the nodes represent the groups of arms and the weighted edges formulate the correlations among groups. To leverage the rich information in AGG, we propose a bandit algorithm, AGG-UCB, where the neural networks are designed to estimate rewards, and we propose to utilize graph neural networks (GNN) to learn the representations of arm groups with correlations. To solve the exploitation-exploration dilemma in bandits, we derive a new upper confidence bound (UCB) built on neural networks (exploitation) for exploration. Furthermore, we prove that AGG-UCB can achieve a near-optimal regret bound with over-parameterized neural networks, and provide the convergence analysis of GNN with fully-connected layers which may be of independent interest. In the end, we conduct extensive experiments against state-of-the-art baselines on multiple public data sets, showing the effectiveness of the proposed algorithm.

# **CCS CONCEPTS**

• Information systems  $\rightarrow$  Personalization; • Theory of computation  $\rightarrow$  Online learning algorithms.

#### **KEYWORDS**

Contextual Bandits; Online Learning; Graph Neural Networks

# ACM Reference Format:

Yunzhe Qi, Yikun Ban, and Jingrui He. 2022. Neural Bandit with Arm Group Graph. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22), August 14–18, 2022, Washington, DC, USA. ACM, New York, NY, USA, 15 pages. https://doi.org/10.1145/3534678.

#### 1 INTRODUCTION

Contextual bandits are a specific type of multi-armed bandit (MAB) problem where the learner has access to the contextual information (contexts) related to arms at each round, and the learner is required to make recommendations based on past contexts and received rewards. A variety of models and algorithms have been proposed

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '22, August 14–18, 2022, Washington, DC, USA © 2022 Association for Computing Machinery. ACM ISBN 978-1-4503-9385-0/22/08...\$15.00 https://doi.org/10.1145/3534678.3539312

and successfully applied on real-world problems, such as online content and advertising recommendation [28, 40], clinical trials [15, 36] and virtual support agents [31].

In this paper, we focus on exploiting the accessible arm information to improve the performance of bandit algorithms. Among different types of contextual bandit algorithms, upper confidence bound (UCB) algorithms have been proposed to balance between exploitation and exploration [3, 11, 34]. For conventional UCB algorithms, they are either under the "pooling setting" [11] where one single UCB model is applied for all candidate arms, or the "disjoint setting" [28] where each arm is given its own estimator without the collaboration across different arms. Both settings have their limitations: applying only one single model may lead to unanticipated estimation error when some arms exhibit distinct behaviors [40, 41]; on the other hand, assigning each arm its own estimator neglects the mutual impacts among arms and usually suffers from limited user feedback [5, 18].

To deal with this challenge, adaptively assigning UCB models to arms based on their group information can be an ideal strategy, i.e., each group of arms has one estimator to represent its behavior. This modeling strategy is linked to "arm groups" existing in real-world applications. For example, regarding the online movie recommendation scenario, the movies (arms) with the same genre can be assigned to one (arm) group. Another scenario is the drug development, where given a new cancer treatment and a patient pool, we need to select the best patient on whom the treatment is most effective. Here, the patients are the arms, and they can be naturally grouped by their non-numerical attributes, such as the cancer types. Such group information is easily accessible, and can significantly improve the performance of bandit algorithms. Although some works [12, 30] have been proposed to leverage the arm correlations, they can suffer from two common limitations. First, they rely on the assumption of parametric (linear / kernel-based) reward functions, which may not hold in real-world applications [46]. Second, they both neglect the correlations among arm groups. We emphasize that the correlations among arm groups also play indispensable roles in many decision-making scenarios. For instance, in online movie recommendation, with each genre being a group of movies, the users who like "adventure movies" may also appreciate "action movies". Regarding drug development, since the alternation of some genes may lead to multiple kinds of tumors [32], different types of cancer can also be correlated to some extent.

To address these limitations, we first introduce a novel model, AGG (Arm Group Graph), to formulate non-linear reward assumptions and arm groups with correlations. In this model, as arm attributes are easily accessible (e.g., movie's genres and patient's cancer types), the arms with the same attribute are assigned into one group, and represented as one node in the graph. The weighted edge between two nodes represents the correlation between these

two groups. In this paper, we assume the arms from the same group are drawn from one unknown distribution. This also provides us with an opportunity to model the correlation of two arm groups by modeling the statistical distance between their associated distributions. Meantime, the unknown non-parametric reward mapping function can be either linear or non-linear.

Then, with the arm group graph, we propose the AGG-UCB framework for contextual bandits. It applies graph neural networks (GNNs) to learn the representations of arm groups with correlations, and neural networks to estimate the reward functions (exploitation). In particular, with the collaboration across arm groups, each arm will be assigned with the group-aware arm representation learned by GNN, which will be fed into a fully-connected (FC) network for the estimation of arm rewards. To deal with the exploitation-exploration dilemma, we also derive a new upper confidence bound based on network-gradients for exploration. By leveraging the arm group information and modeling arm group correlations, our proposed framework provides a novel arm selection strategy for dealing with the aforementioned challenges and limitations. Our main contributions can be summarized as follows:

- First, motivated by real-world applications, we introduce a new graph-based model in contextual bandits to leverage the available group information of arms and exploit the potential correlations among arm groups.
- Second, we propose a novel UCB-based neural framework called AGG-UCB for the graph-based model. To exploit the relationship of arm groups, AGG-UCB estimates the arm group graph with received contexts on the fly, and utilizes GNN to learn group-aware arm representations.
- Third, we prove that AGG-UCB can achieve a near-optimal regret bound in the over-parameterized neural works, and provide convergence analysis of GNN with fully-connected layers, which may be of independent interest.
- Finally, we conduct experiments on publicly available real data sets, and demonstrate that our framework outperforms state-of-the-art techniques. Additional studies are conducted to understand the properties of the proposed framework.

The rest of this paper is organized as following. In Section 2, we briefly discuss related works. Section 3 introduces the new problem settings, and details of our proposed framework AGG-UCB will be presented in Section 4. Then, we provide theoretical analysis for AGG-UCB in Section 5. After presenting experimental results in Section 6, we finally conclude the paper in Section 7. Due to the page limit, readers may refer to our arXiv version of the paper for the supplementary contents (https://arxiv.org/abs/2206.03644).

# 2 RELATED WORKS

In this section, we briefly review the related work on contextual bandits. Lin-UCB [11] first formulates the reward estimation through a linear regression with the received context and builds a confidence bound accordingly. Kernel-UCB [34] further extends the reward mapping to the Reproducing Kernel Hilbert Space (RKHS) for the reward and confidence bound estimation under non-linear settings. Besides, there are algorithms under the non-linear settings. Similarly, CGP-UCB [27] models the reward function through a Gaussian process. GCN-UCB [33] applies the GNN model to learn each

context an embedding for the linear regression. Then, Neural-UCB [46] proposes to apply FC neural network for reward estimations and derive a confidence bound with the network gradient, which is proved to be a success, and similar ideas has been applied to some other models [4, 6, 45]. [8] assigns another FC neural network to learn the confidence ellipsoid for exploration. Yet, as these works consider no collaboration among estimators, they may suffer from the performance bottleneck in the introduction.

To collaborate with different estimators for contextual bandits, various approaches are proposed from different perspectives. User clustering algorithms [5, 7, 19, 29] try to cluster user with alike preferences into user groups for information sharing while COFIBA [30] additionally models the relationship of arms. Then, KMTL-UCB [12] extends Kernel-UCB to multi-task learning settings for a refined reward and confidence bound estimation. However, these works may encounter with performance bottlenecks as they incline to make additional assumptions on the reward mapping by applying parametric models and neglect the available arm group information.

GNNs [17, 25, 26, 39] are a kind of neural models that deal with tasks on graphs, such as community detection [44], recommender systems [38] and modeling protein interactions [16]. GNNs can learn from the topological graph structure information and the input graph signal simultaneously, which enables AGG-UCB to cooperate with different arm groups by sharing information over the arm group neighborhood.

# 3 PROBLEM DEFINITION AND NOTATION

We suppose a fixed pool  $C = \{1, \ldots, N_c\}$  for arm groups with the number of arm groups being  $|C| = N_c$ , and assume each arm group  $c \in C$  is associated with an unknown fixed context distribution  $\mathcal{D}_c$ . At each time step t, we will receive a subset of groups  $C_t \subseteq C$ . For each group  $c \in C_t$ , we will have the set of sampled arms  $X_{c,t} = \{x_{c,t}^{(1)}, \cdots x_{c,t}^{(n_{c,t})}\}$  with the size of  $|X_{c,t}| = n_{c,t}$ . Then,  $\forall i \in [n_{c,t}] = \{1, \ldots, n_{c,t}\}$ , we suppose  $x_{c,t}^{(i)} \sim \mathcal{D}_c$  with the dimensionality  $x_{c,t}^{(i)} \in \mathbb{R}^{d_x}$ . Therefore, in the t-th round, we receive

$$\{X_{c,t}|c\in C_t\} \text{ and } X_{c,t} = \{x_{c,t}^{(1)}, \cdots x_{c,t}^{(n_{c,t})}\}, \forall c\in C_t.$$
 (1)

With  $W^* \in \mathbb{R}^{N_c \times N_c}$  being the unknown affinity matrix encoding the true arm group correlations, the true reward  $r_{c,t}^{(i)}$  for arm  $x_{c,t}^{(i)}$  is defined as

$$r_{c,t}^{(i)} = h(\mathbf{W}^*, \mathbf{x}_{c,t}^{(i)}) + \epsilon_{c,t}^{(i)}$$
 (2)

where  $h(\cdot)$  represents the unknown reward mapping function, and  $\epsilon$  is the zero-mean Gaussian noise. For brevity, let  $x_t$  be the arm we select in round t and t be the corresponding received reward.

Our goal is recommending arm  $x_t$  (with reward  $r_t$ ) at each time step t to minimize the cumulative pseudo-regret  $R(T) = \sum_{t=1}^T \mathbb{E}\left[(r_t^* - r_t)\right]$  where  $\mathbb{E}[r_t^*] = \max_{(c \in C_t, i \in [n_{c,t}])} \left[h(W^*, \mathbf{x}_{c,t}^{(i)})\right]$ . At each time step t, the overall set of received contexts is defined as  $X_t = \{\mathbf{x}_{c,t}^{(i)}\}_{c \in C_t, i \in [n_{c,t}]}$ . Note that one arm is possibly associated with multiple arm groups, such as a movie with multiple genres. In other words, for some  $c, c' \in C_t$ , we may have  $X_{c,t} \cap X_{c',t} \neq \emptyset$ .

In order to model the arm group correlations, we maintain an undirected graph  $G_t = (V, E, W_t)$  at each time step t, where each arm group from C is mapped to a corresponding node in node set V.

#### **ALGORITHM 1:** AGG-UCB

```
1 Input: Number of rounds T, exploration parameter \gamma,
   regularization parameter \lambda, network width m, network
   depth L, neighborhood size k.
```

```
2 Output: Arm recommendation x_t for each time step t.
```

```
3 Initialization: Initialize the arm group graph as a
    connected graph G_1 = (V, E, W_1). Initialize gradient matrix
    Z_0 = \lambda I. Initialize parameter \Theta_0 for the model f(\mathcal{G}, X; \Theta_0).
```

```
4 for t = 1, 2, ..., T do
         Receive a set of arm contexts X_t = \{x_{c,t}^{(i)}\}_{c \in C_t, i \in [n_{c,t}]}.
          Embed the arm set X_t into \widetilde{X}_t w.r.t. Eq.4.
6
         for each embedded arm \widetilde{X}_{c,t}^{(i)} \in \widetilde{X}_t do
                Obtain the point estimate \widehat{r}_{c,t}^{(i)} = f(\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}; \Theta_{t-1}).
8
```

Obtain network gradient  $g_{c,t}^{(i)} \leftarrow \nabla_{\Theta} f(\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}; \Theta_{t-1}).$  Calculate confidence bound as  $\widehat{i}_{c,t}^{(i)} = \sqrt{g_{c,t}^{(i)} {}^{\mathsf{T}} Z_{t-1} g_{c,t}^{(i)} / m}.$  end

11

10

12

14

15

Recommend  $\widetilde{X}_t = \operatorname{argmax}_{\widetilde{X}_{c,t}^{(i)} \in \widetilde{X}_t} (\widehat{r}_{c,t}^{(i)} + \gamma \cdot \widehat{i}_{c,t}^{(i)})$  with the received reward represented as  $r_t$ .

Calculate arm group distances w.r.t. Eq.3, and update 13 the arm group graph  $G_t$  to  $G_{t+1}$ .

Update the model parameter  $\Theta_{t-1}$  to  $\Theta_t$  according to Algorithm 2.

Retrieve the  $\widetilde{X}_t$ 's gradient vector  $g_t$ , and update gradient matrix  $Z_t = Z_{t-1} + q_t \cdot q_t^{\mathsf{T}}$ .

16 end

Then,  $E = \{e(c_i, c_j)\}_{c_i, c_i \in C}$  is the set of edges, and  $W_t$  represents the set of edge weights. Note that by definition,  $G_t$  will stay as a fully-connected graph, and the estimated arm group correlations are modeled by the edge weights connecting pairs of nodes. For a node  $v \in V$ , we denote the augmented k-hop neighborhood  $N_k(v) =$  $\mathcal{N}_k(v) \cup \{v\}$  as the union node set of its *k*-hop neighborhood  $\mathcal{N}_k(v)$ and node v itself. For the arm group graph  $G_t$ , we denote  $A_t \in$  $\mathbb{R}^{N_c \times N_c}$  as the adjacency matrix (with added self-loops) of the given arm group graph and  $D_t \in \mathbb{R}^{N_c \times N_c}$  as its degree matrix. For the notation consistency, we will apply a true arm group graph  $\mathcal{G}^*$ instead of  $W^*$  in Eq. 2 to represent the true arm group correlation.

# PROPOSED AGG-UCB FRAMEWORK

In this section, we start with an overview to introduce our proposed AGG-UCB framework. Then, we will show our estimation method of arm group graph before mentioning the related group-aware arm embedding. Afterwards, the two components of our proposed framework, namely the aggregation module and the reward-estimation module, will be presented.

### Overview of AGG-UCB Framework

In Algorithm 1, we present the pseudo-code of proposed AGG-UCB framework. At each time step  $t \in [T]$ , AGG-UCB would receive a set of input arm contexts  $X_t = \{x_{c,t}^{(i)}\}_{c \in C_t, i \in [n_{c,t}]}$  (line 5).

#### **ALGORITHM 2:** Model Training

- 1 **Input:** Initial parameter  $Θ_0$ , step size η, training steps J, network width m. Updated arm group graph  $\mathcal{G}_{t+1}$ . Selected embedded contexts  $\{X_{\tau}\}_{\tau=1}^{t}$ .
- 2 **Output:** Updated model parameter  $\Theta_t$ .
- $\Theta_t^0 \leftarrow \Theta_0$ .
- 4 Let  $\mathcal{L}(\Theta) = \frac{1}{2} \sum_{\tau=1}^{t} |f(\mathcal{G}_{t+1}, \widetilde{X}_{\tau}; \Theta) r_{\tau}|^2$

- 8 Return new parameter  $\Theta_t^J$ .

Then, we embed the arm set  $X_t$  to  $\widetilde{X}_t$  based on **Eq.**4 from Subsection 4.3 (line 6). For each embedded arm  $\widetilde{X} \in \widetilde{X}_t$ , its estimated reward  $\hat{r}$  and confidence bound  $\hat{i}$  would be calculated (line 8-10) with the model  $f(\cdot)$  in Subsection 4.4. After recommending the best arm  $X_t$ (line 12) and receiving its true reward  $r_t$ , we update the current arm group graph  $G_t$  based on Subsection 4.2 (line 13). Then, the model parameters  $\Theta_{t-1}$  will be trained based on Algorithm 2 (line 14), and we incrementally update the gradient matrix to  $Z_t =$  $Z_{t-1} + g_t \cdot g_t^{\mathsf{T}}$  with the gradient vector  $g_t$  of model  $f(\cdot)$  given the selected arm  $\widetilde{X}_t$  (line 15).

The steps from Algorithm 2 demonstrate our training process for AGG-UCB parameters. With the updated arm group graph  $G_{t+1}$ and the past embedded arm contexts  $\{\widetilde{X}_{ au}\}_{ au=1}^t$  until current time step t, we define the loss function as the straightforward quadratic loss function (line 4). Finally, we run gradient descent (GD) for I steps to derive the new model parameters  $\Theta_t$  (lines 5-7) based on the initial parameters  $\Theta_0$  (initialized in Subsection 4.5). Next, we proceed to introduce the detail of framework components.

#### 4.2 Arm Group Graph Estimation

Recall that at time step t, we model the similar arms into an arm group graph  $G_t = (V, E, W_t)$  where the nodes V are corresponding to the arm groups from C and edges weights  $W_t$  formulate the correlations among arm groups. Given two nodes  $\forall c, c' \in C$ , to measure the similarity between them, inspired by the kernel mean embedding in the multi-task learning settings [9, 12], we define edge weight between c and c' as:

$$w^*(c,c') = \exp(-\|\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c} [\phi_{k_C}(\mathbf{x})] - \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}_{c'}} [\phi_{k_C}(\mathbf{x}')]\|^2 / \sigma_s)$$

where  $\phi_{k_{\mathcal{G}}}(\cdot)$  is the induced feature mapping of a given kernel  $k_{\mathcal{G}}$ , e.g., a radial basis function (RBF) kernel. Unfortunately,  $\forall c \in C$ ,  $\mathcal{D}_c$  is unknown. Therefore, we update the edge weight based on the empirical estimation of arm group correlations. Here, let  $X_c^t =$  $\{x_{c,\tau}^{(i)}\}_{\tau\in[t],i\in[n_{c,\tau}]}$  represent the set of all arm contexts from group  $c \in C$  up to time step t. We define the arm similarity measurement between arms  $c, c' \in C$  through a Gaussian-like kernel as

$$w_t(c, c') = \exp(-\|\Psi_t(\mathcal{D}_c) - \Psi_t(\mathcal{D}_{c'})\|^2 / \sigma_s)$$
 (3)

where  $\Psi_t(\mathcal{D}_c) = \frac{1}{|\mathcal{X}_c^t|} \sum_{x \in \mathcal{X}_c^t} k_{\mathcal{G}}(\cdot, x)$  denotes the kernel mean estimation of  $\mathcal{D}_c$  with a given kernel  $k_G$ ; and  $\sigma_s$  refers to the bandwidth. Then, at time step t and  $\forall c, c' \in C$ , we update the corresponding weight of edge e(c, c') in the weight set  $W_t$  with  $w_t(c, c')$ .

# **Group-Aware Arm Embedding**

To conduct the aggregation operations of GNN, we reconstruct a matrix for each arm context vector. Recall that for an arm group  $c \in C$ , if  $c \in C_t$ , we receive the contexts  $x_{c,t}^{(i)} \in \mathbb{R}^{d_x}$ ,  $i \in [n_{c,t}]$  at time step t. Then, the reconstructed matrix for  $\boldsymbol{x}_{c.t}^{(i)}$  is defined as

$$\widetilde{X}_{c,t}^{(i)} = \begin{pmatrix} (x_{c,t}^{(i)})^{\mathsf{T}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & (x_{c,t}^{(i)})^{\mathsf{T}} & \cdots & \mathbf{0} \\ \vdots & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & (x_{c,t}^{(i)})^{\mathsf{T}} \end{pmatrix} \in \mathbb{R}^{N_c \times d_{\widetilde{x}}}$$
(4)

where  $d_{\widetilde{X}} = d_X \cdot N_c$  is the column dimension of  $\widetilde{X}_{c.t}^{(i)}$ . Here, for the c'-th row in matrix  $\widetilde{X}_{c,t}^{(i)}$ , the  $((c'-1)\cdot d_x+1)$ -th to the  $(c'\cdot d_x)$ -th entries are the transposed original arm context  $(x_{c,t}^{(i)})^\intercal$ , while the other entries are zeros. Receiving a set of arm contexts  $X_t$ , we derive the corresponding embedded arm set as  $\widetilde{X}_t = \{\widetilde{X}_{c,t}^{(i)}\}_{c \in C_t, i \in [n_{c,t}]}$ .

4.3.1 Aggregation of arm group representations. To leverage the estimated arm group graph for downstream reward estimations, we propose to aggregate over the arm group neighborhood for a more comprehensive arm representation through the GNN-based module, named as group-aware arm representation. It has been proven that the local averaging operation on the graph neighborhood can be deemed as applying the low-pass filter on the corresponding node features [23, 39], which would give locally smooth node features within the same neighborhood. Inspired by the SGC model [39], we propose to aggregate over the *k*-hop arm group neighborhood  $\widetilde{\mathcal{N}}_k(\cdot)$ for incorporating arm group correlations to obtain the aggregated group-aware embedding for an embedded arm  $\widetilde{X}_{c,t}^{(i)}$ , denoted by

$$H_{gnn} = \sqrt{\frac{1}{m}} \cdot \sigma(S_t^k \cdot \widetilde{X}_{c,t}^{(i)} \Theta_{gnn}) \in \mathbb{R}^{N_c \times m}$$
 (5)

where  $S_t = D_t^{-\frac{1}{2}} A_t D_t^{-\frac{1}{2}}$  is the symmetrically normalized adjacency matrix, and we have

$$\Theta_{gnn} = \left( egin{array}{c} \Theta_{gnn}^1 \in \mathbb{R}^{d_x imes m} \\ dots \\ \Theta_{gnn}^{c'} \in \mathbb{R}^{d_x imes m} \\ dots \\ \Theta_{gnn}^{N_c} \in \mathbb{R}^{d_x imes m} \end{array} 
ight) \in \mathbb{R}^{d_{\widetilde{x}} imes m}.$$

being the trainable weight matrix with width m. Here,  $\sigma(\cdot)$  denotes the non-linear activation function, which is added after the aggregation operation to alleviate potential concerns when the contexts are not linearly separable [23]. Note that the c'-th row of  $(\widetilde{X}_{c,t}^{(i)} \cdot \Theta_{ann})$ , denoted by  $[\widetilde{X}_{c,t}^{(i)} \cdot \Theta_{gnn}]_{c',:}$ , is the hidden representation of arm x in terms of c'-th arm group in C. Then, these hidden representations will then be aggregated over  $\widetilde{\mathcal{N}}_k(c), c \in C$  by multiplying with  $S_t^k$ to derive the aggregated arm representation for x, i.e.,  $H_{qnn}(x)$ .

4.3.2 Incorporating initial embedded contexts. Moreover, solely aggregating information from neighbors through the GNN-based models can lead to "over-smoothing" problems [42, 43]. Aggregating

from the node neighborhood will end up with identical representations for all the nodes if they form an isolated complete sub-graph, which may not correctly reflect the relationship among these nodes in real-world applications. Therefore, we propose to apply skipconnections to address this potential problem by combining the initial contexts with the aggregated hidden features. Similar ideas have been applied to boost the performance of neural models. For instance, JK-Net [42] and GraphSAGE [20] concatenate hidden features from different levels of node neighborhoods; and ResNet [21] adopts additive residual connections.

Putting these two parts together and setting  $d' = m + d_{\widetilde{X}}$ , we then have  $H_0 \in \mathbb{R}^{N_c \times d'}$  as the output group-aware arm representation for  $\widetilde{X}_{c,t}^{(i)}$ , represented by

$$H_0 = f_{gnn}(\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}; \Theta_{gnn}) = [\sigma(S_t^k \cdot \widetilde{X}_{c,t}^{(i)} \Theta_{gnn}); \widetilde{X}_{c,t}^{(i)}]$$
(6)

where  $[\cdot;\cdot]$  refers to the column-wise concatenation of matrices.

# **Reward Estimation Module**

In this subsection, we estimate the rewards with a FC network of Llayers and width m, based on group-aware arm representation  $H_0$ .

4.4.1 Reward and confidence bound estimation. Here, let  $\Theta_{fc}$  =  $\{\Theta_l\}_{l\in[L]}$  be the set of trainable weight matrices of a fully-connected network, where the specifications are:  $\Theta_1 \in \mathbb{R}^{d' \times m}$ ,  $\Theta_L \in \mathbb{R}^m$  and  $\Theta_l \in \mathbb{R}^{m \times m}, \forall l \in \{2, \dots, L-1\}$ . Then, given the group-aware representation  $H_0$ , we have the reward estimation module as follows

$$H_{l} = \sqrt{\frac{1}{m}} \cdot \sigma(H_{l-1} \cdot \Theta_{l}), \ l \in \{1, \dots, L-1\},$$

$$\widehat{r}_{all} = f_{fc}(H_{0}; \Theta_{fc}) = \sqrt{\frac{1}{m}} \cdot H_{L-1} \cdot \Theta_{L}$$
(7)

where  $\hat{r}_{all} \in \mathbb{R}^{N_c}$  represents the point-estimation vector for the received contexts embedding  $H_0$  with respect to all the arms groups. Given that the arm  $\widetilde{x}_{c,t}^{(i)}$  belonging to c-th group, we will then have the reward estimation  $\widehat{r}_{c,t}^{(i)} = [\widehat{r}_{all}]_c \in \mathbb{R}$  for the embedded context matrix  $\widetilde{X}_{c,t}^{(i)}$ , which is the c-th element of  $\widehat{r}_{all}$ . Finally, combining the aggregation module with the reward estimates

mation module, given arm group graph  $\mathcal{G}_t$  at time step t, the reward estimation for the embedded arm  $\widetilde{X}_{c,t}^{(i)}$  (i.e., the reward estimation given its arm group) can be represented as

$$\widehat{r}_{c,t}^{(i)} = f(\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}; \Theta) = \left[ \left( f_{fc}(\cdot; \Theta_{fc}) \circ f_{gnn}(\cdot; \Theta_{gnn}) \right) (\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}) \right]_c$$

Setting  $p = (2N_c \cdot d) \cdot m + (L-1) \cdot m^2 + m$ , we have  $\Theta \in \mathbb{R}^p$  being the set of all the parameters from these two modules.

4.4.2 Arm pulling mechanism. We obtain confidence bounds for the point estimation with the network gradients as  $\hat{i} = \sqrt{q^{\top} \cdot Z_{t-1} \cdot q/m}$ where  $g = \nabla_{\Theta} f(\mathcal{G}_t, \widetilde{X}_{c,t}^{(i)}; \Theta) \in \mathbb{R}^p$  is the gradient vector, and  $Z_{t-1} = I + \sum_{\tau=1}^{t-1} g_{\tau} \cdot g_{\tau}^{\mathsf{T}}$  with  $g_{\tau}$  being the gradient vector of the embedded arm which is selected at step  $\tau \in \{1, ..., t-1\}$ . After obtaining the reward and confidence bound estimations for all embedded arm in set  $\widetilde{X}_t$ , we choose the best arm as  $\widetilde{X}_t = \operatorname{argmax}_{\widetilde{X}_{c,t} \in \widetilde{\mathcal{X}}_t} (\widehat{r}_{c,t}^{(i)} + \gamma \cdot \widehat{l}_{c,t}^{(i)})$  where  $\gamma$  is the exploration parameter, and the theoretical upper confidence bound will be given in

Section 5. Note that based on our problem definition (Section 3), one arm may associate with multiple arm groups. Here, we will separately estimate rewards and confidence bounds of each arm group it belongs to, and consider them as different arms for selection.

#### 4.5 Model Initialization

For the aggregation module weight matrix  $\Theta_{gnn}$ , each of its entries is sampled from the Gaussian distribution N(0,1). Similarly, the parameters from the first L-1 reward estimation module layers  $([\Theta_1,\ldots,\Theta_{L-1}])$  are also sampled from N(0,1). For the final (L-th) layer, its weight matrix  $\Theta_L$  is initialized by drawing the entry values from the Gaussian distribution N(0,1/m).

# 5 THEORETICAL ANALYSIS

In this section, we provide the theoretical analysis for our proposed framework. For the sake of analysis, at each time step t, we assume each arm group  $c \in C$  would receive one arm  $x_{c,t}$ , which makes  $|X_1^t| = \cdots = |X_{N_c}^t| = t$ . We also apply the adjacency matrix  $A_t$  instead of  $S_t$  for aggregation, and set its elements  $[A_t]_{ij} = \frac{1}{t \cdot N_c} \sum_{\tau=1}^t \phi_{k_{\mathcal{G}}}(x_{c_i,\tau})^\intercal \phi_{k_{\mathcal{G}}}(x_{c_j,\tau})$  for arm group similarity between group  $c_i, c_j \in C$ . Here,  $\phi_{k_{\mathcal{G}}}(\cdot)$  is the kernel mapping given an RBF kernel  $k_{\mathcal{G}}$ . With  $\mathcal{G}^*$  being the unknown true arm group graph, its adjacency matrix elements are  $[A^*]_{ij} = \frac{1}{N_c} \mathbb{E}_{x_i \sim \mathcal{D}_{c_i}, x_j \sim \mathcal{D}_{c_j}} (\phi_{k_{\mathcal{G}}}(x_i)^\intercal \phi_{k_{\mathcal{G}}}(x_j))$ . Note that the norm of adjacency matrices  $||A^*||_2, ||A_t||_2 \le 1$  since  $\langle \phi_{k_{\mathcal{G}}}(x), \phi_{k_{\mathcal{G}}}(x') \rangle \le 1$  for any  $x, x' \in \mathbb{R}^{d_x}$ , which makes it feasible to aggregate over k-hop neighborhood without the explosion of eigenvalues. Before presenting the main results, we first introduce the following background.

LEMMA 5.1 ([4, 46]). For any  $t \in [T]$ , given arm  $\mathbf{x} \in \mathbb{R}^{d_x}$  satisfying  $\|\mathbf{x}\|_2 = 1$  and its embedded context matrix  $\widetilde{\mathbf{X}}$ , there exists  $\boldsymbol{\Theta}_{t-1}^* \in \mathbb{R}^p$  at time step t, and a constant S > 0, such that

$$h(\mathcal{G}^*, \widetilde{X}) = \langle g(\mathcal{G}^*, \widetilde{X}; \Theta_{t-1}), \Theta_{t-1}^* - \Theta_0 \rangle$$
 (8)

where  $\|\Theta_{t-1}^* - \Theta_0\|_2 \le S/\sqrt{m}$ ,  $\forall t \in [T]$ , and  $\mathcal{G}^*$  stands for the unknown true underlying arm group graph.

Note that with sufficient network width m, we will have  $\Theta_{t-1}^* \approx \Theta_0^*$ ,  $\forall t \in [T]$ , and we will include more details in the full version of the paper. Following the analogous ideas from previous works [6, 46], this lemma formulates the expected reward as a linear function parameterized by the difference between randomly initialized network parameter  $\Theta_0$  and the parameter  $\Theta_{t-1}^*$ , which lies in the confidence set with the high probability [1]. Then, regarding the activation function  $\sigma(\cdot)$ , we have the following assumption on its continuity and smoothness.

Assumption 5.2 ( $\zeta$ -Lipschitz continuity and Smoothness [4, 13]). For non-linear activation function  $\sigma(\cdot)$ , there exists a positive constant  $\zeta > 0$ , such that  $\forall x, x' \in \mathbb{R}$ , we have

$$|\sigma(x) - \sigma(x')| \le \zeta \cdot ||x - x'||, \quad |\sigma'(x) - \sigma'(x')| \le \zeta \cdot ||x - x'||$$

with  $\sigma'(\cdot)$  being the derivative of activation function  $\sigma(\cdot)$ .

Note that **Assumption** 5.2 is mild and applicable on many activation functions, such as Sigmoid. Then, we proceed to bound the regret for a single time step t.

# 5.1 Upper Confidence Bound

Recall that at time step t, given an embedded arm matrix  $\widetilde{X}$ , the output of our proposed framework is  $\widehat{r} = f(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1})$  with  $\mathcal{G}_t$ ,  $\Theta_{t-1}$  as the estimated arm group graph and trained parameters respectively. The true function  $h(\mathcal{G}^*, \widetilde{X})$  is given in Lemma 5.1. Supposing there exists the true arm group graph  $\mathcal{G}^*$ , the confidence bound for a single round t will be

$$CB_{t}(\widetilde{X}) = |f(\mathcal{G}_{t}, \widetilde{X}; \Theta_{t-1}) - h(\mathcal{G}^{*}, \widetilde{X})|$$

$$\leq \underbrace{|f(\mathcal{G}_{t}, \widetilde{X}; \Theta_{t-1}) - h(\mathcal{G}_{t}, \widetilde{X})|}_{R_{1}} + \underbrace{|h(\mathcal{G}_{t}, \widetilde{X}) - h(\mathcal{G}^{*}, \widetilde{X})|}_{R_{2}}$$
(9)

where  $R_1$  denotes the error induced by network parameter estimations, and  $R_2$  refers to the error from arm group graph estimations. We will then proceed to bound them separately.

5.1.1 Bounding network parameter error  $R_1$ . For simplicity, the  $\mathcal{G}_t$  notation is omitted for this subsection. To bridge the network parameters after GD with those at random initialization, we define the gradient-based regression estimator  $\widehat{\Theta}_t = Z_t^{-1}b_t$  where  $Z_t = \lambda I + \frac{1}{m} \sum_{\tau=1}^t g(\widetilde{X}_\tau; \Theta_\tau) \cdot g(\widetilde{X}_\tau; \Theta_\tau)^\intercal, b_t = \sum_{\tau=1}^t r_\tau \cdot g(\widetilde{X}_\tau; \Theta_\tau)/\sqrt{m}$ . Then, we derive the bound for  $R_1$  with the following lemma.

LEMMA 5.3. Assume there are constants  $\beta_F > 0$ ,  $1 < \beta_1, \beta_2, \beta_3, \beta_4 < 2$ ,  $\beta_L = \max\{\beta_1, \beta_2, \beta_3, \beta_4\}$ , and

$$\beta_h = \max\{\zeta\beta_1,\ \zeta\beta_2 + \zeta^2\beta_1\beta_2,\ \zeta\beta_L + 1,\ (\zeta\beta_4)^{L-2}(\zeta\beta_2 + \zeta^2\beta_1\beta_2)\}.$$

With a constant  $\delta \in (0,1)$ , and L as the layer number for the FC network, let the network width  $m \geq Poly(t,L,\frac{1}{\beta_F},\frac{1}{\lambda},(\zeta\beta_L)^L,\log(\frac{1}{\delta}))$ , and learning rate  $\eta \leq O((t \cdot L\beta_b^2(2\zeta\beta_L)^{2L})^{-1})$ . Denoting the terms

$$\begin{split} \Upsilon &= \frac{2\sqrt{2}t}{\beta_F} (\beta_h + \Lambda)(\beta_L + 1)^L \zeta^L, \quad \Lambda &= \frac{\zeta \Upsilon \beta_h}{m} \cdot \frac{(2\zeta \beta_L)^L - 1}{2\zeta \beta_L - 1} \\ \widetilde{I}_1 &= \sqrt{t \cdot (L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m)} + \Lambda \cdot \sqrt{t \cdot (9L + m^{-1})}, \\ \widetilde{I}_2 &= \lambda \sqrt{L + 1} \cdot \Upsilon / \sqrt{m}, \end{split}$$

at time step t, given the received contexts and rewards, with probability at least  $1 - \delta$  and the embedded context  $\widetilde{X}$ , we have

$$|h(\widetilde{X}) - f(\widetilde{X}; \Theta_{t-1})| \le B_1 ||g(\widetilde{X}; \Theta_{t-1}) / \sqrt{m}||_{Z_{t-1}^{-1}} + B_2 + B_3$$

with the terms

$$\begin{split} B_1 &= \sqrt{\log(\frac{\det(Z_{t-1})}{\det(\lambda I)}) - 2\log(\delta)} + \lambda^{\frac{1}{2}}S, \\ B_2 &= \Big(\frac{\widetilde{I}_1 \cdot \sqrt{t}B_3 + m \cdot \widetilde{I}_2}{m\lambda} + \sqrt{\frac{t}{m\lambda}}\Big) \\ &\quad \cdot \Big(\Lambda \cdot \sqrt{9L + m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}\Big) \\ B_3 &= m^{-0.5} \Big(\beta_3(\Lambda + \beta_h) + L \cdot \Upsilon \cdot (\Lambda + \beta_h)(\Lambda/\beta_h + 1)\Big). \end{split}$$

**Proof.** Given the embedded context  $\widetilde{X}$ , and following the statement in Lemma 5.1, we have

$$\begin{split} |h(\widetilde{X}) - f(\widetilde{X}; \Theta_{t-1})| \\ & \leq |\langle g(\widetilde{X}; \Theta_{t-1}) / \sqrt{m}, \sqrt{m}(\Theta_{t-1}^* - \Theta_0) \rangle - \langle g(\widetilde{X}; \Theta_{t-1}) / \sqrt{m}, \widehat{\Theta}_{t-1} \rangle| \\ & + |\langle g(\widetilde{X}; \Theta_{t-1}) / \sqrt{m}, \widehat{\Theta}_{t-1} \rangle - f(\widetilde{X}; \Theta_{t-1})| = R_3 + R_4. \end{split}$$

With Theorem 2 from [1], we have  $R_3 \leq B_1 \| g(\widetilde{X}; \Theta_{t-1}) / \sqrt{m} \|_{Z_{t-1}^{-1}}$ . Then, for  $R_4$ , we have  $|f(\widetilde{X}; \Theta_{t-1}) - \langle g(\widetilde{X}; \Theta_{t-1} / \sqrt{m}), \widehat{\Theta}_{t-1} \rangle|$ 

$$\leq R_5 + R_6 = |f(\widetilde{X}; \Theta_{t-1}) - \langle g(\widetilde{X}; \Theta_{t-1}), \Theta_{t-1} - \Theta_0 \rangle|$$
  
+  $|\langle g(\widetilde{X}; \Theta_{t-1}), \Theta_{t-1} - \Theta_0 - \widehat{\Theta}_{t-1} / \sqrt{m} \rangle|$ 

where  $R_5$  can be bounded by  $B_3$  with Lemma A.6. Then, with conclusions from Lemma B.3 and Lemma A.5, we have

$$\begin{aligned} R_6 &\leq \|\Theta_{t-1} - \Theta_0 - \widehat{\Theta}_{t-1} / \sqrt{m} \|_2 \cdot \|g(\widetilde{X}; \Theta_{t-1})\|_2 \\ &\leq B_2 = \left( (\widetilde{I}_1 \cdot \sqrt{t} B_3 + m \cdot \widetilde{I}_2) / (m\lambda) + \sqrt{t/(m\lambda)} \right) \\ &\cdot \left( \Lambda \cdot \sqrt{9L + m^{-1}} + m^{-1} \beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m} \right), \end{aligned}$$

which completes the proof. ■

5.1.2 Bounding graph estimation error  $R_2$ . Regarding the regret term  $R_2$  and for the aggregation module, we have

$$H_{gnn} = \sqrt{\frac{1}{m}} \cdot \sigma(A_t^k \cdot \widetilde{X} \Theta_{gnn}) \in \mathbb{R}^{N_c \times m}$$

as the output where  $\Theta_{gnn}$  refers to the trainable weight matrix. Then, we use the following lemma to bound  $R_2$ .

LEMMA 5.4. At this time step t+1, given any two arm groups  $c_i, c_j \in C$  and their sampled arm contexts  $X_{c_i}^t = \{x_{c_i,\tau}\}_{\tau=1}^t, X_{c_j}^t = \{x_{c_j,\tau}\}_{\tau=1}^t$ , with the notation from Lemma 5.3 and the probability at least  $1-\delta$ , we have

$$\|A^* - A_t\|_{\max} \leq \frac{1}{N_c} \cdot \sqrt{\frac{1}{2t} \log(\frac{N_c^2 - N_c}{\delta})}$$

where  $\|\cdot\|_{max}$  refers to the greatest entry of a matrix. Then, we will have  $R_2 \leq B_4 \sqrt{1/t}$  with

$$B_4 = \frac{S\sqrt{L}k}{\sqrt{m}}(\beta_h + \Lambda)(\zeta\beta_L + \frac{\Upsilon\zeta}{m})^{O(L)}\sqrt{\frac{1}{2}\log(\frac{N_c^2 - N_c}{\delta})},$$

and  $N_c = |C|$  is the number of arm groups.

**Proof.** Recall that for  $c_i, c_j \in C$ , the element of matrix  $[A^*]_{ij} = \frac{1}{N_c} \mathbb{E}_{x_i \sim \mathcal{D}_{c_i}, x_j \sim \mathcal{D}_{c_j}} (\phi_{k_G}(x_i)^\intercal \phi_{k_G}(x_j)), \forall i, j \in [N_c]$ , and  $[A_t]_{ij} = \frac{1}{t \cdot N_c} \sum_{\tau=1}^t \phi_{k_G}(x_{c_i,\tau})^\intercal \phi_{k_G}(x_{c_j,\tau})$ . Here, suppose a distribution  $\mathcal{D}_{ij}$  where  $\mathbb{E}[\mathcal{D}_{ij}] = \frac{1}{N_c} \mathbb{E}_{x_i \sim \mathcal{D}_{c_i}, x_j \sim \mathcal{D}_{c_j}} (\phi_{k_G}(x_i)^\intercal \phi_{k_G}(x_j))$ . Given  $N_c$  arm groups, we have  $N_c(N_c - 1)/2$  different group pairs. For group pair  $c_i, c_j \in C$ , each  $\phi_{k_G}(x_{c_i,\tau})^\intercal \phi_{k_G}(x_{c_j,\tau}), \tau \in [t]$  is a sample drawn from  $\mathcal{D}_{ij}$ , and the element distance  $|[A_t]_{ij} - [A^*]_{ij}|$  can be regarded as the difference between the mean value of samples and the expectation. Applying the Hoeffding's inequality and the union bound would complete the proof. As  $\|\cdot\|_2 \leq n\|\cdot\|_{\max}$  for an  $n \times n$  square matrix, we have the bound for matrix differences.

Then, consider the power of adjacency matrix  $A^k$  (for graph  $\mathcal{G}$ ) as input and fix  $\widetilde{X}$ . Analogous to the idea that the activation function with the Lipschitz continuity and smoothness property will lead to Lipschitz neural networks [2], applying Assumption 5.2 and with Lemma A.2, Lemma A.3, we simply have the gradient

$$\begin{split} g(\mathcal{G},\widetilde{X};\Theta_{t-1}) \text{ being Lipschitz continuous w.r.t. the input graph as} \\ R_2 &\leq \|g(\mathcal{G}^*,\widetilde{X};\Theta_{t-1}) - g(\mathcal{G}_t,\widetilde{X};\Theta_{t-1})\|_2 \cdot \|\Theta_{t-1}^* - \Theta_0\|_2 \\ &\leq \frac{S\sqrt{L}}{\sqrt{m}}(\beta_h + \Lambda)(\zeta\beta_L + \frac{\Upsilon\zeta}{\sqrt{m}})^{O(L)} \cdot |\|(A_t)^k\|_2 - \|(A^*)^k\|_2| \\ &\leq \frac{S\sqrt{L}k}{\sqrt{m}}(\beta_h + \Lambda)(\zeta\beta_L + \frac{\Upsilon\zeta}{\sqrt{m}})^{O(L)} \cdot \|A_t - A^*\|_2 \end{split}$$

where (i) is because  $A_t$ ,  $A^*$  are symmetric and bounded polynomial functions are Lipschitz continuous. Combining the two parts will lead to the conclusion.

5.1.3 Combining  $R_2$  with  $R_1$ . At time step t, with the notation and conclusions from Lemma 5.3 and Lemma 5.4, re-scaling the constant  $\delta$ , we have the confidence bound given embedded arm  $\widetilde{X}$  as

$$CB_t(\widetilde{X}) \le B_1 \|g(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1}) / \sqrt{m}\|_{Z_{t-1}^{-1}} + B_2 + B_3 + B_4 \sqrt{\frac{1}{t}}.$$
 (10)

# 5.2 Regret Bound

With the UCB shown in Eq. 10, we provide the following regret upper bound R(T), for a total of T time steps.

THEOREM 5.5. Given the received contexts and rewards, with the notation from Lemma 5.3, Lemma 5.4, and probability at least  $1 - \delta$ , if m,  $\eta$  satisfy conditions in Lemma 5.3, we will have the regret

$$R(T) \le 2 \cdot \left(2B_4\sqrt{T} + 2 - B_4\right) + 2\sqrt{2\widetilde{d}T\log(1 + T/\lambda)} + 2T$$
$$\cdot \left(\sqrt{\lambda}S + \sqrt{1 - 2\log(\delta/2) + (\widetilde{d}\log(1 + T/\lambda))}\right)$$

where the effective dimension  $\widetilde{d} = \frac{\log \det(I + G(0)/\lambda)}{\log(1 + T/\lambda)}$  with  $G(0) = G_0 G_0^{\mathsf{T}}$  and  $G_0 = \left(g(\widetilde{X}_1; \Theta_0)^{\mathsf{T}}, \dots, g(\widetilde{X}_t; \Theta_0)^{\mathsf{T}}\right)$ .

**Proof.** By definition, we have the regret  $R_t$  for time step t as

$$\begin{split} R_t &= h(\mathcal{G}^*, \widetilde{X}_t^*) - h(\mathcal{G}^*, \widetilde{X}_t) \\ &\leq \mathsf{CB}_t(\widetilde{X}_t^*) + f(\mathcal{G}_t, \widetilde{X}_t^*; \Theta_{t-1}) - h(\mathcal{G}^*, \widetilde{X}_t) \\ &\leq \mathsf{CB}_t(\widetilde{X}_t) + f(\mathcal{G}_t, \widetilde{X}_t; \Theta_{t-1}) - h(\mathcal{G}^*, \widetilde{X}_t) \leq 2 \cdot \mathsf{CB}_t(\widetilde{X}_t) \end{split}$$

where the second inequality is due to our arm pulling mechanism. Then, based on Lemma 5.4, Lemma 5.3, and Eq. 10, we have R(T) =

$$\begin{split} \sum_{t=1}^T R_t &\leq 2 \sum_{t=1}^T \left( B_1 \| g(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1}) / \sqrt{m} \|_{Z_{t-1}^{-1}} + B_2 + B_3 + B_4 \sqrt{\frac{1}{t}} \right) \\ &\leq 2 \cdot (2B_4 \sqrt{T} + 2 - B_4) + 2 \sum_{t=1}^T (B_1 \| g(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1}) / \sqrt{m} \|_{Z_{t-1}^{-1}}) \end{split}$$

with the choice of m for bounding the summation of  $B_2$ ,  $B_3$ , and the bound of  $\sum_{i=1}^{T} [t^{-i/2}]$  in [10]. Then, with Lemma 11 from [1],

$$\begin{split} &\sum_{t=1}^{I} (B_1 \| g(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1}) / \sqrt{m} \|_{Z_{t-1}^{-1}}) \\ &\leq B_1 \sqrt{T \sum_{t=1}^{T} \| g(\mathcal{G}_t, \widetilde{X}; \Theta_{t-1}) / \sqrt{m} \|_{Z_{t-1}^{-1}}^2} \leq \sqrt{T} B_1 \sqrt{2 \log(\frac{\det(Z_T)}{\det(\lambda I)})} \\ &\leq \sqrt{2\widetilde{d}T \log(1 + T/\lambda) + 2T} \left( \sqrt{\lambda}S + \sqrt{1 - 2 \log(\delta/2) + (\widetilde{d}\log(1 + T/\lambda))} \right) \end{split}$$

where (i) is based on Lemma 6.3 in [4] and Lemma 5.4 in [46].  $\blacksquare$ 

Here, the effective dimension  $\widetilde{d}$  measures the vanishing speed of G(0)'s eigenvalues, and it is analogous to that of existing works on neural contextual bandits algorithms [4, 6, 46]. As  $\widetilde{d}$  is smaller than the dimension of the gradient matrix G(0), it is applied to prevent the dimension explosion. Our result matches the state-of-the-art regret complexity [4, 45, 46] under the worst-case scenario.

# 5.3 Model Convergence after GD

For model convergence, we first give an assumption of the gradient matrix after j iterations of GD. First, we define  $G^{(j)}(\Theta_{L-1}) = (g(\widetilde{X}_1;\Theta_{L-1}^{(j)}),\ldots,g(\widetilde{X}_T;\Theta_{L-1}^{(j)}))^{\mathsf{T}}(g(\widetilde{X}_1;\Theta_{L-1}^{(j)}),\ldots,g(\widetilde{X}_T;\Theta_{L-1}^{(j)}))$  where  $g(\widetilde{X};\Theta_{L-1})$  is the gradient vector w.r.t.  $\Theta_{L-1}$ .

Assumption 5.6. With width  $m \ge Poly(T, L, \frac{1}{\beta_F}, \frac{1}{\lambda}, (\zeta \beta_L)^L, \log(\frac{1}{\delta}))$  and for  $j \in [J]$ , we have the minimal eigenvalue of  $G^{(j)}$  as

$$\lambda_{\min}(G^{(j)}(\Theta_{L-1})) \ge \lambda_0/2$$

where  $\lambda_0$  is the minimal eigenvalue of the neural tangent kernel (NTK) [24] matrix induced by AGG-UCB.

Note that Assumption 5.6 is mild and has been proved for various neural architectures in [13]. The NTK for AGG-UCB can be derived following a comparable approach as in [14, 24]. Then, we apply the following lemma and theorem to prove the convergence of AGG-UCB. The proof of Lemma 5.7 is given in the appendix.

Lemma 5.7. After T time steps, assume the network are trained with the J-iterations GD on the past contexts and rewards. Then, with  $\beta_F > 0$  and  $\beta_F \cdot \eta < 1$ , for any  $j \in [J]$ :

$$\|F_T^{(j)} - F_T^{(j+1)}\|_2^2 \le \frac{1}{4}\eta\beta_F \cdot \|F_T^{(j)} - Y_T\|_2^2$$

with network width m defined in Lemma 5.3.

The Lemma 5.7 shows that we are able to bound the difference in network outputs after one step of GD. Then, we proceed to prove the convergence with the theorem below.

Theorem 5.8. After T time steps, assume the model with width m defined in Lemma 5.3 is trained with the J-iterations GD on the contexts  $\{\widetilde{X}_{\tau}\}_{\tau=1}^{T}$  and rewards  $\{r_{\tau}\}_{\tau=1}^{T}$ . With probability at least  $1-\delta$ , a constant  $\beta_{F}$  such that  $\beta_{F} \cdot \eta < 1$ , set the network width  $m \geq Poly(T, L, \frac{1}{\beta_{F}}, \frac{1}{\lambda}, (\zeta\beta_{L})^{L}, \log(\frac{1}{\delta}))$  and the learning rate  $\eta \leq O(T^{-1}L^{-1}\beta_{h}^{-2}(2\zeta\beta_{L})^{-2L})$ . Then, for any  $j \in [J]$ , we have

$$\|F_T^{(j)} - Y_T\|_2^2 \leq (1 - \beta_F \cdot \eta)^j \cdot \|F_T^{(0)} - Y_T\|_2^2$$

where the vector  $\mathbf{F}^{(j)} = [f(\mathcal{G}_T, \widetilde{\mathbf{X}}_\tau; \mathbf{\Theta}^{(j)})]_{\tau=1}^T$ , and  $\mathbf{Y}_T = [r_\tau]_{\tau=1}^T$ .

**Proof.** Following an approach analogous to [13], we apply and induction based method for the proof. The hypothesis is that  $\|F_T^{(j)} - Y_T\|_2^2 \le (1 - \beta_F \cdot \eta)^j \cdot \|F_T^{(0)} - Y_T\|_2^2, j \in [J]$ . With a similar procedure in Condition A.1 of [13], we have

$$\begin{split} \|F_T^{(j+1)} - Y_T\|_2^2 &\leq \|F_T^{(j)} - Y_T\|_2^2 - 2\eta \|F_T^{(j)} - Y_T\|_{G^{(j)}}^2 \\ &- 2(Y_T - F_T^{(j)})^\intercal V^{(j)} + \|F_T^{(j+1)} - F_T^{(j)}\|_2^2 \\ \text{with } V^{(j)} &= (V^{(j)}(\widetilde{X}_1), \dots, V^{(j)}(\widetilde{X}_T))^\intercal . \text{ For } \Theta' \in \{\Theta_{gnn}, \dots, \Theta_{L-1}\}, \\ |V^{(j)}(\widetilde{X})| &= \eta \max_{0 \leq s \leq \eta} \left[ \sum_{C \leq s} \|\nabla \mathcal{L}(\Theta'^{(j)})\|_F \|\nabla f(\Theta'^{(j)}) - \nabla f(\Theta'^{(j)}, s)\|_F \right] \end{split}$$

where  $\nabla f(\Theta'^{(j)}, s) = \nabla f(\Theta'^{(j)} - s \cdot \nabla \mathcal{L}(\Theta'^{(j)}))$ . The notation  $\mathcal{G}, \widetilde{X}$  is omitted for simplicity. Then, based on the conclusions from Lemma C.1, Lemma 5.7 and Assumption 5.6, we can have

$$\begin{split} \|F_T^{(j+1)} - Y_T\|_2^2 &\leq (1 - \eta \lambda_0) \|F_T^{(j)} - Y_T\|_2^2 - 2(Y_T - F_T^{(j)})^{\mathsf{T}} V^{(j)} \\ &+ \|F_T^{(j+1)} - F_T^{(j)}\|_2^2 \leq (1 - \frac{\eta \lambda_0}{2}) \|F_T^{(j)} - Y_T\|_2^2 \end{split}$$

by setting  $\beta_F = \lambda_0/2$ .

This theorem shows that with sufficiently large m and proper  $\eta$ , the GD will converge to the global minimum at a linear rate, which is essential for proving the regret bound.

#### **6 EXPERIMENTS**

In this section, we demonstrate the effectiveness of our proposed framework by comparing its performances with state-of-the-art baselines through experiments on four real data sets. As linear algorithms have been outperformed in previous works [12, 45, 46], we will not include these linear methods in the experiments below. Our six baseline algorithms are:

- KMTL-UCB [12] estimates the "task similarities" with received contextual information. The estimations are based on a variant of kernel ridge regression.
- Kernel-Ind is Kernel-UCB [34] under the "disjoint setting"
   [28] where it learns individual estimators for each arm group.
- Kernel-Pool represents Kernel-UCB under the "pooling setting" where it applies a single estimator for all arm groups.
- Neural-TS stands for Neural Thompson Sampling [45] with group-aware embedding, which enables it to leverage the group information. It applies a neural network for exploitation and Thompson sampling strategy for exploration.
- Neural-Pool is for Neural-UCB [46] with a single neural network to evaluate the reward, and calculate the upper confidence bounds with the network gradients.
- Neural-Ind represents Neural-UCB with group-aware embedding for utilizing the group information.

Note that COFIBA [30] is naturally Kernel-Ind (with linear kernel) given the arm group information and one single user to serve, so we do not include it in our benchmarks. To find the best exploration parameter, we perform grid searches over the range  $\{10^{-1}, 10^{-2}, 10^{-3}\}$  for all algorithms. Similarly, the learning rate for neural algorithms are chosen from  $\{10^{-2}, 10^{-3}, 10^{-4}\}$ . For Neural-UCB, Neural-TS and our reward estimation module, we apply a two-layer FC network with m=500. RBF kernels are applied for KMTL-UCB and Kernel-UCB as well as our graph estimation module. Kernel-Pool and Neural-Pool will not fit into the multi-class classification setting, as we only receive one arm (context) at each time step without the arm group information.

#### 6.1 Real Data Sets

Here, we compare our proposed model with baseline algorithms on four real data sets with different specifications.

MovieLens and Yelp data sets. The first real data set is the "MovieLens 20M rating data set" (grouplens.org/datasets/movielens/20m/). To obtain the user features, we first choose 100 movies and 4000 users with **most reviews** to form the user-movie matrix where the entries are user ratings, and the user features  $v_u \in \mathbb{R}^d$  are obtained through singular value decomposition (SVD) where the

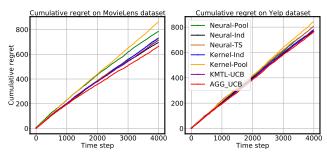


Figure 1: Cumulative regrets for recommendation data sets.

dimension d=20. Then, since the genome-scores of user-specified tags are provided for each movie, we select 20 tags with the highest variance to construct the movie features  $v_i \in \mathbb{R}^d$  with their scores on these tags. Then, these movies are allocated into 19 groups based on their genres (|C|=19). Receiving a user  $u_t$  at each time step t, we follow the idea of Generalized Matrix Factorization (GMF) [22, 47, 48] to encode user information into the contexts as  $\widetilde{\boldsymbol{x}}_{c,t}^{(i)} = [\boldsymbol{v}_{u_t} \odot \boldsymbol{v}_i] \in \mathbb{R}^d$ ,  $c \in C_t$ ,  $i \in [n_{c,t}]$ , and let  $|X_t|=20$ . Finally, we concatenate a constant 0.01 to each  $\widetilde{\boldsymbol{x}}_{c,t}^{(i)}$  to obtain  $\boldsymbol{x}_{c,t}^{(i)} \in \mathbb{R}^{d_x}$ , which makes  $d_x=21$ , before normalizing  $\boldsymbol{x}_{c,t}^{(i)}$ . Rewards  $r_{c,t}^{(i)}$  are user ratings normalized into range [0,1].

Then, for the Yelp data set (https://www.yelp.com/dataset), we choose 4000 users with **most reviews** and restaurants from 20 different categories as arms (|C|=20). Both user features and arm features are obtained through SVD with the dimension d=20. Analogous to the MovieLens data set, we follow the GMF based approach and the fore-mentioned constant concatenation to get the arm context  $\mathbf{x}_{c,t}^{(i)}$  ( $d_x=21$ ,  $|X_t|=20$ ) to encode the user information, and the rewards are the normalized user ratings.

MNIST data set with augmented classes (MNIST-Aug). MNIST is a well-known classification data set with 10 original classes where each sample is labeled as a digit from 0 to 9. Here, we further divide the samples from each class into 5 sub-divisions through K-means clustering, which gives us a total of 50 augmented sub-classes (i.e., arm groups) for the whole data set. Given a sample  $x_t$ , the reward would be  $r_t=1$  if the learner accurately predicts its sub-class; or the learner will receive the partial reward  $r_t=0.5$  when it chooses the wrong sub-class, but this sub-class and the correct one belong to the same digit (original class). Otherwise, the reward  $r_t=0$ .

XRMB data set. XRMB data set [37] is a multi-view classification data set with 40 different labels. Here, we only apply samples from the first 38 classes as there are insufficient samples for the last two classes. The arm contexts  $x_t$  are the first-view features of the samples. Then, learner will receive a reward of  $r_t = 1$  when they predict the right label, and  $r_t = 0$  otherwise.

# 6.2 Experimental Results

Figure 1 shows the cumulative regret results on the two real recommendation data sets where our proposed AGG-UCB outperforms all strong baselines. In particular, we can find that algorithms with group-aware arm embedding tend to perform better than those without the arm group information (Kernel-Pool, Neural-Pool). This confirms the necessity of exploiting arm group information. Nevertheless, these baselines fed with group-aware are outperformed by

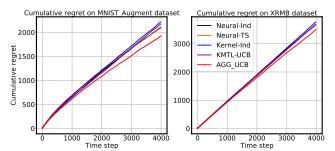


Figure 2: Cumulative regrets for classification data sets.

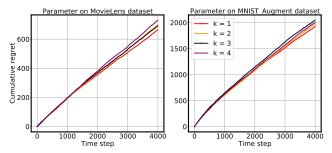


Figure 3: Cumulative regrets on MovieLens and MNIST-Aug data sets with different neighborhood parameter k.

AGG-UCB, which implies the advantages of of our new graph-based model. Meantime, it can be observed that neural algorithms (AGG-UCB, Neural-Ind, Neural-TS) generally perform better compared with other baselines due to the representation power of neural networks. Note that since the user features and arm features of the Yelp data set are directly extracted with SVD, the reward estimation on the Yelp data set is comparably easy compared with others data sets. Therefore, the performances of benchmarks do not differ dramatically with AGG-UCB. In opposite, MovieLens data set with true arm features tends to be a more challenging task where a more complex mapping from arms to their rewards can be involved. This can be reason for AGG-UCB's superiority over the competitors.

Then, Figure 2 shows the cumulative regret results on the two classification data sets where our AGG-UCB achieves the best performance compared with other baselines. In particular, since subclasses from each digit are highly correlated in the MNIST-Aug data set, our proposed AGG-UCB tends to perform significantly better due to its ability of leveraging arm group correlations compared with other neural methods. Thus, these two aspects verify our claim that associating the neural models with arm group relationship modeling can lead to better performance.

### 6.3 Parameter Study

In this section, we conduct our parameter study for the neighborhood parameter k on the MovieLens data set and MNIST-Aug data set with augmented labels, and the results are presented in Figure 3. For the MovieLens data set, we can observe that setting k=1 would give the best result. Although increasing k can enable the aggregation module to propagate the hidden representations for multiple hops, it can potentially fail to focus on local arm group neighbors with high correlations, which is comparable to the aforementioned "over-smoothing" problem. In addition, since the arm group graph of MovieLens data set only has 19 nodes, k=1 would be enough.

Meantime, setting k=1 also achieves the best performance on the MNIST data set. The reason can be that the 1-hop neighborhood of each sub-class can already include all the other sub-classes from the same digit with heavy edge weights within the neighborhood for arm group collaboration. Therefore, unless setting k to considerably large values, the AGG-UCB can maintain robust performances, which reduces the workload for hyperparameter tuning.

#### 7 CONCLUSION

In this paper, motivated by real applications where the arm group information is available, we propose a new graph-based model to characterize the relationship among arm groups. Base on this model, we propose a novel UCB-based algorithm named AGG-UCB, which uses GNN to exploit the arm group relationship and share the information across similar arm groups. Compared with existing methods, AGG-UCB provides a new way of collaborating multiple neural contextual bandit estimators for obtaining the rewards. In addition to the theoretical analysis of AGG-UCB, we empirically demonstrate its superiority on real data sets in comparison with state-of-the-art baselines.

#### **ACKNOWLEDGMENTS**

This work is supported by National Science Foundation under Award No. IIS-1947203, IIS-2117902, IIS-2137468, and IIS-2002540. The views and conclusions are those of the authors and should not be interpreted as representing the official policies of the funding agencies or the government.

#### REFERENCES

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. 2011. Improved algorithms for linear stochastic bandits. NeurIPS 24 (2011). 2312–2320.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. 2019. A convergence theory for deep learning via over-parameterization. In ICML. PMLR, 242–252.
  [2] Peter Augus Nicole Cong Bisseli and Paul Fischer 2002. First time analysis of
- [3] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47, 2-3 (2002), 235–256.
- [4] Yikun Ban and Jingrui He. 2021. Convolutional neural bandit: Provable algorithm for visual-aware advertising. arXiv preprint arXiv:2107.07438 (2021).
   [5] Yikun Ban and Jingrui He. 2021. Local clustering in contextual multi-armed
- 5) Yikun Ban and Jingrui He. 2021. Local clustering in contextual multi-armec bandits. In *Proceedings of the Web Conference* 2021. 2335–2346.
- [6] Yikun Ban, Jingrui He, and Curtiss B Cook. 2021. Multi-facet Contextual Bandits: A Neural Network Perspective. arXiv preprint arXiv:2106.03039 (2021).
- [7] Yikun Ban, Yunzhe Qi, Tianxin Wei, and Jingrui He. 2022. Neural Collaborative Filtering Bandits via Meta Learning. ArXiv abs/2201.13395 (2022).
- [8] Yikun Ban, Yuchen Yan, Arindam Banerjee, and Jingrui He. 2022. EE-Net: Exploitation-Exploration Neural Networks in Contextual Bandits. In ICLR.
- [9] Gilles Blanchard, Gyemin Lee, and Clayton Scott. 2011. Generalizing from several related classification tasks to a new unlabeled sample. In NeurIPS. 2178–2186.
- [10] Edward Chlebus. 2009. An approximate formula for a partial sum of the divergent p-series. Applied Mathematics Letters 22, 5 (2009), 732–737.
- [11] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. 2011. Contextual bandits with linear payoff functions. In AISTATS. 208–214.
- [12] Aniket Anand Deshmukh, Urun Dogan, and Clay Scott. 2017. Multi-task learning for contextual bandits. In NeurIPS. 4848–4856.
- [13] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. 2019. Gradient descent finds global minima of deep neural networks. In ICML. PMLR, 1675–1685.
- [14] Simon S Du, Kangcheng Hou, Barnabás Póczos, Ruslan Salakhutdinov, Ruosong Wang, and Keyulu Xu. 2019. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. arXiv preprint arXiv:1905.13192 (2019).
- [15] Audrey Durand, Charis Achilleos, Demetris Iacovides, Katerina Strati, Georgios D Mitsis, and Joelle Pineau. 2018. Contextual bandits for adapting treatment in a mouse model of de novo carcinogenesis. In Machine learning for healthcare conference. PMLR, 67–82.
- [16] Dongqi Fu and Jingrui He. 2021. DPPIN: A biological repository of dynamic protein-protein interaction network data. arXiv preprint arXiv:2107.02168 (2021).
- [17] Dongqi Fu and Jingrui He. 2021. SDG: A Simplified and Dynamic Graph Neural Network. In SIGIR '21. 2273–2277.

- [18] Claudio Gentile, Shuai Li, Purushottam Kar, Alexandros Karatzoglou, Giovanni Zappella, and Evans Etrue. 2017. On context-dependent clustering of bandits. In ICML. 1253–1262.
- [19] Claudio Gentile, Shuai Li, and Giovanni Zappella. 2014. Online clustering of bandits. In ICML. 757–765.
- [20] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. arXiv preprint arXiv:1706.02216 (2017).
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In CVPR. 770–778.
- [22] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In WWW. 173–182.
- [23] NT Hoang, Takanori Maehara, and Tsuyoshi Murata. 2021. Revisiting Graph Neural Networks: Graph Filtering Perspective. In 2020 25th International Conference on Pattern Recognition (ICPR). IEEE, 8376–8383.
- [24] Arthur Jacot, Franck Gabriel, and Clément Hongler. 2018. Neural tangent kernel: Convergence and generalization in neural networks. arXiv preprint arXiv:1806.07572 (2018).
- [25] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907 (2016).
- [26] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Predict then propagate: Graph neural networks meet personalized pagerank. arXiv preprint arXiv:1810.05997 (2018).
- [27] Andreas Krause and Cheng Soon Ong. 2011. Contextual Gaussian Process Bandit Optimization.. In NeurIPS. 2447–2455.
- [28] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In WWW. 661–670.
- [29] Shuai Li, Wei Chen, Shuai Li, and Kwong-Sak Leung. 2019. Improved Algorithm on Online Clustering of Bandits. In IJCAI. 2923–2929.
- [30] Shuai Li, Alexandros Karatzoglou, and Claudio Gentile. 2016. Collaborative filtering bandits. In SIGIR. 539–548.
- [31] Sandra Sajeev, Jade Huang, Nikos Karampatziakis, Matthew Hall, Sebastian Kochman, and Weizhu Chen. 2021. Contextual Bandit Applications in a Customer Support Bot. In KDD '21. 3522–3530.
- [32] Xin Shao, Ning Lv, Jie Liao, Jinbo Long, Rui Xue, Ni Ai, Donghang Xu, and Xiaohui Fan. 2019. Copy number variation is highly correlated with differential gene expression: a pan-cancer study. BMC medical genetics 20, 1 (2019), 1–14.
- [33] Sohini Upadhyay, Mikhail Yurochkin, Mayank Agarwal, Yasaman Khazaeni, et al. 2020. Online Semi-Supervised Learning with Bandit Feedback. arXiv preprint arXiv:2010.12574 (2020).
- [34] Michal Valko, Nathaniel Korda, Rémi Munos, Ilias Flaounas, and Nelo Cristianini. 2013. Finite-time analysis of kernelised contextual bandits. arXiv preprint arXiv:1309.6869 (2013).
- [35] Roman Vershynin. 2010. Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027 (2010).
- [36] Sofia S Villar, Jack Bowden, and James Wason. 2015. Multi-armed bandit models for the optimal design of clinical trials: benefits and challenges. Statistical science: a review journal of the Institute of Mathematical Statistics 30, 2 (2015), 199.
- [37] Weiran Wang, Raman Arora, Karen Livescu, and Jeff A Bilmes. 2015. Unsupervised learning of acoustic features via deep canonical correlation analysis. In 2015 IEEE ICASSP IEEE.
- [38] Tianxin Wei, Ziwei Wu, Ruirui Li, Ziniu Hu, Fuli Feng, Xiangnan He, Yizhou Sun, and Wei Wang. 2020. Fast adaptation for cold-start collaborative filtering with meta-learning. In ICDM. IEEE, 661–670.
- [39] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In ICML. PMLR, 6861–6871.
- [40] Qingyun Wu, Huazheng Wang, Quanquan Gu, and Hongning Wang. 2016. Contextual bandits in a collaborative environment. In SIGIR. 529–538.
- [41] Qingyun Wu, Huazheng Wang, Yanen Li, and Hongning Wang. 2019. Dynamic Ensemble of Contextual Bandits to Satisfy Users' Changing Interests. In WWW. 2080–2090.
- [42] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. In ICML. PMLR, 5453–5462.
- [43] Keyulu Xu, Mozhi Zhang, Stefanie Jegelka, and Kenji Kawaguchi. 2021. Optimization of graph neural networks: Implicit acceleration by skip connections and more depth. In ICML. PMLR, 11592–11602.
- [44] Jiaxuan You, Rex Ying, and Jure Leskovec. 2019. Position-aware graph neural networks. In ICML. PMLR, 7134–7143.
- [45] Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. 2020. Neural thompson sampling. arXiv preprint arXiv:2010.00827 (2020).
- [46] Dongruo Zhou, Lihong Li, and Quanquan Gu. 2020. Neural Contextual Bandits with UCB-based Exploration. arXiv:1911.04462 [cs.LG]
- [47] Yao Zhou, Haonan Wang, Jingrui He, and Haixun Wang. 2021. From Intrinsic to Counterfactual: On the Explainability of Contextualized Recommender Systems. ArXiv (2021). arXiv:2110.14844

[48] Yao Zhou, Jianpeng Xu, Jun Wu, Zeinab Taghavi Nasrabadi, Evren Körpeoglu, Kannan Achan, and Jingrui He. 2021. PURE: Positive-Unlabeled Recommendation with Generative Adversarial Network. In KDD '21. 2409–2419.

# A LEMMAS FOR INTERMEDIATE VARIABLES AND WEIGHT MATRICES

Due to page limit, we will give the proof sketch for lemmas at the end of each corresponding appendix section. Recall that each input context  $x_{c,t}^{(i)}, i \in [n_{c,t}]$  is embedded to  $\widetilde{X}_{c,t}^{(i)}$  (represented by  $\widetilde{X}$  for brevity). Supposing  $\widetilde{X}$  belongs to the arm group c, denote  $h_A = [A_t^k \widetilde{X}]_c$  as the corresponding row in matrix  $A_t^k \widetilde{X}$  based on index of group c in C (if group c is the c'-th group in C, then  $h_A$  is the c'-th row in  $A_t^k \widetilde{X}$ ). Similarly, we have  $h_{gnn} = [H_{gnn}]_c$  and  $h_l = [H_l]_c$  respectively. Given received contexts  $\{\widetilde{X}_T\}_{\tau=1}^T$  and rewards  $\{r_T\}_{\tau=1}^T$ , the gradient w.r.t. weight matrix  $\Theta_l$ ,  $\forall l \in \{1, \ldots, L-1\}$  will be

$$\frac{\partial \mathcal{L}(\Theta)}{\partial \Theta_l} = m^{-\frac{L-l+1}{2}} \sum_{\tau=1}^{T} |f(\widetilde{X}_{\tau}; \Theta) - r_{\tau}|^2 \left( \mathbf{h}_{l-1} \Theta_L^{\mathsf{T}} \left( \prod_{q=l+1}^{L-1} \Gamma_q \Theta_q^{\mathsf{T}} \right) \Gamma_l \right)$$

where  $\Gamma_q = diag([\sigma'(h_{q-1}\Theta_q)])$  is the diagonal matrix whose entries are the elements from  $\sigma'(h_{q-1}\Theta_q)$ . The coefficient  $\frac{1}{2}$  of the cost function is omitted for simplicity. Then, for  $\Theta_{qnn}$ , we have

$$\frac{\partial \mathcal{L}(\Theta)}{\partial \Theta_{gnn}} = m^{-\frac{L+1}{2}} \sum_{\tau=1}^{T} |f(\widetilde{X}_{\tau}; \Theta) - r_{\tau}|^{2} \left( h_{A} \Theta_{L}^{\mathsf{T}} \left( \prod_{q=2}^{L-1} \Gamma_{q} \Theta_{q}^{\mathsf{T}} \right) \Gamma_{1} \Theta_{1}^{\mathsf{T}} Q \Gamma_{gnn} \right)$$

where 
$$\Gamma_{gnn} = diag([\sigma'(h_A\Theta_{gnn})]). Q = \left(\frac{I \in \mathbb{R}^{m \times m}}{0 \in \mathbb{R}^{(d'-m) \times m}}\right) \in$$

 $\mathbb{R}^{d'\times m}$ . Given the same  $\mathcal{G}_t$ , we provide lemmas to bound the  $R_1$  term of **Eq.** 9. For brevity, the subscript  $\tau \in [T]$  and notation  $\mathcal{G}_t$  are omitted below by default.

Lemma A.1. Given the randomly initialized parameters  $\Theta^{(0)} = \{\Theta_{gnn}^{(0)},\Theta_1^{(0)},\Theta_2^{(0)},\ldots,\Theta_L^{(0)}\}$ , with the probability at least  $1-O(TL)\cdot e^{-\Omega(m)}$  and constants  $1<\beta_1,\beta_2,\beta_3,\beta_4<2$ , we have

$$\begin{split} &\|\boldsymbol{\Theta}_{gnn}^{(0)}\|_{2} \leq \beta_{1}\sqrt{m}, \ \|\boldsymbol{\Theta}_{1}^{(0)}\|_{2} \leq \beta_{2}\sqrt{m}, \ \|\boldsymbol{\Theta}_{L}^{(0)}\|_{2} \leq \beta_{3}, \\ &\|\boldsymbol{h}_{gnn}^{(0)}\|_{2} \leq \zeta \cdot \beta_{1}, \quad \|\boldsymbol{h}_{1}^{(0)}\|_{2} \leq \zeta \cdot \beta_{2} + \zeta^{2} \cdot \beta_{1}\beta_{2}, \\ &|f(\widetilde{X};\boldsymbol{\Theta}^{(0)})| \leq \zeta \cdot \beta_{3} \cdot (\zeta \cdot \beta_{4})^{L-2}(\zeta \cdot \beta_{2} + \zeta^{2} \cdot \beta_{1}\beta_{2})/\sqrt{m}, \\ &\|\boldsymbol{\Theta}_{l}^{(0)}\|_{2} \leq \beta_{4}\sqrt{m}, \ \|\boldsymbol{h}_{l}^{(0)}\|_{2} \leq (\zeta \cdot \beta_{4})^{l-1}(\zeta \cdot \beta_{2} + \zeta^{2} \cdot \beta_{1}\beta_{2}), \\ \forall l \in \{2, \dots, L-1\}. \end{split}$$

**Proof.** Based on the properties of random Gaussian matrices [4, 13, 35], with the probability of at least  $1-e^{-\frac{(\beta_1-\sqrt{d_{\widetilde{X}}/m}-1)^2\cdot m}{2}}=1-e^{-\Omega(m)}$ , we have

$$\|\mathbf{\Theta}_{ann}^{(0)}\|_2 \le \beta_1 \sqrt{m}$$

where  $\beta_1 \geq \sqrt{d_{\widetilde{X}}/m} + 1$  with  $m > d_{\widetilde{X}}$ . Applying the analogous approach for the other randomly initialized matrices would give similar bounds. Regarding the nature of A, we can easily have  $\|h_A\|_2 \leq 1$ . Then,

$$\|\boldsymbol{h}_{qnn}^{(0)}\|_{2} = m^{-\frac{1}{2}}\|\sigma(\boldsymbol{h}_{A}\cdot\boldsymbol{\Theta}_{qnn})\|_{2} \leq \zeta m^{-\frac{1}{2}}\cdot\|\boldsymbol{h}_{A}\|_{2}\|\boldsymbol{\Theta}_{qnn}\|_{2} \leq \zeta\cdot\beta_{1}$$

due to the assumed  $\zeta$ -Lipschitz continuity. Denoting the concatenated input for reward estimation module as  $x' = [h_{qnn}^{(0)}; \widetilde{X}]_c \in$ 

 $\mathbb{R}^{1\times (d_{\widetilde{x}}+m)}$ , we can easily derive that  $\|x'\|_2 \leq \zeta \cdot \beta_1 + 1$ . Thus,

$$\|\boldsymbol{h}_{1}^{(0)}\|_{2} = m^{-\frac{1}{2}} \|\sigma(\boldsymbol{x}' \cdot \boldsymbol{\Theta}_{1})\|_{2} \leq \zeta m^{-\frac{1}{2}} \cdot \|\boldsymbol{x}'\|_{2} \|\boldsymbol{\Theta}_{1}\|_{2}$$
  
$$\leq \zeta \cdot \beta_{2}(\zeta \cdot \beta_{1} + 1) = \zeta \cdot \beta_{2} + \zeta^{2} \cdot \beta_{1}\beta_{2}.$$

Following the same procedure recursively for other intermediate outputs and applying the union bound would complete the proof.

Lemma A.2. After T time steps, run GD for J-iterations on the network with the received contexts and rewards. Suppose  $\|\mathbf{h}_l^{(j)} - \mathbf{h}_l^{(0)}\|_2 \leq \Lambda^{(j)}, \forall j \in [J]$ . With the probability of at least  $1 - O(TL) \cdot e^{-\Omega(m)}$  and  $\Theta \in \{\Theta_{ann}, \Theta_1, \dots, \Theta_L\}$ , we have

$$\|\mathbf{\Theta}^{(j)} - \mathbf{\Theta}^{(0)}\|_F \leq \Upsilon/\sqrt{m}$$

where 
$$\Upsilon = \frac{2\sqrt{2}t}{\beta_F}(\beta_h + \Lambda^{(j)})(\beta_L + 1)^L \zeta^L$$
.

**Proof.** We prove this Lemma following an induction-based procedure [4, 13]. The hypothesis is  $\|\Theta^{(0)} - \Theta^{(j)}\|_F \le \Upsilon/\sqrt{m}$ ,  $\forall \Theta \in \{\Theta_{gnn}, \Theta_1, \dots, \Theta_L\}$ , and let  $\beta_L = \max\{\beta_1, \beta_2, \beta_3, \beta_4\}$ . According to **Algorithm** 2, we have for the j+1-th iteration and  $l \in [L]$ ,

$$\begin{split} \|\boldsymbol{\Theta}_{l}^{(j+1)} - \boldsymbol{\Theta}_{l}^{(j)}\|_{F} &= m^{-\frac{L-l+1}{2}} \boldsymbol{\eta} \cdot \|\sum_{\tau=1}^{T} |f(\widetilde{X}_{\tau}; \boldsymbol{\Theta}^{(j)}) - r_{\tau}|^{2} \cdot \boldsymbol{h}_{l-1}^{(j)}(\boldsymbol{\Theta}_{L}^{(j)})^{\intercal} \\ & \cdot \big(\prod_{q=l+1}^{L-1} \Gamma_{q}^{(j)} \cdot (\boldsymbol{\Theta}_{q}^{(j)})^{\intercal} \big) \cdot \Gamma_{l}^{(j)}\|_{F} \end{split}$$

$$\leq m^{-\frac{L-l+1}{2}}\eta\sqrt{t}\cdot\|F_t^{(j)}-Y_t\|_2\|\boldsymbol{h}_{l-1}^{(j)}\boldsymbol{\Theta}_L^{(j)}\|_F\prod_{q=l+1}^{L-1}\|\boldsymbol{\Theta}_q^{(j)}\|_2\prod_{q=l}^{L-1}\|\boldsymbol{\Gamma}_q^{(j)}\|_2$$

$$\leq m^{-\frac{L-l+1}{2}}\eta\sqrt{t}\cdot\|F_t^{(j)}-Y_t\|_2\|\boldsymbol{h}_{l-1}^{(j)}\|_2\|\boldsymbol{\Theta}_L^{(j)}\|_2\prod_{q=l+1}^{L-1}\|\boldsymbol{\Theta}_q^{(j)}\|_2\prod_{q=l}^{L-1}\|\Gamma_q^{(j)}\|_2$$

by Cauchy inequality. For  $\|\Theta_q^{(j)}\|_2$ , we have

$$\prod_{q=l+1}^{L-1} \|\boldsymbol{\Theta}_{q}^{(j)}\|_{2} \leq \prod_{q=l+1}^{L-1} \left( \|\boldsymbol{\Theta}_{q}^{(0)}\|_{2} + \|\boldsymbol{\Theta}_{q}^{(j)} - \boldsymbol{\Theta}_{q}^{(0)}\|_{2} \right) \\
\leq (\beta_{L} \sqrt{m} + \Upsilon / \sqrt{m})^{L-l-1};$$

while for  $\|\Gamma_q^{(j)}\|_2$ , we have  $\prod_{q=l}^{L-1} \|\Gamma_q^{(j)}\|_2 \le \zeta^{L-l}$ . Combining all the results above and based on Lemma 5.8, it means that for  $l \in [L]$ ,

$$\begin{split} \|\Theta_{l}^{(j+1)} - \Theta_{l}^{(j)}\|_{F} &\leq m^{-\frac{L-l+1}{2}} \eta \sqrt{t} \cdot (1 - \beta_{F} \cdot \eta)^{j/2} \cdot \|F_{t}^{(0)} - Y_{t}\|_{2} \\ & \cdot \|h_{l-1}^{(j)}\|_{2} \cdot \|\Theta_{L}^{(j)}\|_{2} \cdot (\beta_{L} \sqrt{m} + \Upsilon/\sqrt{m})^{L-l-1} \cdot \zeta^{L-l+1} \\ &\leq m^{-\frac{1}{2}} (1 - \beta_{F} \eta)^{j/2} \eta \sqrt{t} \|F_{t}^{(0)} - Y_{t}\|_{2} ((\beta_{h} + \Lambda^{(j)}) (\beta_{L} + \Upsilon/m)^{L-l} \zeta^{L-l} \end{split}$$

where the last inequality is due to Lemma A.3. Then, since we have  $\|\Theta_I^{(j+1)} - \Theta_I^{(j)}\|_F \le \|\Theta_I^{(j+1)} - \Theta_I^{(j)}\|_F + \|\Theta_I^{(j)} - \Theta_I^{(0)}\|_F$ , it leads to

$$\|\Theta_l^{(j+1)} - \Theta_l^{(0)}\|_F \leq \frac{2\sqrt{t}}{\beta_F \sqrt{m}} \|F_t^{(0)} - Y_t\|_2 (\beta_h + \Lambda^{(j)}) (\beta_L + \Upsilon/m)^{L-l} \zeta^{L-l}.$$

For the last layer  $\Theta_L$ , the conclusion can be verified through a similar procedure. Analogously, for  $\Theta_{qnn}$ , we have

$$\begin{split} \|\boldsymbol{\Theta}_{gnn}^{(j+1)} - \boldsymbol{\Theta}_{gnn}^{(j)}\|_{F} \\ &= m^{-\frac{L+1}{2}} \eta \|\sum_{\tau=1}^{T} |f(\widetilde{\boldsymbol{X}}_{\tau};\boldsymbol{\Theta}) - r_{\tau}|^{2} \cdot \left(\boldsymbol{h}_{A} \cdot \boldsymbol{\Theta}_{L}^{\intercal} \cdot \left(\prod_{q=2}^{L-1} \boldsymbol{\Gamma}_{q} \cdot \boldsymbol{\Theta}_{q}^{\intercal}\right) \cdot \boldsymbol{\Gamma}_{1} \boldsymbol{\Theta}_{1}^{\intercal}\right) \\ & \cdot \boldsymbol{Q} \cdot \boldsymbol{\Gamma}_{gnn}\|_{F} \end{split}$$

$$\leq m^{-\frac{L+1}{2}} \eta \sqrt{t} (1 - \beta_F \cdot \eta)^{j/2} \|F_t^{(0)} - Y_t\|_2 \|h_A\|_2 \|\prod_{q=1}^L \Theta_q\|_2 \zeta^L \|Q\|_2$$

$$\leq \sqrt{t} m^{-\frac{1}{2}} (1-\beta_F \cdot \eta)^{j/2} \eta \cdot \|\boldsymbol{F}_t^{(0)} - \boldsymbol{Y}_t\|_2 \cdot \boldsymbol{\zeta}^L \cdot (\beta_L + \boldsymbol{\Upsilon}/m)^L,$$

which leads to

$$\|\Theta_{gnn}^{(j+1)} - \Theta_{gnn}^{(0)}\|_F \leq \frac{2}{\beta_F} \sqrt{t} m^{-\frac{1}{2}} \cdot \|F_t^{(0)} - Y_t\|_2 \cdot \zeta^L \cdot (\beta_L + \Upsilon/m)^L.$$

Since  $||F_t^{(0)} - Y_t||_2 \le \sqrt{2t}$  (Lemma 5.8) and  $\Upsilon/m \le 1$  with sufficiently large m, combining all the results above would give the conclusion.

LEMMA A.3. After T time steps, with the probability of at least  $1-O(TL) \cdot e^{-\Omega(m)}$  and running GD of J-iterations on the contexts and rewards, we have  $\beta_h' = \max\{\zeta \cdot \beta_1, \zeta \cdot \beta_2 + \zeta^2 \cdot \beta_1 \beta_2, (\zeta \cdot \beta_4)^{L-2}(\zeta \cdot \beta_2 + \zeta^2 \cdot \beta_1 \beta_2)\}$  and  $\beta_h = \max\{\zeta \cdot \beta_L + 1, \beta_h'\}$ . With  $h \in \{h_{gnn}, h_1, \ldots, h_{L-1}\}$ , we have

$$\|\boldsymbol{h}^{(j)} - \boldsymbol{h}^{(0)}\|_{2} \le \frac{\zeta \Upsilon}{m} \cdot \beta_{h} \cdot \frac{(2\zeta \beta_{L})^{L} - 1}{2\zeta \beta_{I} - 1} = \Lambda^{(j)}, \ \|\boldsymbol{h}^{(j)}\|_{2} \le \beta_{h} + \Lambda^{(j)}$$

**Proof.** Similar to the proof of **Lemma** A.2, we adopt an induction-based approach. For  $l \in [L-1]$ , we have

$$\begin{split} &\|\boldsymbol{h}_{l}^{(j)} - \boldsymbol{h}_{l}^{(0)}\|_{2} = \sqrt{\frac{1}{m}} \|\sigma(\boldsymbol{h}_{l-1}^{(j)} \cdot \boldsymbol{\Theta}_{l}^{(j)}) - \sigma(\boldsymbol{h}_{l-1}^{(0)} \cdot \boldsymbol{\Theta}_{l}^{(0)})\|_{2} \\ &\leq \sqrt{\frac{1}{m}} \zeta \cdot (\|\boldsymbol{h}_{l-1}^{(j)} \cdot \boldsymbol{\Theta}_{l}^{(j)} - \boldsymbol{h}_{l-1}^{(0)} \cdot \boldsymbol{\Theta}_{l}^{(j)}\|_{2} + \|\boldsymbol{h}_{l-1}^{(0)} \cdot \boldsymbol{\Theta}_{l}^{(j)} - \boldsymbol{h}_{l-1}^{(0)} \cdot \boldsymbol{\Theta}_{l}^{(0)}\|_{2}) \\ &\leq \sqrt{\frac{1}{m}} \zeta \cdot (\|\boldsymbol{\Theta}_{l}^{(0)}\|_{2} + \|\boldsymbol{\Theta}_{l}^{(j)} - \boldsymbol{\Theta}_{l}^{(0)}\|_{F}) \cdot \|\boldsymbol{h}_{l-1}^{(j)} - \boldsymbol{h}_{l-1}^{(0)}\|_{2} + \\ &\sqrt{\frac{1}{m}} \zeta \cdot \|\boldsymbol{h}_{l-1}^{(0)}\|_{2} \cdot \|\boldsymbol{\Theta}_{l}^{(j)} - \boldsymbol{\Theta}_{l}^{(0)}\|_{F} \\ &\leq \sqrt{\frac{1}{m}} \zeta \cdot (\beta_{L} \sqrt{m} + \Upsilon/\sqrt{m}) \cdot \|\boldsymbol{h}_{l-1}^{(j)} - \boldsymbol{h}_{l-1}^{(0)}\|_{2} + \zeta \cdot \beta_{h}' \cdot \Upsilon/m \\ &\leq \zeta \cdot (\beta_{L} + \Upsilon/m) \cdot \zeta \frac{\Upsilon}{m} \cdot \Lambda_{l-1}^{(j)} + \zeta \cdot \beta_{h}' \cdot \Upsilon/m \\ &\leq \zeta \frac{\Upsilon}{m} \cdot (\beta_{h} + 2\zeta\beta_{L} \cdot \Lambda_{l-1}^{(j)}) = \zeta \frac{\Upsilon}{m} \cdot \Lambda_{l}^{(j)} \end{split}$$

where the last two inequalities are derived by applying **Lemma** A.2 and the hypothesis. For the aggregation module output  $h_{ann}^{(0)}$ ,

$$\begin{aligned} & \| \boldsymbol{h}_{gnn}^{(j)} - \boldsymbol{h}_{gnn}^{(0)} \|_{2} = \sqrt{\frac{1}{m}} \| \sigma(\boldsymbol{\Theta}_{gnn}^{(j)} \cdot \boldsymbol{h}_{S}) - \sigma(\boldsymbol{\Theta}_{gnn}^{(0)} \cdot \boldsymbol{h}_{S}) \|_{2} \\ & \leq \frac{\zeta}{\sqrt{m}} \| \boldsymbol{\Theta}_{gnn}^{(j)} - \boldsymbol{\Theta}_{gnn}^{(0)} \|_{F} \cdot \| \boldsymbol{h}_{S} \|_{2} \leq \frac{\zeta \Upsilon}{m} \beta_{h}. \end{aligned}$$

Then, for the first layer l = 1, we have

$$\begin{aligned} & \|\boldsymbol{h}_{1}^{(j)} - \boldsymbol{h}_{1}^{(0)}\|_{2} = \sqrt{\frac{1}{m}} \|\sigma(\boldsymbol{x}' \cdot \boldsymbol{\Theta}_{1}^{(j)}) - \sigma(\boldsymbol{x}' \cdot \boldsymbol{\Theta}_{1}^{(0)})\|_{2} \\ & \leq \frac{\zeta}{\sqrt{m}} \|\boldsymbol{\Theta}_{gnn}^{(j)} - \boldsymbol{\Theta}_{1}^{(0)}\|_{F} \cdot \|\boldsymbol{x}'\|_{2} \leq \frac{\zeta \Upsilon}{m} \cdot (\zeta \cdot \beta_{L} + 1) \leq \frac{\zeta \Upsilon}{m} \cdot \beta_{h}. \end{aligned}$$

Combining all the results, for  $h \in \{h_{qnn}, h_1, ..., h_{L-1}\}$ , it has

$$\|\boldsymbol{h}^{(j)} - \boldsymbol{h}^{(0)}\|_{2} \le \frac{\zeta \Upsilon}{m} \cdot \beta_{h} \cdot \frac{(2\zeta \beta_{L})^{L} - 1}{2\zeta \beta_{L} - 1} = \Lambda^{(j)},$$

which completes the proof.

Lemma A.4. With initialized network parameters  $\Theta$  and the probability of at least  $1 - O(TL) \cdot e^{-\Omega(m)}$ , we have

 $\|\nabla_{\Theta}f(\widetilde{X};\Theta^{(0)})\|_F \leq \beta_h\beta_3 \cdot (\beta_L\zeta)^L/m, \|\nabla_{\Theta_L}f(\widetilde{X};\Theta^{(0)})\|_F \leq \beta_h/\sqrt{m},$  and the norm of gradient difference

$$\begin{split} \|\nabla_{\Theta}f(\widetilde{X};\Theta^{(0)}) - \nabla_{\Theta}f(\widetilde{X};\Theta^{(j)})\|_{F} &\leq 3 \cdot \Lambda^{(j)}, \\ \|\nabla_{\Theta_{L}}f(\widetilde{X};\Theta^{(0)}) - \nabla_{\Theta_{L}}f(\widetilde{X};\Theta^{(j)})\|_{F} &\leq \Lambda^{(j)}/\sqrt{m}. \\ with \Theta &\in \{\Theta_{qnn},\Theta_{1},\ldots,\Theta_{L-1}\}. \end{split}$$

**Proof.** First, for  $l \in [L-1]$ , we have

$$\begin{split} \|\nabla_{\Theta_{l}}f(\widetilde{X};\Theta^{(0)})\|_{F} &= m^{-\frac{L-l+1}{2}} \|\boldsymbol{h}_{l-1}^{(0)}(\Theta_{L}^{(0)}(\prod_{q=l+1}^{L-1}\Gamma_{q}\cdot\Theta_{q}^{(0)})\cdot\Gamma_{l})\|_{F} \\ &\leq m^{-\frac{L-l+1}{2}}\cdot\|\boldsymbol{h}_{l-1}^{(0)}\Theta_{L}^{(0)}\|_{F}\cdot\|\prod_{q=l}^{L-1}\Gamma_{q}\|_{2}\cdot\|\prod_{q=l+1}^{L-1}\Theta_{q}^{(0)}\|_{2} \\ &\leq m^{-\frac{L-l+1}{2}}\cdot\|\boldsymbol{h}_{l-1}^{(0)}\|_{2}\cdot\|\Theta_{L}^{(0)}\|_{2}\cdot\|\prod_{q=l}^{L-1}\Gamma_{q}\|_{2}\cdot\|\prod_{q=l+1}^{L-1}\Theta_{q}^{(0)}\|_{2} \\ &\leq m^{-\frac{L-l+1}{2}}\cdot\beta_{h}\beta_{3}\cdot\zeta^{L-l}\cdot(\beta_{L}\sqrt{m})^{L-l-1}\leq\beta_{h}\beta_{3}\cdot(\beta_{L}\zeta)^{L}/m. \end{split}$$

For  $\Theta_{qnn}$ , we can also derive similar results. For  $\Theta_L$ ,

$$\begin{split} \|\nabla_{\Theta_L} f(\widetilde{X}; \Theta^{(0)})\|_F &= m^{-0.5} \cdot \|\boldsymbol{h}_{L-1}^{(0)}\|_2 \leq \beta_h / \sqrt{m} \\ \text{Then, with } \nabla_l^{(j)} &= m^{-\frac{L-l+1}{2}} \cdot (\Theta_L^{(j)})^\intercal \cdot \left(\prod_{q=l+1}^{L-1} \Gamma_q \cdot (\Theta_q^{(j)})^\intercal\right) \cdot \Gamma_l, \\ \text{we have the norm of gradient difference} \\ \|\nabla_{\Theta_l} f(\widetilde{X}; \Theta^{(j)}) - \nabla_{\Theta_l} f(\widetilde{X}; \Theta^{(0)})\|_F &= \|\boldsymbol{h}_{l-1}^{(0)} \cdot \nabla_l^{(0)} - \boldsymbol{h}_{l-1}^{(j)} \cdot \nabla_l^{(j)} \|_F \\ &\leq \|\boldsymbol{h}_{l-1}^{(0)} \cdot \nabla_l^{(0)} - \boldsymbol{h}_{l-1}^{(j)} \cdot \nabla_l^{(0)} \|_F + \|\boldsymbol{h}_{l-1}^{(j)} \cdot \nabla_l^{(0)} - \boldsymbol{h}_{l-1}^{(j)} \cdot \nabla_l^{(j)} \|_F \\ &\leq \|\boldsymbol{h}_{l-1}^{(j)}\|_F \cdot \|\nabla_l^{(0)} - \nabla_l^{(j)}\|_F + \|\nabla_l^{(0)}\|_F \cdot \|\boldsymbol{h}_{l-1}^{(0)} - \boldsymbol{h}_{l-1}^{(j)} \|_F \\ &\leq (\beta_h + \Lambda^{(j)}) \cdot \|\nabla_l^{(0)} - \nabla_l^{(j)} \|_F + \Lambda^{(j)} \cdot \|\nabla_l^{(0)} \|_F. \end{split}$$

Here, for the difference of  $\nabla$ , we have

$$\begin{split} &\|\nabla_{l}^{(0)} - \nabla_{l}^{(j)}\|_{F} \\ &= m^{-\frac{L-l+1}{2}} \|\Theta_{L}^{(0)} \cdot \big(\prod_{q=l+1}^{L-1} \Gamma_{q} \cdot \Theta_{q}^{(0)}\big) \Gamma_{l} - \Theta_{L}^{(j)} \big(\prod_{q=l+1}^{L-1} \Gamma_{q} \Theta_{q}^{(j)}\big) \Gamma_{l}\|_{F} \\ &= m^{-\frac{1}{2}} \cdot \|\nabla_{l+1}^{(0)} \Gamma_{l} \cdot \Theta_{l+1}^{(0)} - \nabla_{l+1}^{(j)} \Gamma_{l} \cdot \Theta_{l+1}^{(j)}\|_{F} \\ &\leq \frac{\zeta}{\sqrt{m}} \cdot (\|\nabla_{l+1}^{(0)} \cdot \Theta_{l+1}^{(0)} - \nabla_{l+1}^{(0)} \cdot \Theta_{l+1}^{(j)}\|_{F} + \|\nabla_{l+1}^{(0)} \cdot \Theta_{l+1}^{(j)} - \nabla_{l+1}^{(j)} \cdot \Theta_{l+1}^{(j)}\|_{F}) \\ &\leq \frac{\zeta}{L_{CL}} \cdot (\|\nabla_{l+1}^{(0)}\|_{F} \|\cdot\Theta_{l+1}^{(0)} - \Theta_{l+1}^{(j)}\|_{F} + \|\Theta_{l+1}^{(j)}\|_{F} \|\nabla_{l+1}^{(0)} - \nabla_{l+1}^{(j)}\|_{F}). \end{split}$$

To continue the proof, we need to bound the term  $\|\nabla_I^{(0)}\|_F$  as

$$\|\nabla_{l}^{(0)}\|_{F} = m^{-0.5} \|\Gamma_{l}\Theta_{l+1}^{(0)} \cdot \nabla_{l+1}^{(0)}\|_{F} \le \zeta \beta_{L} \cdot \|\nabla_{l+1}^{(0)}\|_{F}.$$

Since for l = L - 1 we have

$$\|\nabla_{L-1}^{(0)}\|_F \le \frac{\zeta \beta_3}{m},$$

we can derive

$$\|\nabla_l^{(0)}\|_F \le \frac{\beta_3}{m} \cdot (\zeta \cdot \beta_L)^L \le 1$$

with sufficiently large m, and this bound also applies to  $\|\nabla_{gnn}^{(0)}\|_F$ . For

$$\|\nabla_{L}^{(0)} - \nabla_{L}^{(j)}\|_{F} = m^{-0.5} \|\boldsymbol{h}_{L-1}^{(0)} - \boldsymbol{h}_{L-1}^{(j)}\|_{F} \le \Lambda^{(j)} / \sqrt{m}$$

Therefore, we have

$$\|\nabla_{l}^{(0)} - \nabla_{l}^{(j)}\|_{F} \le \frac{\zeta \Upsilon}{m} + \zeta \cdot (\beta_{L} + \Upsilon/m) \|\nabla_{l+1}^{(0)} - \nabla_{l+1}^{(j)}\|_{F}.$$

By following a similar approach as in Lemma A.3, we will have

$$\|\nabla_l^{(0)} - \nabla_l^{(j)}\|_F \le \frac{\zeta \Upsilon}{m} \cdot \frac{(2\zeta \beta_L)^L - 1}{2\zeta \beta_L - 1} = \frac{\Lambda^{(j)}}{\beta_h}.$$

Therefore, we will have

$$\begin{split} \|\nabla_{\Theta_l} f(\widetilde{X}; \Theta^{(j)}) - \nabla_{\Theta_l} f(\widetilde{X}; \Theta^{(0)}) \|_F &\leq \left(\beta_h + \Lambda^{(j)}\right) \cdot \frac{\Lambda^{(j)}}{\beta_h} + \Lambda^{(j)} \\ &\leq \frac{\Lambda^{(j)}}{\beta_h} \cdot (2\beta_h + 1) = \Lambda^{(j)} \cdot (2 + \frac{1}{\beta_h}) \leq 3 \cdot \Lambda^{(j)} \end{split}$$

with sufficiently large m. This bound can also be derived for  $\|\nabla_{\Theta_{gnn}} f(\widetilde{X}; \Theta^{(0)})\|_2$  with a similar procedure. For L-th layer, we have

$$\|\nabla_{\Theta_{L}} f(\widetilde{X}; \Theta^{(j)}) - \nabla_{\Theta_{L}} f(\widetilde{X}; \Theta^{(0)})\|_{F} \le m^{-0.5} \cdot \|\boldsymbol{h}_{L-1}^{(0)} - \boldsymbol{h}_{L-1}^{(j)}\|_{F}$$

$$\le \Lambda^{(j)} / \sqrt{m},$$

which completes the proof.

Lemma A.5. With the probability of at least  $1 - O(TL) \cdot e^{-\Omega(m)}$ , we have the gradient for all the network as

$$\begin{split} & \|g(\widetilde{X};\Theta^{(0)})\|_2 \leq m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}, \\ & \|g(\widetilde{X};\Theta^{(j)})\|_2 \leq \Lambda^{(j)} \cdot \sqrt{9L + m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m} \\ & \|g(\widetilde{X};\Theta^{(0)}) - g(\widetilde{X};\Theta^{(j)})\|_2 \leq \Lambda^{(j)} \cdot \sqrt{9L + m^{-1}}. \end{split}$$

**Proof.** First, for the gradient before GD, we have

$$\begin{split} & \|g(\widetilde{X};\Theta^{(0)})\|_{2} = \sqrt{\|\nabla_{\Theta_{gnn}}f(\widetilde{X};\Theta^{(0)})\|_{2}^{2} + \sum_{l=1}^{L} \|\nabla_{\Theta_{l}}f(\widetilde{X};\Theta^{(0)})\|_{2}^{2}} \\ & \leq m^{-1}\beta_{h} \cdot \sqrt{L \cdot \beta_{2}^{2} \cdot (\beta_{L}\zeta)^{2L} + m}. \end{split}$$

Then, for the norm of gradients,  $\Theta \in \{\Theta_{gnn}, \Theta_1, \dots, \Theta_{L-1}\}$ , we have

$$\begin{split} &\|g(\widetilde{X};\Theta^{(0)}) - g(\widetilde{X};\Theta^{(j)})\|_2 \\ &= \sqrt{\sum_{\Theta}} \|\nabla_{\Theta} f(\widetilde{X};\Theta^{(0)}) - \nabla_{\Theta} f(\widetilde{X};\Theta^{(j)})\|_2^2 \\ &\leq \sqrt{9L \cdot (\Lambda^{(j)})^2 + (\Lambda^{(j)})^2/m} = \Lambda^{(j)} \cdot \sqrt{9L + m^{-1}}. \end{split}$$

Then, for the network gradient after GD, we have

$$\begin{split} &\|g(\widetilde{X};\Theta^{(j)})\|_{2} \leq \|g(\widetilde{X};\Theta^{(0)}) - g(\widetilde{X};\Theta^{(j)})\|_{2} + \|g(\widetilde{X};\Theta^{(0)})\|_{2} \\ &\leq \Lambda^{(j)} \cdot \sqrt{9L + m^{-1}} + m^{-1}\beta_{h} \cdot \sqrt{L \cdot \beta_{3}^{2} \cdot (\beta_{L}\zeta)^{2L} + m} \end{split}$$

LEMMA A.6. With the probability of at least  $1 - O(TL) \cdot e^{-\Omega(m)}$ , for the initialized parameter  $\Theta^{(0)}$ , we have

$$|f(\widetilde{X}; \Theta^{(j)}) - \langle g(\widetilde{X}; \Theta^{(0)}), \Theta^{(j)} - \Theta^{(0)} \rangle|$$

$$\leq m^{-0.5} \cdot (\Lambda^{(j)} (1 + \beta_3) + \beta_3 \beta_h + L \cdot \beta_h \Upsilon),$$

and for the network parameter after GD,  $\Theta^{(j)}$ , we have

$$|f(\widetilde{X}; \Theta^{(j)}) - \langle g(\widetilde{X}; \Theta^{(j)}), \Theta^{(j)} - \Theta^{(0)} \rangle| \le B_3$$
  
=  $m^{-0.5} (\beta_3(\Lambda^{(j)} + \beta_h) + L \cdot \Upsilon \cdot (\Lambda^{(j)} + \beta_h)(\Lambda^{(j)}/\beta_h + 1)).$ 

**Proof.** For the sake of enumeration, we let  $\Theta_0 = \Theta_{gnn}$ ,  $\nabla_0 = \nabla_{gnn}$ ,  $h_0 = h_{gnn}$  and  $h_{-1} = h_S$ . Then, we can derive

$$\begin{split} |f(\widetilde{X};\Theta^{(j)}) - \langle g(\widetilde{X};\Theta^{(0)}),\Theta^{(j)} - \Theta^{(0)} \rangle| &= |\frac{1}{\sqrt{m}} \langle \boldsymbol{h}_{L-1}^{(j)},\Theta_L^{(j)} \rangle \\ &- \frac{1}{\sqrt{m}} \langle \boldsymbol{h}_{L-1}^{(0)},\Theta_L^{(0)} - \Theta_L^{(j)} \rangle - \sum_{l=0}^{L-1} (\boldsymbol{h}_{l-1}^{(0)})^{\mathsf{T}} (\Theta_l^{(0)} - \Theta_l^{(j)}) \nabla_l^{(0)})| \\ &\leq m^{-0.5} \|\boldsymbol{h}_L^{(j)} - \boldsymbol{h}_L^{(0)}\|_2 \|\Theta_L^{(j)}\|_2 + m^{-0.5} \|\boldsymbol{h}_{L-1}^{(0)}\|_2 \|\Theta_L^{(0)}\|_2 \\ &+ \sum_{l=0}^{L-1} \|\boldsymbol{h}_{l-1}^{(0)}\|_2 \|\Theta_l^{(0)} - \Theta_l^{(j)}\|_F \|\nabla_l^{(0)}\|_F \\ &\leq m^{-0.5} \Lambda^{(j)} (\Upsilon/\sqrt{m} + \beta_3) + m^{-0.5} \beta_3 \beta_h + L \cdot \beta_h \frac{\Upsilon}{\sqrt{m}} \\ &\leq m^{-0.5} \cdot \left(\Lambda^{(j)} (1 + \beta_3) + \beta_3 \beta_h + L \cdot \beta_h \Upsilon\right). \end{split}$$

On the other hand, for network parameter after GD, we can have

$$\begin{split} |f(\widetilde{X};\Theta^{(j)}) - \langle g(\widetilde{X};\Theta^{(j)}),\Theta^{(j)} - \Theta^{(0)} \rangle| &= |\frac{1}{\sqrt{m}} \langle \pmb{h}_{L-1}^{(j)},\Theta_L^{(j)} \rangle \\ &- \frac{1}{\sqrt{m}} \langle \pmb{h}_{L-1}^{(j)},\Theta_L^{(j)} - \Theta_L^{(0)} \rangle - \sum_{l=0}^{L-1} (\pmb{h}_{l-1}^{(j)})^{\intercal} (\Theta_l^{(0)} - \Theta_l^{(j)}) \nabla_l^{(j)})| \\ &\leq |m^{-0.5} \langle \pmb{h}_{L-1}^{(j)},\Theta_L^{(0)} \rangle - \sum_{l=0}^{L-1} (\pmb{h}_{l-1}^{(j)})^{\intercal} (\Theta_l^{(0)} - \Theta_l^{(j)}) \nabla_l^{(j)}| \\ &\leq m^{-0.5} ||\pmb{h}_{L-1}^{(j)}||_2 ||\Theta_L^{(0)}||_2 + \sum_{l=0}^{L-1} ||\pmb{h}_{l-1}^{(j)}||_2 ||\Theta_l^{(0)} - \Theta_l^{(j)}||_F ||\nabla_l^{(j)}||_F \\ &\leq m^{-0.5} \beta_3 (\Lambda^{(j)} + \beta_h) + L \cdot (\Lambda^{(j)} + \beta_h) (\Upsilon/\sqrt{m}) (\Lambda^{(j)}/\beta_h + 1) \\ &\leq m^{-0.5} (\beta_3 (\Lambda^{(j)} + \beta_h) + L \cdot \Upsilon \cdot (\Lambda^{(j)} + \beta_h) (\Lambda^{(j)}/\beta_h + 1)). \end{split}$$

This completes the proof.

**Proof sketch for Lemmas A.1-A.6.** First we derive the conclusions in Lemma A.1 with the property of Gaussian matrices. Then, Lemmas A.2 and A.3 are proved through the induction after breaking the target into norms of individual terms (variables, weight matrices) and applying Lemma A.1. Finally, for Lemmas A.4-A.6, we also decompose targets into norms of individual terms. Then, applying Lemmas A.1-A.3 the to bound these terms (at random initialization / after GD) would give the result. ■

#### **B** LEMMAS FOR GRADIENT MATRICES

Inspired by [4, 46] and with sufficiently large network width m, the trained network parameter can be related to ridge regression estimator where the context is embedded by network gradients. With the received contexts and rewards up to time step t, we have the estimated parameter  $\widehat{\Theta}$  as  $\widehat{\Theta}_0 = (Z_0)^{-1} \cdot b_0$  where  $Z_0 = \lambda I + \frac{1}{m} \sum_{\tau=1}^t g(\widetilde{X}_\tau; \Theta_0) g(\widetilde{X}_\tau; \Theta_0)^\intercal, b_0 = \frac{1}{\sqrt{m}} \sum_{\tau=1}^t r_\tau \cdot g(\widetilde{X}_\tau; \Theta_0)$ . We also define the gradient matrix w.r.t. the network parameters as

$$G^{(j)} = (g(\widetilde{X}_1; \Theta^{(j)}), \dots, g(\widetilde{X}_t; \Theta^{(j)}))$$

$$f^{(j)} = (f(\widetilde{X}_1; \Theta^{(j)}), \dots, f(\widetilde{X}_t; \Theta^{(j)})), \quad \mathbf{r} = (r_1, \dots, r_t)$$

$$\Theta^{(j+1)} = \Theta^{(j)} - \eta \cdot ((G^{(j)})^{\mathsf{T}} (f^{(j)} - \mathbf{r})).$$

where the t notation is omitted by default. Then, we use the following Lemma to bound the above matrices.

Lemma B.1. After j iterations, with the probability of at least  $1 - O(L) \cdot e^{-\Omega(m)}$ , we have

$$\begin{split} &\|G^{(0)}\|_{F} \leq G_{1} = m^{-1}\beta_{h} \cdot \sqrt{t \cdot (L \cdot \beta_{3}^{2} \cdot (\beta_{L}\zeta)^{2L} + m)}, \\ &\|G^{(0)} - G^{(j)}\|_{F} \leq \Lambda^{(j)} \cdot \sqrt{t \cdot (9L + m^{-1})}, \\ &\|G^{(j)}\|_{F} \leq \widetilde{I}_{1} = \sqrt{t \cdot (L \cdot \beta_{3}^{2} \cdot (\beta_{L}\zeta)^{2L} + m)} + \Lambda^{(j)}\sqrt{t \cdot (9L + m^{-1})} \\ &\|f^{(j)} - (G^{(j)})^{\intercal} (\widehat{\Theta}^{(j)} - \widehat{\Theta}^{(0)})\|_{2} \leq \sqrt{t} \cdot B_{3} \\ &= \sqrt{t} \cdot m^{-0.5} (\beta_{3}(\Lambda^{(j)} + \beta_{b}) + L \cdot \Upsilon \cdot (\Lambda^{(j)} + \beta_{b})(\Lambda^{(j)}/\beta_{b} + 1)) \end{split}$$

**Proof.** For the gradient matrix after random initialization, we have

$$\|G^{(0)}\|_F = \sqrt{\sum_{\tau=1}^t \|g(\widetilde{X}_\tau; \Theta^{(0)})\|_2^2} \le m^{-1} \beta_h \cdot \sqrt{t \cdot L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}$$

with the conclusion from Lemma A.5. Then,

$$\begin{split} \|G^{(0)} - G^{(j)}\|_F &= \sqrt{\sum_{\tau=1}^t \|g(\widetilde{X}_\tau; \Theta^{(0)}) - g(\widetilde{X}_\tau; \Theta^{(j)})\|_2^2} \\ &\leq \Lambda^{(j)} \cdot \sqrt{t \cdot (9L + m^{-1})}. \end{split}$$

For the third inequality in this Lemma, we have

$$\begin{split} &\|f^{(j)} - (G^{(j)})^{\intercal}(\widehat{\Theta}^{(j)} - \widehat{\Theta}^{(0)})\|_2 \\ &= \sqrt{\sum_{\tau=1}^{t}} |f(\widetilde{X}_{\tau}; \Theta^{(j)}) - \langle g(\widetilde{X}_{\tau}; \Theta^{(j)}), \Theta^{(j)} - \Theta^{(0)} \rangle)|^2 \\ &\leq \sqrt{t} \cdot m^{-0.5} \big(\beta_3(\Lambda^{(j)} + \beta_h) + L \cdot \Upsilon \cdot (\Lambda^{(j)} + \beta_h)(\Lambda^{(j)}/\beta_h + 1)\big) \\ \text{based on Lemma A.6.} \end{split}$$

Analogous to [4, 46], we define another auxiliary sequence to bound the parameter difference. With  $\widetilde{\Theta}^{(0)} = \Theta^{(0)}$ , we have  $\widetilde{\Theta}^{(j+1)} = \Theta^{(0)}$ 

$$\widetilde{\Theta}^{(j)} - \eta \cdot \left( G^{(j)} \big( (G^{(j)})^{\mathsf{T}} (\widetilde{\Theta}^{(j)} - \widetilde{\Theta}^{(0)}) - r \big) + m \lambda (\widetilde{\Theta}^{(j)} - \widetilde{\Theta}^{(0)}) \right).$$

Lemma B.2. After j iterations, with the probability of at least  $1 - O(L) \cdot e^{-\Omega(m)}$ , we have

$$\|\widetilde{\Theta}^{(j)} - \Theta^{(0)} - \widehat{\Theta}_t / \sqrt{m}\|_2 \le \sqrt{t/(m\lambda)}$$

**Proof.** The proof is analogous to Lemma 10.2 in [4] and Lemma C.4 in [46]. Switching  $G_0$  to  $G_j$  would give the result.

Then, we can have the following lemma to bridge the difference between the regression estimator  $\widehat{\Theta}$  and the network parameter  $\Theta$ .

Lemma B.3. At this time step t, with the notation defined in Lemma 5.3 and the probability at least  $1 - O(L) \cdot e^{-\Omega(m)}$ , we will have

$$\|\Theta_t - \Theta_0 - \widehat{\Theta}_t / \sqrt{m}\|_2 \le (\widetilde{I}_1 \cdot \sqrt{t}B_3 + m \cdot \widetilde{I}_2)/(m\lambda) + \sqrt{t/(m\lambda)}$$
 with proper  $m, \eta$  as in Lemma 5.3

**Proof.** With an analogous approach from Lemma 6.2 in [4], we can have

$$\begin{split} &\|\widetilde{\Theta}^{(j+1)} - \Theta^{(j+1)}\|_2 \\ &\leq \eta \|G^{(j)}\|_2 \|f^{(j)} - (G^{(j)})^\intercal (\Theta^{(j)} - \Theta^{(0)})\|_2 + \eta m \lambda \|\Theta^{(0)} - \Theta^{(j)}\|_2 \\ &+ \|I - \eta \cdot (m \lambda I + G^{(j)} (G^{(j)})^\intercal)\|_2 \|\Theta^{(j)} - \widetilde{\Theta}^{(j)}\|_2 = I_1 + I_2 + I_3. \end{split}$$
 With Lemma B.1, we can bound them as

$$\begin{split} I_1 &\leq \eta \cdot \widetilde{I}_1 \cdot \sqrt{t} B_3 \\ I_2 &\leq \eta m \lambda \sqrt{\sum_{i=0}^{L} \|\boldsymbol{\Theta}_l^{(0)} - \boldsymbol{\Theta}_l^{(j)}\|_F^2} \leq \eta m \cdot \widetilde{I}_2 = \eta m \lambda \sqrt{L+1} \cdot \Upsilon / \sqrt{m}. \end{split}$$

based on the conclusion from Lemma A.2. For  $I_3$ , we have

$$\begin{split} \eta \cdot (m\lambda I + G^{(0)}(G^{(0)})^\intercal) &\leq \eta \cdot I \\ & \cdot \left( m\lambda + (m^{-1}\beta_h \cdot \sqrt{t \cdot (L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m)})^2 \right) \leq I \end{split}$$

with proper choice of m and  $\eta$ . It leads to

$$\begin{split} \|\widetilde{\Theta}^{(j+1)} - \Theta^{(j+1)}\|_2 &\leq (1 - \eta m \lambda) \|\widetilde{\Theta}^{(j)} - \Theta^{(j)}\|_2 + \widetilde{I}_1 \cdot \sqrt{t} B_3 + \eta m \cdot \widetilde{I}_2 \end{split}$$
 which by induction and  $\widetilde{\Theta}^{(0)} = \Theta^{(0)}$ , we have

$$\|\widetilde{\boldsymbol{\Theta}}^{(j)} - \boldsymbol{\Theta}^{(j)}\|_2 \le (\widetilde{I}_1 \cdot \sqrt{t}B_3 + m \cdot \widetilde{I}_2)/(m\lambda).$$

Finally,

$$\begin{split} &\|\Theta_t - \Theta_0 - \widehat{\Theta}_t / \sqrt{m}\|_2 \le \|\widetilde{\Theta}^{(j)} - \Theta^{(j)}\|_2 + \|\widetilde{\Theta}_t - \Theta_0 - \widehat{\Theta}_0 / \sqrt{m}\|_2 \\ &\le (\widetilde{I}_1 \cdot \sqrt{t}B_3 + m \cdot \widetilde{I}_2) / (m\lambda) + \sqrt{t/(m\lambda)}, \end{split}$$

which completes the proof.

Lemma B.4. At this time step t, with the probability at least  $1-O(L)\cdot e^{-\Omega(m)}$ , we will have

$$||Z_t||_2 < \lambda +$$

$$\begin{split} &\frac{t(L+1)}{m} \left(\Lambda^{(j)} \cdot \sqrt{9L+m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2} \cdot (\beta_L \zeta)^{2L} + m\right)^2, \\ &\|G_t^{\mathsf{T}} G_t - G_0^{\mathsf{T}} G_0\|_F \leq 2t \cdot m^{-1} (\Lambda^{(j)} \cdot \sqrt{9L+m^{-1}}) \\ &\cdot \left(\Lambda^{(j)} \sqrt{9L+m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}\right) = B_G/m \\ &\text{with proper } m, \eta \text{ as in Lemma 5.3.} \end{split}$$

**Proof.** For the gradient matrix of ridge regression, we have

$$\begin{split} \|Z_t\|_2 &\leq \lambda + m^{-1} \sum_{\tau=1}^t \|g(\widetilde{X}_\tau; \Theta_t)\|_2^2 \leq \lambda + \\ &\frac{t(L+1)}{m} \left(\Lambda^{(j)} \cdot \sqrt{9L + m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}\right)^2 \end{split}$$

with the results from Lemma A.5. Then,

$$\begin{split} &\|G_t^{\mathsf{T}}G_t - G_0^{\mathsf{T}}G_0\|_F \leq m^{-1} \cdot \\ &\sqrt{\sum_{i,j=1}^{t} \|g(\widetilde{X}_i;\Theta_t) + g(\widetilde{X}_j;\Theta_0)\|_2^2 + \|g(\widetilde{X}_i;\Theta_t) - g(\widetilde{X}_j;\Theta_0)\|_2^2} \\ &\leq 2t \cdot m^{-1} \cdot \left(\Lambda^{(j)}\sqrt{9L + m^{-1}} + m^{-1}\beta_h \cdot \sqrt{L \cdot \beta_3^2 \cdot (\beta_L \zeta)^{2L} + m}\right) \\ &(\Lambda^{(j)} \cdot \sqrt{9L + m^{-1}}) = B_G/m. \end{split}$$

The proof is then completed.

**Proof sketch for Lemmas B.1-B.4.** Analogous to lemmas in Section A, Lemma B.1 is proved by Lemmas A.5, A.6 by breaking the target into the product of norms. The proof of Lemma B.2 is analogous to Lemma 10.2 in [4] and Lemma C.4 in [46], then replacing  $G_0$  with  $G_j$  would give the result. Then, based on Lemma B.2 results, Lemma B.3 will be proved with after bounding  $\|\widetilde{\Theta}^{(j+1)} - \Theta^{(j+1)}\|_2$  by induction. Finally, Lemma B.4 is proved by decomposing the norm into sum of individual terms, and bounding these terms with bounds on gradients in Lemma A.5.  $\blacksquare$ 

# C LEMMAS FOR MODEL CONVERGENCE

LEMMA C.1. After T time steps, assume the model with width m defined in Lemma 5.3 are trained with the J-iterations GD on the past contexts and rewards. Then, there exists a constant  $\beta_F$ , such that  $\beta_F \cdot \eta < 1$ , for any  $j \in [J]$ :

$$\|V^{(j)}\|_2 \le \frac{1}{4}\eta \beta_F \cdot \|F_T^{(j)} - Y_T\|_2$$

where 
$$F^{(j)} = [f(\mathcal{G}_T, \widetilde{X}_\tau; \Theta^{(j)})]_{\tau=1}^T$$
, and  $Y_T = [r_\tau]_{\tau=1}^T$ .

**Proof.** We prove this lemma following an analogous approach as Lemma B.6 in [13]. Given  $\widetilde{X}$ , we denote  $\nabla \mathcal{L}(\Theta^{(j)}) = \frac{\partial \mathcal{L}(\Theta^{(j)})}{\partial \Theta}$ , and  $\nabla f(\Theta^{(j)}) = \frac{\partial f(\mathcal{G}_T, \widetilde{X}; \Theta^{(j)})}{\partial \Theta}$ , where  $\Theta \in \{\Theta_{gnn}, \Theta_1, \ldots, \Theta_L\}$ . By the definition of  $\|V^{(j)}\|$ , we have its element  $|V^{(j)}(\widetilde{X})|$ 

$$\leq \eta \cdot \max_{0 \leq s \leq \eta} \bigg[ \sum_{\Theta} \|\nabla \mathcal{L}(\Theta^{(j)})\|_F \|\nabla f(\Theta^{(j)}) - \nabla f(\Theta^{(j)}, s)\|_F \bigg].$$

With the notation and conclusion from Lemma A.2, we have

$$\|\nabla \mathcal{L}(\Theta^{(j)})\|_{F} \leq m^{-\frac{1}{2}} 2\sqrt{T} \|F_{T}^{(j)} - Y_{T}\|_{2} \cdot \zeta^{L} \cdot (2\beta_{L})^{L} \beta_{h}$$

$$\begin{split} &\text{Meantime, } \|\nabla f(\Theta_l^{(j)}) - \nabla f(\Theta_l^{(j)},s)\|_F = m^{-\frac{L-l+1}{2}} \\ &\|\boldsymbol{h}_{l-1}^{(j)}(\Theta_L^{(j)})^\intercal \big(\prod_{q=l+1}^{L-1} \Gamma_q^{(j)}(\Theta_q^{(j)})^\intercal \big) \cdot \Gamma_l^{(j)} - \boldsymbol{h}_{l-1}^{(j),s}(\Theta_L^{(j),s})^\intercal \\ &\big(\prod_{q=l+1}^{L-1} \Gamma_q^{(j),s} \cdot (\Theta_q^{(j),s})^\intercal \big) \cdot \Gamma_l^{(j),s}\|_F. \text{ A similar form can also be derived for } \boldsymbol{\Theta}_{qnn}. \end{split}$$

With  $\Upsilon/\sqrt{m} \le 1$  and  $\Lambda^{(j)} \le \beta_h$  and a similar procedure as in Lemma A.3 and Lemma A.2, we have

$$\begin{split} \|\Theta^{(j+1)} - \Theta^{(j)}\|_{F} &\leq \eta \frac{\Upsilon^{(j)}}{\sqrt{m}}, \quad \|\Theta^{(j)}\|_{F} \leq 2\beta_{L}\sqrt{m} \\ \|\boldsymbol{h}^{(j+1)} - \boldsymbol{h}^{(j)}\|_{2} &\leq \eta \frac{2\zeta\beta_{h}}{\sqrt{m}} (2\zeta\beta_{L})^{L} \Upsilon^{(j)}, \quad \|\boldsymbol{h}^{(j)}\|_{2} \leq 2\beta_{h}, \\ \|\Gamma^{(j+1)} - \Gamma^{(j)}\|_{F} &\leq 2\eta\zeta^{2}\beta_{h} (2\zeta\beta_{L})^{L} \Upsilon^{(j)}, \quad \|\Gamma^{(j)}\|_{2} \leq \zeta \end{split}$$

With Lemma G.1 from [13], for  $\Theta \in \{\Theta_{qnn}, \Theta_1, \dots, \Theta_L\}$ ,

$$\|\nabla f(\boldsymbol{\Theta}^{(j)}) - \nabla f(\boldsymbol{\Theta}^{(j)}, s)\|_F \leq \frac{4\zeta}{\sqrt{m}} \eta \Upsilon^{(j)} \beta_h L(2\zeta\beta_L)^{2L}.$$

Combining with  $\|\nabla \mathcal{L}(\Theta'^{(j)})\|_F$ , we have

$$|V^{(j)}(\widetilde{X})| \leq \eta^2 \frac{4T}{m} (L+2)^2 \beta_h^3 \cdot \|F_T^{(j)} - Y_T\|_2^2 (2\zeta\beta_L)^{4L}.$$

Since this inequality holds for an arbitrary  $\widetilde{X} \in \{\widetilde{X}_{\tau}\}_{\tau \in [T]}$  and  $\|F_T^{(0)} - Y_T\|_2 = O(\sqrt{T})$ , given network width m, we finally have

$$\|V^{(j)}\| \leq \frac{1}{4}\eta \beta_F \|F_T^{(j)} - Y_T\|_2^2.$$

with the choice of learning rate  $\eta \leq O(T^{-1}L^{-1}\beta_h^{-2}(2\zeta\beta_L)^{-2L})$ .

**Proof of Lemma 5.7.** We prove this lemma following an analogous approach as Lemma B.7 in [13]. By the model definition and substituting  $\Upsilon^{(j)}/\sqrt{m}$  with  $m^{-\frac{1}{2}}2\sqrt{T}\|F_T^{(j)}-Y_T\|_2\cdot\zeta^L(2\beta_L)^L\beta_h$  as the upper bound based on Lemma A.2, with  $\Lambda^{(j)}\leq\beta_h$ , we have

$$\begin{split} &\|\boldsymbol{F}_{T}^{(j)} - \boldsymbol{F}_{T}^{(j+1)}\|_{2}^{2} = \frac{1}{m} \sum_{\tau=1}^{T} \left( (\boldsymbol{h}_{L-1,\tau}^{(j+1)})^{\mathsf{T}} \boldsymbol{\Theta}_{L}^{(j+1)} - (\boldsymbol{h}_{L-1,\tau}^{(j)})^{\mathsf{T}} \boldsymbol{\Theta}_{L}^{(j)} \right)^{2} \\ &\leq \frac{2}{m} \left( \|\boldsymbol{\Theta}_{L}^{(j+1)} - \boldsymbol{\Theta}_{L}^{(j)}\|_{2}^{2} \sum_{\tau=1}^{T} \|\boldsymbol{h}_{L-1,\tau}^{(j+1)}\|_{2}^{2} + \|\boldsymbol{\Theta}_{L}^{(j)}\|_{2}^{2} \sum_{\tau=1}^{T} \|\boldsymbol{h}_{L-1,\tau}^{(j+1)} - \boldsymbol{h}_{L-1,\tau}^{(j)}\|_{2}^{2} \right) \\ &\leq \frac{2}{m} \left( \frac{T}{m} \eta^{2} (2\beta_{h})^{4} \|\boldsymbol{F}_{T}^{(j)} - \boldsymbol{Y}_{T}\|_{2}^{2} + T (2\beta_{3})^{2} (\eta \frac{2\zeta\beta_{h}}{\sqrt{m}} (2\zeta\beta_{L})^{L} \boldsymbol{\Upsilon}^{(j)})^{2} \right) \\ &\leq \frac{1}{4} \eta \beta_{F} \|\boldsymbol{F}_{T}^{(j)} - \boldsymbol{Y}_{T}\|_{2}^{2} \end{split}$$

where the last inequality is due to sufficiently large m and the choice of learning rate  $\eta$ .