

# TOWARDS DOMAIN ADAPTIVE NEURAL CONTEXTUAL BANDITS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Contextual bandit algorithms are essential for solving real-world decision making problems. In practice, collecting a contextual bandit’s feedback from different domains may involve different costs. For example, measuring drug reaction from mice (as a source domain) and humans (as a target domain). Unfortunately, adapting a contextual bandit algorithm from a source domain to a target domain with distribution shift still remains a major challenge and largely unexplored. In this paper, we introduce the first general domain adaptation method for contextual bandits. Our approach learns a bandit model for the target domain by collecting feedback from the source domain. Our theoretical analysis shows that our algorithm maintains a sub-linear regret bound even adapting across domains. Empirical results show that our approach outperforms the state-of-the-art contextual bandit algorithms on real-world datasets.

## 1 INTRODUCTION

Contextual bandit (CB) algorithms have shown great promise for naturally handling exploration/exploitation trade-off problems with optimal statistical properties. Notably, LinUCB (Li et al., 2010) and its various adaptations (Yue & Guestrin, 2011; Agarwal et al., 2014; Li et al., 2016; Kveton et al., 2017; Foster et al., 2018; Korda et al., 2016; Mahadik et al., 2020; Zhou et al., 2019), have been shown to be able learn the optimal strategy when all data come from the same domain. However, these methods fall short when applied to data from a new domain. For example, a drug reaction prediction model trained by collecting feedback from mice (the source domain) may not work for humans (the target domain).

So, how does one effectively explore a high-cost target domain by only collecting feedback from a low-cost source domain, e.g., exploring drug reaction in humans by collecting feedback from mice or exploring real-world environments by collecting feedback from simulated environments? The challenge of effective cross-domain exploration is multifaceted: (1) The need for *effective exploration* in both the source and target domains depends on the *quality of representations* learned in the source domain, which in turn requires *effective exploration*; this leads to a chicken-and-egg dilemma. (2) Aligning source-domain representations with target-domain representations is nontrivial in bandit settings, where ground-truth may still be unknown if the action is incorrect. (3) Balancing these aspects – effective exploration and accurate alignment in bandit settings – is also nontrivial.

To address these challenges, as the first step, we allow our method to simultaneously perform *effective exploration* and *representation alignment*, leveraging unlabeled data from both the source and target domains. Interestingly, our theoretical analysis reveals that naively doing so leads to sub-optimal accuracy/regret (verified by our empirical results) and naturally leads to additional terms in the target-domain regret bound. We then follow the regret bound derived from our analysis to develop an algorithm that adaptively collect feedback from the source domain while aligning representations from the source and target domains. Our contributions are outlined as follows:

- We identify the problem of contextual bandits across domains and propose domain-adaptive contextual bandits (DABand) as the first general method to explore a high-cost target domain while only collecting feedback from a low-cost source domain.
- Our theoretical analysis shows that our method can achieve a sub-linear regret bound in the target domain.

- Our empirical results on real-world datasets show our DABand significantly improve performance over the state-of-the-art contextual bandit methods when adapting across domains.

## 2 RELATED WORK

**Contextual Bandits.** In the realm of adaptive decision-making, contextual bandit algorithms, epitomized by LinUCB (Li et al., 2010), have carved a niche in efficiently balancing the exploitation-exploration paradigm across a spectrum of use cases, notably in complex adaptive systems such as recommender systems (Li et al., 2010; Wang et al., 2022). These algorithmic frameworks, along with their myriad adaptations (Yue & Guestrin, 2011; Agarwal et al., 2014; Li et al., 2016; Korda et al., 2016; Kveton et al., 2017; Foster et al., 2018; Zhou et al., 2019; Mahadik et al., 2020; Xu et al., 2020), have decisively outperformed conventional bandit models that operate devoid of contextual awareness (Auer, 2002). This superiority is underpinned by theoretical models, paralleling the insights in (Auer, 2002), where LinUCB variants are validated to conform to optimal regret boundaries in targeted scenarios (Chu et al., 2011). However, a discernible gap in these methodologies is their reduced efficacy in a new domain, particularly when getting feedback involves high cost. Our proposed DABand bridges this gap by adeptly aligning both source-domain and target-domain representations in the latent space, thereby reducing performance drop when transfer across different domains.

**Domain Adaptation.** The landscape of domain adaptation has been extensively explored, as evidenced by a breadth of research (Pan & Yang, 2009; Pan et al., 2010; Long et al., 2018; Saito et al., 2018; Sankaranarayanan et al., 2018; Zhang et al., 2019; Peng et al., 2019; Chen et al., 2019; Dai et al., 2019; Wang et al., 2020; Nguyen-Meidine et al., 2021; Xu et al., 2023). The primary objective of these studies has been the alignment of source and target domain distributions to facilitate the effective generalization of models trained on labeled source data to unlabeled target data. This alignment is conventionally attained either through the direct matching of distributional statistics (Pan et al., 2010; Tzeng et al., 2014; Sun & Saenko, 2016; Peng et al., 2019; Nguyen-Meidine et al., 2021) or via the integration of an adversarial loss (Ganin et al., 2016b; Zhao et al., 2017; Tzeng et al., 2017; Zhang et al., 2019; Kuroki et al., 2019; Chen et al., 2019; Dai et al., 2019; Wang et al., 2020). The latter, known as adversarial domain adaptation, has surged in prominence, bolstered by its theoretical foundation (Goodfellow et al., 2014; Zhao et al., 2018; Zhang et al., 2019; Zhao et al., 2019), the adaptability of end-to-end training in neural network architectures, and its empirical efficacy. However, these methods only work in offline settings, and assumes complete observability of labels in the source domain. Therefore they are not applicable to our online bandit settings.

**Domain Adaptation Related to Bandits.** There is also work related to both domain adaptation and bandits. Specifically, Guo et al. (2020) propose a domain adaptation method using a bandit algorithm to select which domain to use during training.

We note that their goal and setting is different from our DABand’s. (1) Guo et al. (2020) focuses on improving accuracy in a typical, offline domain adaptation setting. In contrast, our DABand focuses on minimizing regret in an online bandit setting. (2) Guo et al. (2020) assumes complete access to ground-truth labels in the source domain, while in our bandit setting, ground-truth may still be unknown if the action is incorrect. Therefore their work is not applicable to our setting.

## 3 THEORY

In this section, we formalize the problem of contextual bandits across domains and derive the corresponding regret bound. We then develop our DABand inspired by this bound in Sec. 4. **All proofs of lemmas, theorems can be found in the Appendix.**

**Notation.** We use  $[k]$  to denote the set  $\{1, 2, \dots, k\}$ , for  $k \in \mathbb{N}^+$ . The Euclidean norm of a vector  $\mathbf{x} \in \mathbb{R}^d$  is denoted as  $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$ . We denote the operator norm and Frobenius norm of a matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$  as  $\|\mathbf{W}\|$  and  $\|\mathbf{W}\|_F$ , respectively. Given a semi-definite matrix  $A \in \mathbb{R}^{d \times d}$  and a vector  $\mathbf{x} \in \mathbb{R}^d$ , we denote the Mahalanobis norm as  $\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^T A \mathbf{x}}$ . For a function  $f(T)$  of a parameter  $T$ , we denote as  $\mathcal{O}(f(T))$  the terms growing in the order of  $f(T)$ , ignoring constant factors. We assume all the action spaces are identical and with cardinality of  $K$ , i.e.,  $|\mathcal{A}_1| = |\mathcal{A}_2| = \dots = |\mathcal{A}_N| = K$ .

We use  $\langle \cdot, \cdot \rangle$  to denote the inner product of two vectors. We denote as  $a_{\mathcal{D},i}$  the action  $i$  in domain  $\mathcal{D}$ , and omit the subscript  $\mathcal{D}$  when the context is clear, e.g.,  $x_{i,a_i}^{\mathcal{D}} \equiv x_{i,a_i}^{\mathcal{D}}$ .

### 3.1 PRELIMINARIES

**Typical Contextual Bandit Setting.** Suppose we have  $N$  samples from an unknown domain  $\mathcal{D}$ , which is represented as  $\{\mathbf{x}_{i,a}^{\mathcal{D}}, a_i^*\}_{i \in [N], a \in [K]}$ . Here,  $\mathbf{x}_{i,a} \in \mathcal{D}$  denotes the context for action  $a \in \mathcal{A}_i$ , and  $a_i^* \in \mathcal{A}_i$  is the ground-truth action that will receive the optimal reward. We denote as  $r_{i,a}$  the received reward in round  $i$  after performing action  $a$  and  $\hat{a}_i \in \mathcal{A}_i$  the chose action by a bandit algorithm in round  $i$ . The goal in typical contextual bandits is to learn a policy of choosing actions  $\hat{a}_i$  in each round to minimize the regret after  $N$  rounds:

$$R = \sum_{i=1}^N r_{i,a_i^*} - \sum_{i=1}^N r_{i,\hat{a}_i}. \quad (1)$$

**LinUCB.** LinUCB [Li et al. \(2010\)](#) is a classic method for the typical contextual bandit setting. It assumes that the reward is linear to its input context, i.e.,  $\mathbf{x}_{i,a}$ . During each round  $i$ , the agent selects an arm according to historical contexts and their corresponding rewards,  $r_{i,a}$ :

$$\hat{a}_i = \operatorname{argmax}_{a \in \mathcal{A}_i} \left\{ \mathbf{x}_{i,a}^T \theta_i + \alpha \|\mathbf{x}_{i,a}\|_{A_{i-1}^{-1}} \right\},$$

where  $A_{i-1} = \gamma \mathbf{I}_d + \sum_{j=1}^{i-1} \mathbf{x}_{j,\hat{a}_j} \mathbf{x}_{j,\hat{a}_j}^T$  with  $\gamma > 0$  and  $\mathbf{x}_{j,\hat{a}_j}^T$  as the historical selected action's context,  $\mathbf{x}_{i,a}$  is the context for the candidate action  $a$  in round  $i$ ,  $\theta_i$  is the bandit parameter in round  $i$  and  $\alpha > 0$  is a hyperparameter that adjusts the exploration rate.

**LinUCB with Representation Learning.** In this paper, we use Neural LinUCB [Xu et al. \(2020\)](#); [Wang et al. \(2022\)](#) as a backbone model to handle high-dimensional context with representation learning. We assume that there exist a ground-truth encoder  $\phi^*$  to encode the raw context  $\mathbf{x}_{i,a}$  into a latent-space representation (encoding)  $\phi^*(\mathbf{x}_{i,a})$ . Subsequently, the reward for the context  $\mathbf{x}_{i,a}$  is then  $r(\mathbf{x}_{i,a}) = \langle \theta^*, \phi^*(\mathbf{x}_{i,a}) \rangle$  plus some stochastic noise  $\epsilon$ . Furthermore, we use the same setting as in [Xu et al. \(2020\)](#) where the ground-truth rewards are restricted to the range of  $[0, 1]$ . Additionally, we further bound predicted rewards by normalizing  $\|\hat{\theta}\| = 1$  and  $\|\hat{\phi}(\mathbf{x}_{i,a})\| = 1$ .

In the single-domain typical setting, the goal is to learn an encoder  $\hat{\phi}$  and the contextual bandit parameter  $\hat{\theta}$ , such that our estimated reward  $\hat{r}(\mathbf{x}_{i,a}) = \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a}) \rangle$  can be close to the ground-truth reward  $r(\mathbf{x}_{i,a})$ . One can then use this to form derive policies to minimize the regret in Eqn. (1).

For simplicity, in Definition 3.1 below, we further denote as  $f_{\mathcal{D}}$  the (ground-truth) labeling function for domain  $\mathcal{D}$  and as  $h \in \mathcal{H}$  a hypothesis such that  $f$  is parameterized by  $\phi_{\mathcal{D}}^*, \theta_{\mathcal{D}}^*$ , and  $h$  is parameterized by  $\hat{\phi}, \hat{\theta}$ .

### 3.2 OUR CROSS-DOMAIN CONTEXTUAL BANDIT SETTING

Typical contextual bandits operate in a single domain  $\mathcal{D}$ . In contrast, our cross-domain bandit setting involves a source domain  $\mathcal{S}$  and a target domain  $\mathcal{T}$ . In this setting, one can only collect feedback (reward) from the source domain, but not from the target domain. Specifically: (1) We assume a low-cost *source domain* (experiments on mice), where for each round  $i$ , one has access to the contexts  $\{\mathbf{x}_{i,a}^{\mathcal{S}}\}_{a \in [K]}$  for each candidate actions, chooses one action  $\hat{a}_i$ , and receives reward  $r_{i,\hat{a}_i}^{\mathcal{S}}$ . (2) Additionally, we assume a high-cost *target domain* where collecting feedback (reward) is expensive (experiments on humans); therefore, one only has access to the contexts  $\{\mathbf{x}_{i,a}^{\mathcal{T}}\}_{a \in [K]}$  for each candidate actions for each round, but **cannot collect feedback (reward)  $r_{i,\hat{a}_i}^{\mathcal{T}}$  for any action in the target domain**. The goal in our cross-domain contextual bandits is to learn a policy of choosing actions  $\hat{a}_i$  in the target domain to minimize the target-domain **zero-shot regret** for  $N$  (hypothetical) future rounds:

$$R_{\mathcal{T}} = \sum_{i=1}^N r_{i,a_i^*}^{\mathcal{T}} - \sum_{i=1}^N r_{i,\hat{a}_i}^{\mathcal{T}}. \quad (2)$$

**Difference between Eqn. (1) and Eqn. (2).** In the typical setting, Eqn. (1) uses different updated policies for each round; this is also true for the source domain. In contrast, the target regret in Eqn. (2)

use the same fixed policy obtained from the source domain for all  $N$  (hypothetical) rounds because one cannot collect feedback from the target domain. See the difference between Definition 3.4 and Definition 3.5 below for a more formal comparison.

### 3.3 FORMAL DEFINITIONS OF ERROR

**Definition 3.1 (Labeling Function and Hypothesis).** We define the labeling function  $f_{\mathcal{D}}$  and the hypothesis  $h$  for domain  $\mathcal{D}$  as follows:

$$\begin{aligned} f_{\mathcal{D}}(\mathbf{x}_{i,a}^{\mathcal{D}}) &= r(\mathbf{x}_{i,a}^{\mathcal{D}}) = \langle \theta_{\mathcal{D}}^*, \phi_{\mathcal{D}}^*(\mathbf{x}_{i,a}^{\mathcal{D}}) \rangle + \epsilon_i \\ h(\mathbf{x}_{i,a}^{\mathcal{D}}) &= \hat{r}(\mathbf{x}_{i,a}^{\mathcal{D}}) = \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a}^{\mathcal{D}}) \rangle + \epsilon_i \end{aligned}$$

where  $\theta_{\mathcal{D}}^*$  is the optimal predictor for domain  $\mathcal{D}$ ,  $\hat{\theta}$  denotes the estimated predictor, and  $\epsilon_i$  is the random noise.

Next, we analyze how well our hypothesis  $h$  estimates the labeling function  $f$  (i.e., ground-truth hypothesis). Definition 3.2 below quantifies the closeness between two hypotheses by calculating the absolute difference in their estimated rewards.

**Definition 3.2 (Estimated Error between Two Hypotheses).** Assuming all contexts  $\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}$  are from domain  $\mathcal{D}$ , the error between two hypotheses  $h_1, h_2 \in \mathcal{H}$  on domain  $\mathcal{D}$  given selected actions  $\{\hat{a}_i\}_{i \in [N]} \in [K]$  is

$$\epsilon_{\mathcal{D}}(h_1, h_2) = \sum_{i=1}^N \left( |h_1(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}) - h_2(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}})| \right) = \sum_{i=1}^N \left( \left| \langle \hat{\theta}_1, \hat{\phi}_1(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}) \rangle - \langle \hat{\theta}_2, \hat{\phi}_2(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}) \rangle \right| \right). \quad (3)$$

Note that the domain  $\mathcal{D}$  above can be the source domain  $\mathcal{S}$  or the target domain  $\mathcal{T}$ . We then define the error of our estimated hypothesis in these domains.

**Definition 3.3 (Source- and Target-Domain Error).** With  $N$  samples from the source domain  $\mathcal{S}$ , and  $f_{\mathcal{S}}$  denoting the labeling function for  $\mathcal{S}$ , the source-domain error is then

$$\epsilon_{\mathcal{S}}(f_{\mathcal{S}}, h) = \sum_{i=1}^N \left( |f_{\mathcal{S}}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) - h(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}})| \right) = \sum_{i=1}^N \left( \left| \langle \theta_{\mathcal{S}}^*, \phi_{\mathcal{S}}^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right),$$

where  $\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}$  comes from  $\mathcal{S}$ . Furthermore,  $\hat{a}_i = \arg \max_a h(\mathbf{x}_{i,a}^{\mathcal{S}})$  and  $a_i^* = \arg \max_a f_{\mathcal{S}}(\mathbf{x}_{i,a}^{\mathcal{S}})$ . For simplicity, we shorten the notation  $\epsilon_{\mathcal{S}}(f_{\mathcal{S}}, h)$  to  $\epsilon_{\mathcal{S}}(h)$  and similarly use  $\epsilon_{\mathcal{T}}(h)$  for domain  $\mathcal{T}$ . Assuming  $\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}$  comes from  $\mathcal{T}$ , the estimated error for the target domain is then:

$$\epsilon_{\mathcal{T}}(f_{\mathcal{T}}, h) = \sum_{i=1}^N \left( |f_{\mathcal{T}}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) - h(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}})| \right) = \sum_{i=1}^N \left( \left| \langle \theta_{\mathcal{T}}^*, \phi_{\mathcal{T}}^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right).$$

### 3.4 SOURCE REGRET AND TARGET REGRET

Below define the regret for the source and target domains.

**Definition 3.4 (Source Regret).** Assuming  $\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}$  comes from domain  $\mathcal{S}$ , the source regret (i.e., the regret in the source domain) is

$$R_{\mathcal{S}} = \sum_{i=1}^N \left( f_{\mathcal{S}}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) - f_{\mathcal{S}}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \right) = \sum_{i=1}^N \left( \langle \theta_{\mathcal{S}}^*, \phi_{\mathcal{S}}^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \theta_{\mathcal{S}}^*, \phi_{\mathcal{S}}^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right). \quad (4)$$

The goal in our cross-domain contextual bandits is to learn a policy of choosing actions  $\hat{a}_i$  in the high-cost target domain (e.g., human experiments) by collecting feedback only in the source domain (e.g., mouse experiments). Formally, we would like to minimize the target regret as defined below.

**Definition 3.5 (Target Regret and Problem Formulation).** Denoting the estimated hypothesis as  $\hat{h} = \{\hat{\phi}, \hat{\theta}\}$ , the target regret we aim to minimize is defined as

$$\begin{aligned} R_{\mathcal{T}} &= \sum_{i=1}^N \left( f_{\mathcal{T}}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) - f_{\mathcal{T}}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right) \\ \text{s.t. } \hat{a}_i &= \arg \max_a \hat{h}(\mathbf{x}_{i,a}^{\mathcal{T}}) + \alpha \|\hat{\phi}(\mathbf{x}_{i,a}^{\mathcal{T}})\|_{[AS]-1}, \quad \hat{h} = \arg \min_h \epsilon_{\mathcal{S}}(h), \end{aligned}$$

where  $A^S = \gamma \mathbf{I} + \sum_{i=1}^N [\hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S)] [\hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S)]^T$  denotes the context matrix accumulated by the selected context features in the source domain.

### 3.5 CROSS-DOMAIN ERROR BOUND FOR REGRESSION

Prior to introducing our final regret bound, it is necessary to define an additional component. A fundamental challenge in general domain adaptation problems is to manage the divergence between source and target domains. Unfortunately, previous domain adaptation theory only covers *classification* tasks Ben-David et al. (2010); Zhang et al. (2019). In contrast, the problem of contextual bandits is essentially a reward *regression* problem. To address this challenge, we introduce the following new definition to bound the error between the source and target domains.

**Definition 3.6 ( $\mathcal{H}\Delta\mathcal{H}$  Hypothesis Space for Regression).** For a hypothesis space  $\mathcal{H}$ , the symmetric difference hypothesis space  $\mathcal{H}\Delta\mathcal{H}$  is the set of hypotheses s.t.

$$g \in \mathcal{H}\Delta\mathcal{H} \iff g(\mathbf{x}) = |h(\mathbf{x}) - h'(\mathbf{x})|.$$

We then can further define the Divergence for  $\mathcal{H}\Delta\mathcal{H}$  Hypothesis Space:

$$\hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) = 2 \sup_{h, h' \in \mathcal{H}} |\mathbb{E}_{\mathbf{x} \sim \mathcal{S}}[|h(\mathbf{x}) - h'(\mathbf{x})|] - \mathbb{E}_{\mathbf{x} \sim \mathcal{T}}[|h(\mathbf{x}) - h'(\mathbf{x})|]|$$

**Definition 3.7 (Optimal Hypothesis).** The optimal hypothesis, denoted as  $h^*$ , can balance the estimated error for both source and target domain. Formally,

$$h^* = \arg \min_{h \in \mathcal{H}} \epsilon_S(h) + \epsilon_T(h),$$

and we define the minimum estimated error for the optimal hypothesis  $h^*$  as

$$\psi = \epsilon_S(h^*) + \epsilon_T(h^*).$$

### 3.6 FINAL REGRET BOUND

With all the definitions of source/target regret in Definition 3.4 and Definition 3.5, we can then bound the target regret using Theorem 3.1 below.

**Theorem 3.1 (Target Regret Bound).** Denoting the contexts from the source domain  $\mathcal{S}$  as  $\{\mathbf{x}_{t,a}^S\}_{i \in [N], a \in [K]}$ , the associated ground-truth action as  $\{a_i^*\}_{i=1}^N$ , and the contexts from the target domain  $\mathcal{T}$  as  $\{\mathbf{x}_{t,a}^T\}_{i \in [N], a \in [K]}$ , the upper bound for our target regret  $R_T$  is

$$\begin{aligned} R_T &\triangleq \sum_{i=1}^N \left( \left| \langle \theta_T^*, \phi_T^*(\mathbf{x}_{i,a_i^*}^T) \rangle - \langle \theta_T^*, \phi_T^*(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) \\ &\leq \underbrace{R_S}_{\text{Source Regret}} + \underbrace{2 \cdot \epsilon_S(h)}_{\text{Regression Error}} + \underbrace{N \cdot \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T})}_{\text{Data Divergence}} + \underbrace{\psi + C}_{\text{Constant}} \\ &\quad + \underbrace{\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right)}_{\text{Predicted Rewards}}, \end{aligned} \quad (5)$$

where  $R_S$ ,  $\epsilon_S(h)$ , and  $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}$  are the source regret, source-domain error, and  $\mathcal{H}\Delta\mathcal{H}$  divergence defined in Definition 3.4, Definition 3.3, and Definition 3.6, respectively.  $\psi$  is a constant independent to the problem and  $C$  is a constant which can be ignored (see the Appendix for more details).

Here we provide several observations (for more analysis, please refer to the Appendix):

- (1) Since  $R_S$  is sub-linear w.r.t. the number of samples in the source domain, so is  $R_T$ .
- (2) Theorem 3.1 bounds the target regret using the source regret, enabling the exploration of the target domain by only collecting feedback (reward) from the source domain.
- (3) Typical generalization bounds in DA naively minimize the source regret  $R_S$  while aligning source and target data in the latent space (i.e., minimizing  $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}$ ). This is **not sufficient**, as we need two additional terms, i.e., the regression error  $\epsilon_S(h)$  and the predicted reward  $\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right)$ , as regularization terms.



- (4) Minimizing the regression error  $\epsilon_S(h)$  encourages accurate prediction of source rewards.
- (5) Minimizing the predicted reward  $\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right)$  regularizes the model to avoid overestimating rewards.

Note that Theorem 3.1 is **nontrivial**. While it does resemble the generalization bound in domain adaptation, there are key differences. As mentioned in Observation (3) above, our target regret bound includes two additional crucial terms not found in domain adaptation. Specifically:

- **Regression Error in the Source Domain.**  $\sum_{i=1}^N \left( \left| \langle \theta_S^*, \phi_S^*(x_{i,\hat{a}_i}^S) \rangle - \langle \hat{\theta}, \hat{\phi}(x_{i,\hat{a}_i}^S) \rangle \right| \right)$  (i.e.,  $\epsilon_S(h)$  in Theorem 3.1) defines the difference between the true reward from selecting action  $\hat{a}_i$  and the estimated reward for this action.
- **Predicted Reward.**  $\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right)$  serves as a regularization term to regularize the model to avoid overestimating rewards.

The results of the ablation study in Table 4 in Sec. 5 highlight the significance of these two terms. See Appendix G.4 for more discussion on novelty as well as key differences between our DABand and classic domain adaptation.

## 4 METHOD

With Theorem 3.1, we can then design our DABand algorithm to obtain optimal target regret. We discuss how to translate each term in Eqn. (5) to a differential loss term in Sec. 4.1 and then put them together in Sec. 4.2.

### 4.1 FROM THEORY TO PRACTICE: TRANSLATING THE BOUND IN EQN. (5) TO DIFFERENTIABLE LOSS TERMS

**Source Regret.** Inspired by Neural-LinUCB (Xu et al., 2020), we use LinUCB Li et al. (2010) to update  $\hat{\theta}$  and the following empirical loss function when updating the encoder parameters  $\hat{\phi}$  by back-propagation in round  $i$ :

$$\mathcal{L}_i^{RS} = \sum_{a=1}^K \left( \hat{\theta}^T \hat{\phi}(\mathbf{x}_{i,a}) - r(\mathbf{x}_{i,a}) \right)^2, \quad (6)$$

where  $\theta$  is the contextual bandit parameter, and  $\hat{\phi}$  is the encoder shared by the source and target domains, i.e., we set  $\hat{\phi}_S = \hat{\phi}_T$  during training.

**Insights.** From Eqn. (6), a fundamental difference between classification problems and bandit problems becomes apparent. In classification tasks, both the ground-truth label and the estimated label are accessible. However, in bandit problems, we only know the estimated label and reward. If the estimated rewards align with our estimated label’s correctness, it suggests knowledge of the ground-truth label for that feature. If not, we can only deduce that the estimated label is inaccurate, without identifying the correct label among the other candidates. This uncertainty significantly amplifies the complexity of contextual bandit problems.

**Regression Error.** We use  $L_1$  loss to optimize the regression error, defined as the difference between the true reward from selecting action  $\hat{a}_i$  and the estimated reward for this action (see Definition 3.3 for a detailed explanation of this regression error). Minimizing the *regression error* term directly encourages accurate prediction of source rewards using  $L_1$  loss. Specifically, we use

$$\mathcal{L}_i^{\epsilon_S(h)} = \left| \langle \theta_S^*, \phi_S^*(\mathbf{x}_{i,\hat{a}_i}^S) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right|, \quad (7)$$

where  $\langle \theta_S^*, \phi_S^*(\mathbf{x}_{i,\hat{a}_i}^S) \rangle$  is the reward received by the agent in the source domain. Minimizing the regression error term above directly encourages accurate prediction of source rewards.

**Data Divergence.** The data divergence term in Eqn. (5) leads to following loss term:

$$\mathcal{L}^{div} = \max_g N \cdot (\mathbb{E}^S[\mathcal{L}_D(g(\hat{\phi}(\mathbf{x})), 0)] + \mathbb{E}^T[\mathcal{L}_D(g(\hat{\phi}(\mathbf{x})), 1)]), \quad (8)$$

**Algorithm 1** DABand Training Algorithm

---

```

1: Input: regularization parameter  $\gamma > 0$ , number of total steps  $N$ , episode length  $H$ , exploration
   parameters  $\alpha$ .
2: Output: parameters of the model:  $\hat{\phi}_N, \hat{\theta}_N, A_N, b_N$ .
3: Initialization:  $A_0 = \gamma \mathbf{I}$ ,  $b_0 = \mathbf{0}$ , and for  $\hat{\theta}_0$  and  $\hat{\phi}_0$  are initialized following (Xu et al., 2020).
4: for  $i = 1, \dots, N$  do
5:   Obtain  $\{\mathbf{x}_{i,a}^S\}_{a \in [K]}$  and  $\{\mathbf{x}_{i,a}^T\}_{a \in [K]}$  from the source and the target domain, respectively.
6:   Choose an action  $\hat{a}_i = \operatorname{argmax}_{a \in [K]} \left[ \hat{\phi}_{i-1}(\mathbf{x}_{i,a}^S) \right] \hat{\theta}_{i-1} + \alpha \left\| \left[ \hat{\phi}_{i-1}(\mathbf{x}_{i,a}^S) \right] \right\|_{A_{i-1}^{-1}}$ .
7:   Get the reward  $r_i = \mathbf{r}(\mathbf{x}_{i,\hat{a}_i})$  based on the selected action  $\hat{a}_i$ .
8:   Update bandit parameters:
9:   
$$\begin{cases} A_i = A_{i-1} + \left[ \hat{\phi}_{i-1}(\mathbf{x}_{i,\hat{a}_i}^S) \right] \left[ \hat{\phi}_{i-1}(\mathbf{x}_{i,\hat{a}_i}^S) \right]^T, \\ b_i = b_{i-1} + r_i \left[ \hat{\phi}_{i-1}(\mathbf{x}_{i,\hat{a}_i}^S) \right], \\ \hat{\theta}_i = A_i^{-1} b_i \end{cases}$$

10:  if  $\operatorname{mod}(i, H) = 0$  then
11:    for  $j = 1, \dots, H$  do
12:      Calculate  $\mathcal{L}^{DABand}$  in Eqn. (11).
13:      Use the Adam optimizer (Diederik, 2014) to update encoder  $\hat{\phi}_i$  and discriminator  $g$  by
        back-propagation to solve the minimax optimization in Eqn. (10).
14:    end for
15:  else
16:     $\hat{\phi}_i = \hat{\phi}_{i-1}$ 
17:  end if
18: end for
19: Output:  $\hat{\phi}_N, \hat{\theta}_N, A_N, b_N$ .

```

---

where  $\mathbb{E}^S$  and  $\mathbb{E}^T$  denote expectations over the data distributions of  $(\mathbf{x}, a)$  in the source and target domains, respectively.  $g$  is a discriminator that classifies whether  $\mathbf{x}$  is from the source domain or the target domain.  $\mathcal{L}_D$  is the binary classification accuracy, where labels 0 and 1 indicate the source and target domains, respectively. In practice, we use the cross-entropy loss as a differentiable surrogate loss. As in Ganin et al. (2016b), solving the minimax optimization  $\min_{\hat{\phi}} \max_g \mathcal{L}^{div}$  is equivalent to aligning source- and target-domain data distributions in the latent (encoding) space induced by the encoder  $\hat{\phi}$ .

**Predicted Reward.** The predicted reward term in Eqn. (5), i.e.,  $\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right)$ , can be translated to the loss term  $\sum_{i=1}^N \mathcal{L}_i^{reg}$ , where for each round  $i$  we have

$$\mathcal{L}_i^{reg} = \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^T) \rangle \right| + \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^S) \rangle \right| \right). \quad (9)$$

Minimizing the predicted reward  $\mathcal{L}_i^{reg}$  regularizes the model to avoid overestimating rewards.

**Constant Term.** In Eqn. (5),  $\psi + C$  is the constant term, where  $\psi$  is defined in Definition 3.7, and  $C$  contains some small constants which can be ignored (see the Appendix for more details).

## 4.2 PUTTING EVERYTHING TOGETHER: DABAND TRAINING ALGORITHM

Putting all non-constant loss terms above together, we have the final minimax optimization problem

$$\min_{\hat{\phi}} \max_g \mathcal{L}^{DABand}, \quad (10)$$

where the value function (objective function):

$$\mathcal{L}^{DABand} = \sum_{i=1}^N (\mathcal{L}_i^{Rs} + 2 \cdot \mathcal{L}_i^{es(h)} + \mathcal{L}_i^{reg}) + \lambda \cdot \mathcal{L}_i^{div}. \quad (11)$$

Table 1: Accuracy on the target domain for the DIGIT dataset. Note that in this zero-shot target regret setting, the accuracy  $ACC = 1 - \frac{1}{N} R_{\mathcal{T}}$ , where  $R_{\mathcal{T}}$  is the target regret.

Metrics	Acc Per Class $\uparrow$										Acc $\uparrow$
Test Accuracy	0	1	2	3	4	5	6	7	8	9	Average
LINUCB (Li et al., 2010)	0.3667 $\pm$ 0.03	0.4045 $\pm$ 0.02	0.5102 $\pm$ 0.03	0.3811 $\pm$ 0.01	0.3157 $\pm$ 0.04	0.3445 $\pm$ 0.01	0.4229 $\pm$ 0.03	0.4595 $\pm$ 0.03	0.3489 $\pm$ 0.02	0.2449 $\pm$ 0.01	0.3808 $\pm$ 0.02
LINUCB-P (Li et al., 2010)	0.3645 $\pm$ 0.02	0.6398 $\pm$ 0.04	0.3151 $\pm$ 0.02	0.2709 $\pm$ 0.03	0.2258 $\pm$ 0.01	0.1995 $\pm$ 0.02	0.1752 $\pm$ 0.01	0.3786 $\pm$ 0.04	0.1273 $\pm$ 0.02	0.2992 $\pm$ 0.03	0.3060 $\pm$ 0.02
NLINUCB (Xu et al., 2020)	0.3781 $\pm$ 0.04	0.3228 $\pm$ 0.03	0.3569 $\pm$ 0.03	0.3381 $\pm$ 0.04	0.3663 $\pm$ 0.04	0.4312 $\pm$ 0.05	0.4766 $\pm$ 0.05	0.4004 $\pm$ 0.05	0.4034 $\pm$ 0.02	0.3613 $\pm$ 0.03	0.3816 $\pm$ 0.04
NLINUCB-P (Xu et al., 2020)	<b>0.8588</b> $\pm$ 0.03	0.2224 $\pm$ 0.02	0.3880 $\pm$ 0.03	0.3029 $\pm$ 0.02	0.2978 $\pm$ 0.03	0.0000 $\pm$ 0.01	0.1869 $\pm$ 0.02	0.2243 $\pm$ 0.01	0.2148 $\pm$ 0.01	0.0000 $\pm$ 0.00	0.2706 $\pm$ 0.02
DABAND (OURS)	0.6891 $\pm$ 0.10	<b>0.6662</b> $\pm$ 0.03	<b>0.6188</b> $\pm$ 0.01	<b>0.5703</b> $\pm$ 0.05	<b>0.6008</b> $\pm$ 0.05	<b>0.5534</b> $\pm$ 0.03	<b>0.6010</b> $\pm$ 0.02	<b>0.5709</b> $\pm$ 0.01	<b>0.6266</b> $\pm$ 0.04	<b>0.5023</b> $\pm$ 0.02	<b>0.6002</b> $\pm$ 0.02
OURS vs. NLINUCB	<b>+0.3110</b>	<b>+0.3434</b>	<b>+0.2619</b>	<b>+0.2322</b>	<b>+0.2345</b>	<b>+0.1222</b>	<b>+0.1244</b>	<b>+0.1705</b>	<b>+0.2232</b>	<b>+0.1410</b>	<b>+0.2186</b>

Table 2: Accuracy on the target domain for the VisDA17 dataset. Note that in this zero-shot target regret setting, the accuracy  $ACC = 1 - \frac{1}{N} R_{\mathcal{T}}$ , where  $R_{\mathcal{T}}$  is the target regret.

Metrics	Acc Per Class $\uparrow$												Acc $\uparrow$
Test Accuracy	airplane	bicycle	bus	car	horse	knife	motorcycle	person	plant	skateboard	train	truck	Average
LINUCB (Li et al., 2010)	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
LINUCB-P (Li et al., 2010)	0.0510 $\pm$ 0.02	0.0190 $\pm$ 0.00	0.1360 $\pm$ 0.05	0.0440 $\pm$ 0.02	0.0390 $\pm$ 0.01	<b>0.5580</b> $\pm$ 0.06	0.0410 $\pm$ 0.02	0.2050 $\pm$ 0.05	0.0850 $\pm$ 0.04	0.0920 $\pm$ 0.13	0.0580 $\pm$ 0.02	0.0460 $\pm$ 0.12	0.1145 $\pm$ 0.06
NLINUCB (Xu et al., 2020)	0.2870 $\pm$ 0.03	0.0220 $\pm$ 0.02	0.0000 $\pm$ 0.00	0.1250 $\pm$ 0.02	0.0000 $\pm$ 0.00	0.0130 $\pm$ 0.02	0.0470 $\pm$ 0.02	0.0000 $\pm$ 0.00	0.0460 $\pm$ 0.02	0.0050 $\pm$ 0.00	0.6550 $\pm$ 0.04	0.0000 $\pm$ 0.00	0.1001 $\pm$ 0.02
NLINUCB-P (Xu et al., 2020)	0.0060 $\pm$ 0.01	0.0020 $\pm$ 0.01	<b>0.9222</b> $\pm$ 0.02	0.0010 $\pm$ 0.00	0.0080 $\pm$ 0.01	0.0990 $\pm$ 0.03	0.0030 $\pm$ 0.00	0.0020 $\pm$ 0.00	0.0180 $\pm$ 0.02	0.0420 $\pm$ 0.01	0.0000 $\pm$ 0.00	0.0150 $\pm$ 0.01	0.0932 $\pm$ 0.01
DABAND (OURS)	<b>0.5504</b> $\pm$ 0.02	<b>0.3562</b> $\pm$ 0.04	0.4462 $\pm$ 0.02	<b>0.3844</b> $\pm$ 0.01	<b>0.5582</b> $\pm$ 0.04	0.2632 $\pm$ 0.03	<b>0.6474</b> $\pm$ 0.03	<b>0.2396</b> $\pm$ 0.03	<b>0.4882</b> $\pm$ 0.04	<b>0.4062</b> $\pm$ 0.02	<b>0.7685</b> $\pm$ 0.05	<b>0.1968</b> $\pm$ 0.02	<b>0.4644</b> $\pm$ 0.03
OURS vs. NLINUCB	<b>+0.2634</b>	<b>+0.3342</b>	<b>+0.4462</b>	<b>+0.2594</b>	<b>+0.5582</b>	<b>+0.2502</b>	<b>+0.6004</b>	<b>+0.2396</b>	<b>+0.4422</b>	<b>+0.4012</b>	<b>+0.1135</b>	<b>+0.1968</b>	<b>+0.3643</b>

Note that we need to alternate between updating  $\hat{\theta}$  using LinUCB Li et al. (2010) and updating the encoder  $\hat{\phi}$  with the discriminator  $g$  in Eqn. (10).

Formally, our method, described in Alg. 1, operates as follows: In each iteration, we access the context in the original feature space (i.e.,  $\mathbb{R}^d$ ) for each action from both source and target domains, denoted as  $\{x_{i,a}^S\}_{a \in [K]}$  and  $\{x_{i,a}^T\}_{a \in [K]}$  respectively. Subsequently, we compute the latent representation using the encoder  $\phi(\cdot)$  and select the action  $\hat{a}_i$  for iteration  $i$  based on the LinUCB selection rule. The bandit parameters (i.e.,  $\theta, A, b$ ) are updated in each iteration. The weights in  $\phi$  are updated every episode of length  $H$  (i.e., a batch with  $H$  past iterations) using our objective function in Eqn. (10), following the same rule as in Neural-LinUCB (Xu et al., 2020).

## 5 EXPERIMENTS

In this section, we compare DABand with existing methods on real-world datasets.

**Datasets.** To demonstrate the effectiveness of our DABand, We evaluate our methods in terms of prediction accuracy and **zero-shot** target regret on three datasets, i.e., DIGIT (Ganin et al., 2016a), VisDA17 (Peng et al., 2017), and S2RDA49 (Tang & Jia, 2023). See details for each dataset in Appendix B.

**Baselines.** We compare our DABand with both classic and state-of-the-art contextual bandit algorithms, including **LinUCB** (Li et al., 2010), LinUCB with principle component analysis (PCA) pre-processing (i.e., **LinUCB-P**), Neural-LinUCB (Xu et al., 2020) (**NLinUCB**), and Neural-LinUCB variant that incorporates PCA (i.e., **NLinUCB-P**). Details for each baselines and discussion are in Appendix C. Note that domain adaptation baselines are **not applicable** to our setting because it only works in offline settings and assumes complete observability of labels in the source domain (see Appendix C for details).

**Zero-Shot Target Regret (Accuracy).** We evaluate different methods on the zero-shot target regret setting in Definition 3.5, where all methods can only collect feedback from the source domain, but not the target domain. In this setting, accuracy is equal to 1 minus the average target regret, i.e.,  $ACC = 1 - \frac{1}{N} R_{\mathcal{T}}$ , where  $R_{\mathcal{T}}$  is defined in Eqn. (2) and Definition 3.5. Therefore higher accuracy indicates lower target regret and better performance.

We report the per-class accuracy and average accuracy of different methods in Table 1, Table 2 and Table 3 for DIGIT, VisDA17 and S2RDA49, respectively. Our DABand demonstrates favorable performance against NLinUCB, a leading contextual bandit algorithm with representation learning. Despite NLinUCB’s superior representation learning capabilities through neural networks (NN) over LinUCB – yielding impressive performance in single-domain applications – its efficacy in cross-domain tasks is heavily contingent on the disparity between the source and target domains.



Table 3: Accuracy on the target domain for the S2RDA49 dataset. Note that in this zero-shot target regret setting, the accuracy  $ACC = 1 - \frac{1}{N}R_T$ , where  $R_T$  is the target regret.

Metrics	Acc Per Class $\uparrow$										Acc $\uparrow$
Test Accuracy	acropplane	bicycle	bus	car	knife	motorcycle	plant	skateboard	train	truck	Average
LINUCB (Li et al., 2010)	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
LINUCB-P (Li et al., 2010)	0.1430 $\pm$ 0.02	0.2600 $\pm$ 0.03	0.0810 $\pm$ 0.01	0.0240 $\pm$ 0.01	0.1070 $\pm$ 0.02	0.0570 $\pm$ 0.01	0.1060 $\pm$ 0.01	0.0390 $\pm$ 0.01	0.0960 $\pm$ 0.02	0.1260 $\pm$ 0.03	0.1039 $\pm$ 0.02
NLINUCB (Xu et al., 2020)	0.0350 $\pm$ 0.00	0.1210 $\pm$ 0.02	0.0510 $\pm$ 0.01	0.0870 $\pm$ 0.02	0.0270 $\pm$ 0.00	0.0240 $\pm$ 0.01	0.2720 $\pm$ 0.04	0.0090 $\pm$ 0.00	<b>0.1730</b> $\pm$ 0.02	0.0309 $\pm$ 0.01	0.1108 $\pm$ 0.02
NLINUCB-P (Xu et al., 2020)	0.0000 $\pm$ 0.00	0.0000 $\pm$ 0.00	0.0000 $\pm$ 0.00	<b>0.2990</b> $\pm$ 0.03	<b>0.7000</b> $\pm$ 0.02	0.0000 $\pm$ 0.00	0.0000 $\pm$ 0.00	0.0000 $\pm$ 0.00	0.0000 $\pm$ 0.00	0.1160 $\pm$ 0.01	0.1115 $\pm$ 0.01
DABAND (OURS)	<b>0.5501</b> $\pm$ 0.03	<b>0.6418</b> $\pm$ 0.03	<b>0.5844</b> $\pm$ 0.03	0.2346 $\pm$ 0.01	0.0410 $\pm$ 0.01	<b>0.7580</b> $\pm$ 0.04	<b>0.5689</b> $\pm$ 0.04	<b>0.1345</b> $\pm$ 0.02	0.1118 $\pm$ 0.02	<b>0.3126</b> $\pm$ 0.02	<b>0.3923</b> $\pm$ 0.03
OURS vs. NLINUCB	<b>+0.5151</b>	<b>+0.5208</b>	<b>+0.5334</b>	<b>+0.1476</b>	<b>+0.0140</b>	<b>+0.7340</b>	<b>+0.2969</b>	<b>+0.1255</b>	<b>-0.0612</b>	<b>+0.2817</b>	<b>+0.2815</b>

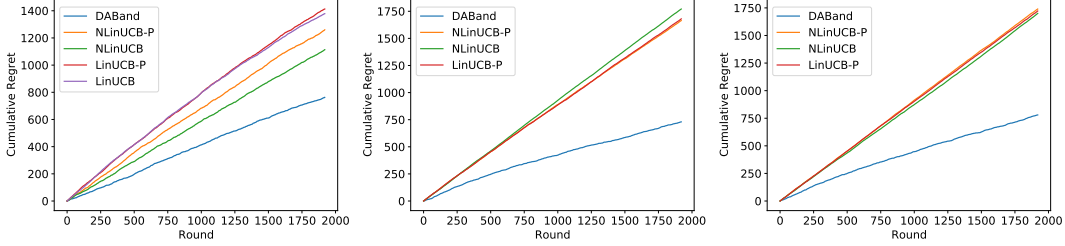


Figure 1: The cumulative regrets of different methods for 1,920 rounds on DIGIT (left), VisDA17 (middle), and S2RDA49 (right). Results are averaged over 5 runs. LinUCB is not reported for VisDA17 and S2RDA49 due to out-of-memory issues.

Typically, the strengths of NN, which include enhanced representation learning (feature extraction), may inadvertently amplify this domain divergence, potentially leading to overfitting on the source domain. Conversely, our DABand algorithm not only improves accuracy but also adeptly addresses the challenges posed by domain divergence. Note that LinUCB encounters an out-of-memory issue due to the high-dimensionality of the data in VisDA17 and S2RDA49, marked as “N/A” in Table 2 and Table 3, respectively.

**Performance Analysis of LinUCB and NLinUCB.** Our target regret bound in Eqn. (5) consists of 5 terms, i.e., source regret, regression error, data divergence, constant, and predicted reward. After applying linear contextual bandits on the source domain, the source regret enjoys a sub-linear rate. However, in the linear model, the encoder  $\hat{\phi}(\mathbf{x}) = \mathbf{x}$  is an identity function and therefore fail to align source-domain and target-domain contexts, leading to an unbounded, large data divergence term in Eqn. (5). This is the main reason why linear contextual bandits, such as LinUCB, perform poorly in the cross-domain contextual bandit setting.

Note that even contextual bandits using a nonlinear encoder, e.g., Neural-LinUCB (NLinUCB), fail in this case because their encoders  $\hat{\phi}(\mathbf{x})$  are not learned to align source-domain and target-domain contexts; they therefore will still lead to an unbounded, large data divergence term in Eqn. (5), and subsequently poor performance in the cross-domain contextual bandit setting.

**Continued Training in Target Domains and Cumulative Regret.** Besides zero-shot target regret, we also evaluate how different methods perform when one continues their training in the high-cost target domain and starts to collect feedback (reward) after the model is trained in the source domain. A good cross-domain bandit model can provide a head-start in this setting, therefore enjoying significantly lower cumulative regret in the target domain and saving substantial cost in collecting feedback in the high-cost target domain.

Table 4: Results of the ablation study in terms of accuracy (higher is better). Note that the accuracy  $ACC = 1 - \frac{1}{N}R_T$ , where  $R_T$  is the target regret. “R” and “P” is short for “Regression Error” and “Predicted Reward”, respectively.

Datasets	w/o R & P	w/o R	w/o P	DABAND (FULL)
DIGIT	0.5676	0.5682	0.5768	<b>0.6002</b>
VisDA17	0.4088	0.4096	0.4304	<b>0.4644</b>
S2RDA49	0.3691	0.3694	0.3719	<b>0.3923</b>

Fig. 1 shows the cumulative regrets of different methods on DIGIT, VisDA17 and S2RDA49, respectively. We show results averaged over five runs with different random seeds. Our DABand consistently

surpasses all baselines for all rounds in terms of both cumulative regrets and their increase rates, demonstrating DABand’s potential to significantly reduce the cost in high-cost target domains.

**Ablation Study.** To verify the effectiveness of each term in our DABand’s regret bound (Theorem 3.1, we report the accuracy of our proposed DABand after removing the Regression Error term and/or the Predicted Reward term (during training) in Table 4 for DIGIT, VisDA17 and S2RDA49. Results on all datasets show that removing either term will lead to performance drop. For example, in VisDA17, our full DABand achieves accuracy of 0.4642; the accuracy drops to 0.4238 after removing the Regression Error (R) term, and further drops to 0.4088 after removing the Predicted Reward (P) term. These results verify the effectiveness of these two terms in our DABand.

## 6 CONCLUSION

In this paper, we introduce a novel domain adaptive neural contextual bandit algorithm, DABand, which adeptly combines effective exploration with representation alignment, utilizing unlabeled data from both source and target domains. Our theoretical analysis demonstrates that DABand can attain a sub-linear regret bound within the target domain. This marks the first regret analysis for domain adaptation in contextual bandit problems incorporating deep representation, shallow exploration, and adversarial alignment. We show that all these elements are instrumental in the domain adaptive bandit setting on real-world datasets. Moving forward, interesting future research could be to uncover more innovative techniques for aligning the source and target domains (mentioned in Appendix F), particularly within the constraints of: (1) bandit settings, as opposed to classification settings, and (2) the high-dimensional and dense nature of the target domain, contrasted with the sparse and simplistic nature of the source domain.

## REFERENCES

- Alekh Agarwal, Daniel J. Hsu, Satyen Kale, John Langford, Lihong Li, and Robert E. Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *ICML*, pp. 1638–1646, 2014.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *JMLR*, 3:397–422, 2002.
- Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems*, 32, 2019.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- Ziliang Chen, Jingyu Zhuang, Xiaodan Liang, and Liang Lin. Blending-target domain adaptation by adversarial meta-adaptation networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2248–2257, 2019.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *AISTATS*, pp. 208–214, 2011.
- Shuyang Dai, Kihyuk Sohn, Yi-Hsuan Tsai, Lawrence Carin, and Manmohan Chandraker. Adaptation across extreme variations using unlabeled domain bridges. *arXiv preprint arXiv:1906.02238*, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- P Kingma Diederik. Adam: A method for stochastic optimization. (*No Title*), 2014.

- Dylan J. Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert E. Schapire. Practical contextual bandits with regression oracles. In *ICML*, pp. 1534–1543, 2018.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016a.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016b.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pp. 2672–2680, 2014.
- Han Guo, Ramakanth Pasunuru, and Mohit Bansal. Multi-source domain adaptation for text classification via distancenet-bandits. In *Proceedings of the AAAI conference on artificial intelligence*, 2020.
- Nathan Korda, Balazs Szorenyi, and Shuai Li. Distributed clustering of linear bandits in peer to peer networks. In *ICML*, pp. 1301–1309, 2016.
- Seiichi Kuroki, Nontawat Charoenphakdee, Han Bao, Junya Honda, Issei Sato, and Masashi Sugiyama. Unsupervised domain adaptation based on source-guided discrepancy. In *AAAI*, pp. 4122–4129, 2019.
- Branislav Kveton, Csaba Szepesvári, Anup Rao, Zheng Wen, Yasin Abbasi-Yadkori, and S. Muthukrishnan. Stochastic low-rank bandits. *CoRR*, abs/1712.04644, 2017.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *WWW*, pp. 661–670, 2010.
- Shuai Li, Alexandros Karatzoglou, and Claudio Gentile. Collaborative filtering bandits. In *SIGIR*, pp. 539–548, 2016.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I. Jordan. Conditional adversarial domain adaptation. In *NIPS*, pp. 1647–1657, 2018.
- Kanak Mahadik, Qingyun Wu, Shuai Li, and Amit Sabne. Fast distributed bandits for online recommendation systems. In *SC*, pp. 1–13, 2020.
- Le Thanh Nguyen-Meidine, Atif Belal, Madhu Kiran, Jose Dolz, Louis-Antoine Blais-Morin, and Eric Granger. Unsupervised multi-target domain adaptation through knowledge distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1339–1347, 2021.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *TKDE*, 22(10):1345–1359, 2009.
- Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *TNN*, 22(2):199–210, 2010.
- Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1406–1415, 2019.
- Viraj Prabhu, Shivam Khare, Deeksha Kartik, and Judy Hoffman. Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8558–8567, 2021.

- Esteban Real, Jonathon Shlens, Stefano Mazzocchi, Xin Pan, and Vincent Vanhoucke. Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5296–5305, 2017.
- Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, pp. 3723–3732, 2018.
- Swami Sankaranarayanan, Yogesh Balaji, Carlos D. Castillo, and Rama Chellappa. Generate to adapt: Aligning domains using generative adversarial networks. In *CVPR*, pp. 8503–8512, 2018.
- Baochen Sun and Kate Saenko. Deep CORAL: correlation alignment for deep domain adaptation. In *ICCV workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, pp. 443–450, 2016.
- Hui Tang and Kui Jia. A new benchmark: On the utility of synthetic data with blender for bare supervised learning and downstream domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15954–15964, 2023.
- Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, pp. 7167–7176, 2017.
- Hao Wang, Hao He, and Dina Katabi. Continuously indexed domain adaptation. *arXiv preprint arXiv:2007.01807*, 2020.
- Hao Wang, Yifei Ma, Hao Ding, and Yuyang Wang. Context uncertainty in contextual bandits with applications to recommender systems. In *AAAI*, 2022.
- Pan Xu, Zheng Wen, Handong Zhao, and Quanquan Gu. Neural contextual bandits with deep representation and shallow exploration. *arXiv preprint arXiv:2012.01780*, 2020.
- Zihao Xu, Guangyuan Hao, Hao He, and Hao Wang. Domain indexing variational bayes: Interpretable domain index for domain adaptation. In *ICLR*, 2023.
- Yisong Yue and Carlos Guestrin. Linear submodular bandits and their application to diversified retrieval. In *NIPS*, pp. 2483–2491, 2011.
- Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael I Jordan. Bridging theory and algorithm for domain adaptation. *arXiv preprint arXiv:1904.05801*, 2019.
- Han Zhao, Shanghang Zhang, Guanhang Wu, José M. F. Moura, João Paulo Costeira, and Geoffrey J. Gordon. Adversarial multiple source domain adaptation. In *NIPS*, pp. 8568–8579, 2018.
- Han Zhao, Remi Tachet des Combes, Kun Zhang, and Geoffrey J. Gordon. On learning invariant representations for domain adaptation. In *ICML*, pp. 7523–7532, 2019.
- Mingmin Zhao, Shichao Yue, Dina Katabi, Tommi S. Jaakkola, and Matt T. Bianchi. Learning sleep stages from radio signals: A conditional adversarial architecture. In *ICML*, pp. 4100–4109, 2017.
- Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with upper confidence bound-based exploration. *arXiv preprint arXiv:1911.04462*, 2019.

## A PROOFS

### A.1 PROOFS FOR LEMMAS

**Lemma A.1.** For any hypotheses  $h, h' \in \mathcal{H}$  and  $a \in \mathcal{A}$  is selected by either  $h$  or  $h'$ ,

$$|\epsilon_{\mathcal{S}}(h, h') - \epsilon_{\mathcal{T}}(h, h')| \leq \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}).$$

*Proof.* By definition of  $\mathcal{H}\Delta\mathcal{H}$  (as defined in Def. 3.6), we have

$$\begin{aligned} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) &\stackrel{\text{Def. 3.6}}{=} 2 \sup_{h, h' \in \mathcal{H}} |\mathbb{E}_{\mathbf{x} \sim \mathcal{S}}[|h(\mathbf{x}) - h'(\mathbf{x})|] - \mathbb{E}_{\mathbf{x} \sim \mathcal{T}}[|h(\mathbf{x}) - h'(\mathbf{x})|]| \\ &= \frac{2}{N} \sup_{h, h' \in \mathcal{H}} \left| N \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{S}}[|h(\mathbf{x}) - h'(\mathbf{x})|] - N \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{T}}[|h(\mathbf{x}) - h'(\mathbf{x})|] \right| \\ &\stackrel{\text{Def. 3.3}}{=} \frac{2}{N} \sup_{h, h' \in \mathcal{H}} |\epsilon_{\mathcal{S}}(h, h') - \epsilon_{\mathcal{T}}(h, h')| \\ &\geq \frac{2}{N} |\epsilon_{\mathcal{S}}(h, h') - \epsilon_{\mathcal{T}}(h, h')|. \end{aligned}$$

□

**Lemma A.2.** Given three hypotheses  $h_1, h_2, h_3 \in \mathcal{H}$ , we have the following triangle inequality:

$$\epsilon_{\mathcal{D}}(h_1, h_2) \leq \epsilon_{\mathcal{D}}(h_1, h_3) + \epsilon_{\mathcal{D}}(h_2, h_3). \quad (12)$$

*Proof.* The hypothesis difference

$$\begin{aligned} \epsilon_{\mathcal{D}}(h_1, h_2) &= \sum_{i=1}^N |h_1(x^{\mathcal{D}}) - h_2(x^{\mathcal{D}})| \\ &\stackrel{\text{Def. 3.2}}{=} \sum_{i=1}^N |h_1(x^{\mathcal{D}}) - h_2(x^{\mathcal{D}}) + h_3(x^{\mathcal{D}}) - h_3(x^{\mathcal{D}})| \\ &\leq \sum_{i=1}^N |h_1(x^{\mathcal{D}}) - h_3(x^{\mathcal{D}})| + \sum_{i=1}^N |h_3(x^{\mathcal{D}}) - h_2(x^{\mathcal{D}})| \\ &= \epsilon_{\mathcal{S}}(h_1, h_3) + \epsilon_{\mathcal{S}}(h_2, h_3). \end{aligned}$$

□

**Lemma A.3.** Let  $\mathcal{H}$  be a hypothesis space. Then for every  $h \in \mathcal{H}$ :

$$\epsilon_{\mathcal{T}}(h) \leq \epsilon_{\mathcal{S}}(h) + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \psi, \quad (13)$$

where  $\psi$  is defined in Definition 3.7.

*Proof.* By definition in Definition 3.3 in main paper, we start from  $\epsilon_{\mathcal{T}}(h) = \epsilon_{\mathcal{T}}(h, f_{\mathcal{T}})$ , then we have

$$\begin{aligned} \epsilon_{\mathcal{T}}(h) &= \epsilon_{\mathcal{T}}(f_{\mathcal{T}}, h) \\ &\stackrel{\text{Lm. A.2}}{\leq} \epsilon_{\mathcal{T}}(f_{\mathcal{T}}, h^*) + \epsilon_{\mathcal{T}}(h, h^*) \\ &\leq \epsilon_{\mathcal{T}}(h^*) + \epsilon_{\mathcal{S}}(h, h^*) + |\epsilon_{\mathcal{T}}(h, h^*) - \epsilon_{\mathcal{S}}(h, h^*)| \\ &\stackrel{\text{Lm. A.1}}{\leq} \epsilon_{\mathcal{T}}(h^*) + \epsilon_{\mathcal{S}}(h, h^*) + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) \\ &\leq \epsilon_{\mathcal{T}}(h^*) + \epsilon_{\mathcal{S}}(h) + \epsilon_{\mathcal{S}}(h^*) + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) \\ &= \epsilon_{\mathcal{S}}(h) + \psi + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}). \end{aligned}$$

□



**Lemma A.4.** Denoting the contexts from the target domain  $\mathcal{T}$  as  $\{\mathbf{x}_{t,a}^{\mathcal{T}}\}_{i \in [N], a \in [K]}$  and the associated ground-truth action as  $\{a_i^*\}_{i=1}^N$ , the estimated regret (i.e., using the trained model  $\hat{\theta}$  from the source domain) for the target domain is

$$\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \leq \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \alpha \cdot \kappa_i^{\mathcal{T}},$$

where  $\kappa_i^{\mathcal{T}}$  is  $\max \left\{ 2 \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} + \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right\}$ .

*Proof.* By definition, we have

$$\begin{aligned} & \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\ &= \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} - \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} - \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right| \right) \\ &= \sum_{i=1}^N \left( \left| \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) - \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) \right| \right) \\ &\quad + \sum_{i=1}^N \alpha \left( \left| \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} - \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right| \right) \\ &\leq \sum_{i=1}^N \left( \left| \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) - \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) \right| \right) \\ &\quad + \sum_{i=1}^N \alpha \max \left\{ \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right\} \end{aligned} \quad (14)$$

$$\leq \sum_{i=1}^N \left( \left| \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) \right| \right) + \sum_{i=1}^N \alpha \max \left\{ \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right\} \quad (15)$$

$$\begin{aligned} &\leq \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \alpha \cdot \max \left\{ 2 \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} + \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right\} \\ &:= \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \alpha \cdot \kappa_i^{\mathcal{T}}, \end{aligned} \quad (16)$$

where the inequality in Eqn. (14) is based on the fact that for any  $a, b > 0$ ,  $|a - b| \leq \max\{a, b\}$ . The inequality in Eqn. (15) is derived from the definition of  $\hat{a}_i$ , since  $\hat{a}_i$  is selected from the maximum value of  $\left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right)$ , therefore we guarantee that  $\left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right) \geq \left( \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle + \alpha \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right)$ . In addition, Eqn. (16) holds by definition as we define  $\kappa_i^{\mathcal{T}}$  as  $\max \left\{ 2 \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} + \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \right\|_{A_{i-1}^{-1}} \right\}$ .  $\square$

Next, we discuss the property of  $\kappa_i^{\mathcal{T}}$  in the following lemma.

**Lemma A.5.** For any arbitrary domain  $\mathcal{D}$ , we denote

$\kappa_i^{\mathcal{D}} = \max \left\{ 2 \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}}, \left\| \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}} + \left\| \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}} \right\}$ . Then, by our DABand algorithm, as the time step  $i \rightarrow \infty$ , we have  $\kappa_i^{\mathcal{D}} \rightarrow 0$ .

*Proof.* By definition of  $A$ , its singular value increase as  $i$  increase, which means that the singular values in its inverse matrix,  $A^{-1}$ , will decrease to 0. Therefore, for any new array  $u$ ,  $\|u\|_{A^{-1}} =$

$\sqrt{uA^{-1}u^T} \rightarrow 0$ . Then when  $i \rightarrow \infty$ , we have

$$\lim_{i \rightarrow \infty} \kappa_i^{\mathcal{D}} = \max \left\{ 2 \left\| \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}}, \left\| \widehat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}} + \left\| \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{D}}) \right\|_{A_{i-1}^{-1}} \right\} \approx 0.$$

□

**Lemma A.6.** Denoting the contexts from the source domain  $\mathcal{S}$  as  $\{\mathbf{x}_{t,a}^{\mathcal{S}}\}_{i \in [N], a \in [K]}$ , the associated ground-truth action as  $\{a_i^*\}_{i=1}^N$ , we have the following inequality:

$$\sum_{i=1}^N \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \leq \epsilon_{\mathcal{S}}(h) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| + \alpha \cdot \kappa_i^{\mathcal{S}} \right),$$

where  $\kappa_i^{\mathcal{S}}$  is defined in Lemma A.5.

*Proof.* By definition, we have

$$\begin{aligned} & \sum_{i=1}^N \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \\ = & \sum_{i=1}^N \left[ \mathbb{1}[a_i^* = \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) + \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \right] \\ \leq & \sum_{i=1}^N \left[ \mathbb{1}[a_i^* = \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - 1 \right| \right) + \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - 0 \right| \right) \right] \\ \leq & \sum_{i=1}^N \left[ \mathbb{1}[a_i^* = \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - 1 \right| \right) + \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle + \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \right] \\ \leq & \sum_{i=1}^N \left[ \mathbb{1}[a_i^* = \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - 1 \right| \right) + \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - 0 \right| \right) + \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \right] \\ \stackrel{\text{L.m. A.4}}{\leq} & \sum_{i=1}^N \left( \left| \langle \theta_{\mathcal{S}}^*, \phi_{\mathcal{S}}^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle - \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| + \alpha \cdot \kappa_i^{\mathcal{S}} \right) \\ \equiv & \epsilon_{\mathcal{S}}(h) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| + \alpha \cdot \kappa_i^{\mathcal{S}} \right), \end{aligned}$$

□

## A.2 PROOF FOR THE TARGET REGRET BOUND

**Theorem A.1 (Target Regret Bound).** Denoting the contexts from the source domain  $\mathcal{S}$  as  $\{\mathbf{x}_{t,a}^{\mathcal{S}}\}_{i \in [N], a \in [K]}$ , the associated ground-truth action as  $\{a_i^*\}_{i=1}^N$ , and the contexts from the target domain  $\mathcal{T}$  as  $\{\mathbf{x}_{t,a}^{\mathcal{T}}\}_{i \in [N], a \in [K]}$ , the upper bound for our target regret  $R_{\mathcal{T}}$  is

$$\begin{aligned} R_{\mathcal{T}} & \triangleq \sum_{i=1}^N \left( \left| \langle \theta_{\mathcal{T}}^*, \phi_{\mathcal{T}}^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{T}}) \rangle - \langle \theta_{\mathcal{T}}^*, \phi_{\mathcal{T}}^*(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\ & \leq \underbrace{R_{\mathcal{S}}}_{\text{Source Regret}} + \underbrace{2 \cdot \epsilon_{\mathcal{S}}(h)}_{\text{Regression Error}} + \underbrace{N \cdot \widehat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T})}_{\text{Data Divergence}} + \underbrace{\psi + C}_{\text{Constant}} \\ & \quad + \underbrace{\sum_{i=1}^N \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \widehat{a}_i] \left( \left| \langle \widehat{\theta}, \widehat{\phi}(\mathbf{x}_{i,\widehat{a}_i}^{\mathcal{S}}) \rangle \right| \right)}_{\text{Predicted Rewards}}, \end{aligned}$$

where  $R_{\mathcal{S}}$ ,  $\epsilon_{\mathcal{S}}(h)$ , and  $\widehat{d}_{\mathcal{H}\Delta\mathcal{H}}$  are the source regret, source-domain error, and  $\mathcal{H}\Delta\mathcal{H}$  divergence defined in Definition 3.4, Definition 3.3, and Definition 3.6, respectively.  $\psi$  is a constant independent to the problem and  $C$  is a constant which can be ignored.

*Proof.* Please check the full proof in the next page.

By the definition of the Target Regret Bound, we have

$$\begin{aligned}
R_{\mathcal{T}} &= \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\
&= \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\
&\leq \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\
&\leq \epsilon_{\mathcal{T}}(h) + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle + \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle \right| \right) \\
&\leq \epsilon_{\mathcal{T}}(h) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle \right| \right) \\
&\leq \epsilon_{\mathcal{T}}(h) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle \right| \right) \\
&\quad + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle \right| \right) - \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle \right| \right) \\
&\leq \epsilon_{\mathcal{T}}(h) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle \right| \right) \\
&\quad + \left| \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle \right| \right) - \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle \right| \right) \right| \\
&\leq \epsilon_{\mathcal{T}}(h) + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\
&\quad + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle + \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \\
&\stackrel{\text{L.m. A.3}}{\leq} \epsilon_{\mathcal{S}}(h) + \psi + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \frac{N}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{T}}) \rangle - \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) \\
&\quad + \sum_{i=1}^N \left( \left| \langle \theta^*, \phi^*(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \\
&\stackrel{\text{L.m. A.4}}{\leq} \epsilon_{\mathcal{S}}(h) + \psi + N \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \alpha \cdot \kappa_i^{\mathcal{T}} \\
&\quad + R_{\mathcal{S}} + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,a_i^*}^{\mathcal{S}}) \rangle - \langle \theta^*, \phi^*(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \\
&\stackrel{\text{L.m. A.6}}{\leq} \epsilon_{\mathcal{S}}(h) + \psi + N \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T}) + \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \alpha \cdot \kappa_i^{\mathcal{T}} \\
&\quad + R_{\mathcal{S}} + \epsilon_{\mathcal{S}}(h) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| + \alpha \cdot \kappa_i^{\mathcal{S}} \right) \\
&\stackrel{\text{L.m. A.5}}{=} \underbrace{R_{\mathcal{S}}}_{\text{Source Regret}} + \underbrace{2 \cdot \epsilon_{\mathcal{S}}(h)}_{\text{Regression Error}} + \underbrace{N \cdot \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T})}_{\text{Data Divergence}} + \underbrace{\psi + C}_{\text{Constant}} \\
&\quad + \underbrace{\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right) \right)}_{\text{Predicted Rewards}},
\end{aligned}$$

where  $C = \sum_{i=1}^N \alpha \cdot \kappa_i^T + \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \sum_{i=1}^N \alpha \cdot \kappa_i^S \right)$  is a small number which can be ignored as  $i \rightarrow \infty$  (i.e., see Lemma A.5).  $\square$

## B DATASETS

To evaluate the effectiveness of our DABand, for each dataset, we treat the “easy” domain as the low-cost source domain and the more challenging one as the high-cost target domain. This allows us to demonstrate the efficacy of our method by adapting from a simpler to a more complex domain within a contextual bandit setting.

*DIGIT.* Our DIGIT dataset consists of MNIST and MNIST-M. MNIST is a gray-scale hand-written digit dataset, while MNIST-M (Ganin et al., 2016a) features color digits. In this paper, MNIST serves as our source domain, and MNIST-M as our target domain, offering a more challenging adaptation path. The MNIST digits are gray-scale and uniform in size, aspect ratio, and intensity range, in stark contrast to the colorful and varied digits of MNIST-M. Therefore, adapting from MNIST to MNIST-M presents a greater challenge than adapting from MNIST-M to MNIST, where the target domain is less complex.

*VisDA17.* The VisDA-2017 image classification challenge (Peng et al., 2017) addresses a domain adaptation problem across 12 classes, involving three distinct datasets. The training set consists of 3D rendering images, whereas the validation and test sets feature real images from the COCO (Lin et al., 2014) and YouTube Bounding Boxes (Real et al., 2017) datasets, respectively. Ground truth labels were provided only for the training and validation sets. Scores for the test set were calculated by a server operated by the competition organizers. For our purposes, we use the training set as the source domain and the validation set as the target domain.

*S2RDA49.* The S2RDA49 (Synthetic-to-Real) is a new benchmark dataset (Tang & Jia, 2023) constructed in 2023. This dataset contains 49 classes. The source domain (i.e., the synthetic domain) is synthesized by rendering 3D models from ShapeNet (Chang et al., 2015). The used 3D models are in the same label space as the target/real domain, and each class has 12K rendered RGB images. The target domain (i.e., the real domain) of S2RDA49 contains 60535 images from 49 classes, collected from the ImageNet validation set (Deng et al., 2009), ObjectNet (Barbu et al., 2019), VisDA2017 validation set (Peng et al., 2017), and the web. In this paper, we select the only 10 class that matches the VisDA17 dataset.

## C BASELINES

We compare our DABand with both classic and state-of-the-art contextual bandit algorithms. We start with **LinUCB** (Li et al., 2010), a typical baseline for linear contextual bandit problems. To enable comprehensive comparisons, we extend our evaluation to include LinUCB augmented with pre-processed features using principle component analysis (PCA), referred to as **LinUCB-P**. In LinUCB-P, we first apply PCA to fit on training data (without labels/feedback) in both source and target domain. Then, with transformed, lower-dimensional context data, LinUCB is performed to see whether this pre-alignment procedure can transfer the knowledge to the target domain. Another pivotal baseline in our study is Neural-LinUCB (Xu et al., 2020) (**NLinUCB**), which utilizes several layers of fully connected neural networks to dynamically process the original features at the beginning of every iteration. Similar to LinUCB-P, we introduce a **NLinUCB** variant that incorporates PCA, i.e., **NLinUCB-P**. Note that domain adaptation baselines are **not applicable** to our setting. Specifically, domain adaptation methods only work in offline settings, and assume complete observability of labels in the source domain. In contrast, contextual bandit is an online setting where the oracle is revealed **only when correctly predicted**. Therefore domain adaptation methods are not applicable to our online bandit settings.

## D IMPLEMENTATION DETAILS

In this section, we provide detailed insights into the implementation of our approach, applied to two distinct datasets: DIGIT and VisDA17. We use Pytorch to implement our method, and all experiments are run on servers with NVIDIA A5000 GPUs.

**DIGIT.** Within the DIGIT dataset framework, we use the MNIST dataset as the source domain and MNIST-M as the target domain. To ensure compatibility between the datasets, we standardize the channel size of images in the source domain ( $c_S = 1$ ) to align with that in the target domain ( $c_T = 3$ ). Each image undergoes normalization and is resized to  $28 \times 28$  pixels with 3 channels to accommodate the format requirements of both domains. Then, an encoder is utilized to diminish the data’s dimensionality to a more manageable latent space. Following this reduction, the data is processed through two fully connected neural network layers, ending in the final latent space necessary for loss computation as delineated in main paper. For the optimal hyperparameters, we set the learning rate to  $1 \times e^{-5}$ , with  $\lambda$  is chosen from  $\{1.0, 5.0, 10.0, 15.0, 20.0\}$  and kept the same for all experiments. Additionally, we set the exploration rate  $\alpha$  to 0.05.

**VisDA17.** We use the VisDA17 dataset’s training set as the source domain, with the validation set functioning as the target domain. We adhere to preprocessing steps established by (Prabhu et al., 2021) to ensure uniformity across domains: Each image is normalized and resized to  $224 \times 224$  pixels with 3 channels, matching the requisite specifications for both source and target domains. For hyperparameters, the learning rate of  $1 \times e^{-5}$  is applied, with  $\lambda$  is chosen from  $\{1.0, 5.0, 10.0, 15.0, 20.0\}$  and then kept the same for all experiments. We set the exploration rate  $\alpha$  to 0.05.

**S2RDA49.** We selects 10 classes from the original 49 classes since it matches the target domain samples in VisDA17 (Peng et al., 2017). We adhere to preprocessing steps established by (Prabhu et al., 2021) to ensure uniformity across domains: Each image is normalized and resized to  $224 \times 224$  pixels with 3 channels. For hyperparameters, the learning rate of  $1 \times e^{-3}$  is applied, with  $\lambda$  chosen from  $\{1.0, 5.0, 10.0, 15.0, 20.0\}$  and then kept the same for all experiments. We set the exploration rate  $\alpha$  to 0.01.

Table 5: Results of the ablation studies in terms of accuracy (higher is better). Note that the accuracy  $ACC = 1 - \frac{1}{N} R_T$ , where  $R_T$  is the target regret. “R”, “P” and “D” are short for “Regression Error”, “Predicted Reward” and “Data Divergence” (i.e., adversarial loss and the discriminator), respectively.

Datasets	w/o R&P&D	w/o R&P	w/o R&D	w/o P&D	w/o R	w/o P	w/o D	DABAND (FULL)
DIGIT	0.3816 $\pm$ 0.04	0.5676 $\pm$ 0.02	0.3793 $\pm$ 0.02	0.3544 $\pm$ 0.01	0.5682 $\pm$ 0.01	0.5768 $\pm$ 0.01	0.3649 $\pm$ 0.02	<b>0.6002</b> $\pm$ 0.02
VisDA17	0.1001 $\pm$ 0.02	0.4088 $\pm$ 0.03	0.1010 $\pm$ 0.02	0.0936 $\pm$ 0.01	0.4096 $\pm$ 0.01	0.4304 $\pm$ 0.01	0.1098 $\pm$ 0.01	<b>0.4644</b> $\pm$ 0.03
S2RDA49	0.1108 $\pm$ 0.02	0.3691 $\pm$ 0.02	0.0918 $\pm$ 0.01	0.1032 $\pm$ 0.01	0.3694 $\pm$ 0.03	0.3719 $\pm$ 0.02	0.1121 $\pm$ 0.02	<b>0.3923</b> $\pm$ 0.03

## E TOTAL ABLATION STUDIES

The full ablation study is shown in Table 5. Furthermore, we ran the corresponding hypothesis tests, and the p values are in the range of  $(3.201 \times 10^{-21}, 1.504 \times 10^{-2})$ , much lower than the threshold of 0.05 and therefore verifying the significance of DABand’s performance improvement.

## F LIMITATION

This work contains several limitations. Specifically, we highlight below:

**Intuitive Perspective.** The bandit algorithm indeed enjoys interpretability and the accuracy of its closed-form updates, but this is true only because it typically employs a linear model.

Unfortunately, linear models do not work very well for real-world data, which is often high-dimensional. For example, the empirical results for LinUCB and LinUCB-P in Table 1 and Table 2 show that such linear contextual bandit algorithms significantly underperform state-of-the-art neural bandit algorithms, which use back-propagation (similar results are shown in (Zhou et al., 2019; Xu et al., 2020)). In summary deep learning models with back-propagation is necessary due to the following reasons:

- **High-Dimensional Data.** To enhance performance in real-world high-dimensional data (e.g., images), the integration of deep learning (and back-propagation) is necessary, though this may sacrifice a portion of the algorithm’s explanatory power.



- **Covariate Shift and Aligning Source and Target Domains in the Latent Space.** In our settings where there is covariance shift between the source and target domains, a deep (nonlinear) encoder is required to transform the original context into a latent space where source-domain encodings and target-domain encoders can align. This also necessitates deep learning models with back-propagation.

Indeed, there is a trade-off between interpretability and performance. This would certainly be an interesting future direction, but it is out of the scope of this paper.

**Theoretical Perspective.** Our Theorem 3.1 is sharp, as all the inequalities are based on lemmas in the paper and the Cauchy inequality. Identifying the criteria under which the target regret bound reaches equality as well as how one can achieve it in practice would be interesting future work.

**Empirical Perspective.** The performance of DABand largely relies on the alignment quality between the source and target domains. Therefore, for two domains that cannot be aligned (for example, most domain adaptation tasks are predefined, and the data for both domains is pre-processed and cleaned, not original real-world data), it is challenging to evaluate how DABand can still transfer knowledge across different domains. Furthermore, it is still unknown whether, if we increase the domain shift in a dataset from another domain, our DABand can still perform well. These issues are beyond the scope of this paper, but they represent interesting areas for future work.

## G DISCUSSIONS

### G.1 SUBLINEARITY FOR THE TARGET REGRET BOUND

To see the data divergence term is sub-linear, note that our target regret bound in Eqn. (5) can be divided into two parts:

- **Total Source Regret:** In Neural-LinUCB (NLinUCB) (Xu et al., 2020) [1], this is shown to be  $O(\sqrt{N})$ , which is sub-linear.
- **Sum of Other Terms (Including the Data Divergence):** These terms can be \*directly optimized\* to convergence and minimized to a small value using SGD variants such as Adam. As shown in Corollary 4.2 of the Adam paper [2], the Adam optimizer enjoys an average regret  $R(N)/N$  of  $O(\frac{1}{\sqrt{N}})$  (here  $N$  is equivalent to  $T$  in the Adam paper), which leads to a sub-linear total regret  $O(\sqrt{N})$ . Since all other terms are directly optimized using Adam, they also enjoy a sub-linear regret.

Therefore, our target regret bound is also sub-linear. Moreover:

- **Key Difference between the Source Regret and Other Terms.** Note that in our target regret bound in Eqn. (5), the source regret term  $R_S$  increases monotonically with respect to  $N$  and **cannot** be directly minimized. In contrast, all other terms, e.g., the data divergence term  $N \cdot \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T})$  **can** be directly minimized since all related data is training data and is already known. This is why we minimize the corresponding loss terms like Eqn. (7), Eqn. (8) and Eqn. (9) during training.
- The data divergence term can be minimized to a small value as  $N \rightarrow \infty$ . For example, if the source and target domains are perfectly aligned after our training, the data divergence becomes 0. Therefore, the cumulative sum of the data divergence term over all rounds will be sub-linear.

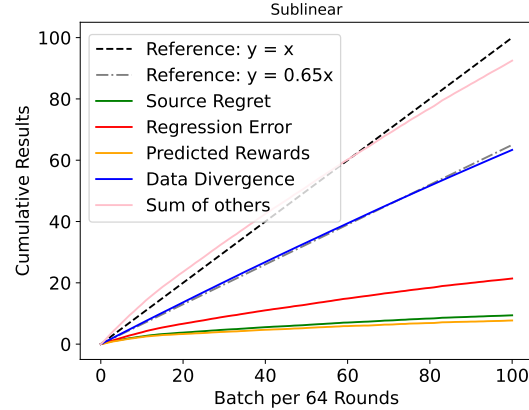


Figure 2: Sub-linearity for each term in Eqn. (5).

## G.2 DIFFERENCE BETWEEN BANDIT AND CLASSIFICATION SETTINGS

In the classification setting, given a sample  $\mathbf{x}$ , a model predicts its label  $\hat{y}$ . Since its ground-truth label is known, one can directly apply the cross-entropy loss to perform back-propagation and update our model. In contrast, in bandit settings, for each sample  $\mathbf{x}$ , our model predicts an action  $\hat{a}$  by optimizing the estimated rewards. One then submits this predicted action to the environment and receives feedback that only indicates whether the predicted action  $\hat{a}$  is the optimal action  $a^*$  or not. If not, we are informed that our prediction is incorrect, yet we do not receive information on what the correct (optimal) action should be. Moreover, during training, we only train on 300 episodes for DIGIT and VisDA17 and 100 episodes for S2RDA49. Each episode contains 64 samples. This implies that compared with those DA methods which train on multiple epochs, in our settings, all of the models only see a sample **once**. Such complexity makes the bandit setting much more challenging.

**Why Simultaneously Training Source and Target Domains is Needed and Helpful.** It is also noteworthy that our DABand is an **end-to-end** model that enables **two-way feedback between source and target domains**:

- **Source to Target:** Collecting **source**-domain rewards helps DABand to learn a better encoder (which transforms raw contexts into encodings in the latent space) because the reward signals help the encoder extract **more relevant encodings** (embeddings) from **target-domain** contexts, thereby **reducing the target-domain regret**.
- **Target to Source:** After aligning source-domain and **target-domain** encodings in the shared latent space (i.e., minimizing the data divergence term in Eqn. (5)), the information from the **target-domain** context can then help the **source** domain to explore arms that are **more relevant to the target domain**. Ultimately, this also helps reduce the target-domain regret even without collecting target-domain rewards.

## G.3 IMPORTANCE OF BANDITS AND DABAND

Bandit algorithms are designed to navigate environments where feedback is sparse, costly, and indirect. By efficiently learning from limited feedback – identifying not just when a prediction is wrong but adapting without explicit guidance on the right choice – bandit algorithms offer a strategic advantage in dynamically evolving settings. Our DABand exemplifies this by leveraging a low-cost source domain to improve performance in a high-cost target domain, thereby reducing cumulative regret (improving accuracy) while minimizing operational costs. DABand not only reduces the expense associated with acquiring and labeling vast datasets but also capitalizes on the intrinsic adaptability of bandit algorithms to learn and optimize in complex, uncertain environments.

## G.4 NOVELTIES RESTATEMENT

**Theoretical Novelty.** In general, DABand is the **first** work to perform contextual bandit in a domain adaptation setting, i.e., adapt from a source domain with feedback to a target domain without feedback. Moreover, our theoretical analysis presents two major technical challenges/novelty below:

- **Generalization of  $\mathcal{H}\Delta\mathcal{H}$  Distance to Regression.** In the original multi-domain generalization bound (Ben-David et al., 2010), the  $\mathcal{H}\Delta\mathcal{H}$  divergence between source and target domains is derived for classification models. However, contextual bandit is essentially a reward regression problem, and therefore necessitates generalizing the  $\mathcal{H}\Delta\mathcal{H}$  distance from the classification case to the regression case. This presents a significant technical challenge, and leads to a series of modifications in the proof.
- **Decomposing the Target Regret into the Source Regret with Other Terms.** The subsequent major challenge our DABand addresses involves decomposing the regret bound for the target domain (i.e., the target regret) into a source regret term and other terms to upper-bound the target regret. This is necessary because we do not have access to feedback/reward from the target domain (we only have access to feedback/reward in the source domain). This process requires a nuanced understanding of the interplay between source and target domain dynamics within our model’s framework.

**Other Technical Novelties.** Moreover, our contribution goes beyond the theoretical analysis itself. Specifically,

- We identify the problem of contextual bandits across domains and propose domain-adaptive contextual bandits (DABand) as the first general method to explore a high-cost target domain while only collecting feedback from a low-cost source domain.
- Our theoretical analysis shows that our method can achieve a sub-linear regret bound in the target domain.
- Our empirical results on real-world datasets show our DABand significantly improves performance over the state-of-the-art contextual bandit methods when adapting across domains.

## G.5 DISCUSSION ON ZERO-SHOT TARGET REGRET BOUND

### G.5.1 SIGNIFICANCE OF THEOREM 3.1

Note that Theorem 3.1 is **nontrivial**. While it does resemble the generalization bound in domain adaptation, there are key differences. As mentioned in Observation (3) in Sec. 3.6, our target regret bound includes two additional crucial terms not found in domain adaptation. Specifically:

- **Regression Error in the Source Domain.**  $\sum_{i=1}^N \left( \left| \langle \theta_{\mathcal{S}}^*, \phi_{\mathcal{S}}^*(x_{i,\hat{a}_i}^{\mathcal{S}}) \rangle - \langle \hat{\theta}, \hat{\phi}(x_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right)$ , which defines the difference between the true reward from selecting action  $\hat{a}_i$  and the estimated reward for this action.
- **Predicted Reward.**  $\sum_{i=1}^N \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{T}}) \rangle \right| \right) + \sum_{i=1}^N \mathbb{1}[a_i^* \neq \hat{a}_i] \left( \left| \langle \hat{\theta}, \hat{\phi}(\mathbf{x}_{i,\hat{a}_i}^{\mathcal{S}}) \rangle \right| \right)$ , which serves as a regularization term to regularize the model to avoid overestimating rewards.

The results of the ablation study in Table 4 in Sec. 5 highlight the significance of these two terms. Please also refer to “Technical Novelty” of Sec. G.4 above for discussion on key novelty/challenges in deriving Theorem 3.1.

### G.5.2 TESTING PHASE RATHER THAN TRADITIONAL BANDIT LEARNING SETTINGS

Our regret bound in Theorem 3.1 is a zero-shot regret bound for the target domain (with empirical results in the **Zero-Shot Target Regret** paragraph of Sec. 5), which corresponds to a testing phase.

However, we would like to clarify that extending this bound to handle continued training in the target domain (corresponding to the **Continued Training in Target Domains and Cumulative Regret** paragraph of Sec. 5) is straightforward. At a high level, we can have

$$R_{\mathcal{T}}^{\text{total}} = R_{\mathcal{T}}^{\text{zero-shot}} + R_{\mathcal{T}}^{\text{continued}},$$

with the first term handled by our DABand’s Theorem 3.1 and the second term handled by typical single-domain contextual bandit.

### G.5.3 SCALE OF SOURCE REGRET AND PREDICTED REWARDS

Similar to Neural-LinUCB (NLinUCB), in our DABand, the true reward is restricted to the range of  $[0, 1]$  (as we mentioned in Sec. 3.1); therefore it will not make the source regret unbounded.

Furthermore, our DABand algorithm tries to minimize the source regret while aligning the source and target domains. The minimization of the source regret is theoretically guaranteed, as discussed in recent work such as Neural-LinUCB (NLinUCB) (Xu et al., 2020).

## G.6 ARE THE COMPARISONS FAIR?

**Leveraging Target-Domain Contexts.** A lot of our baselines **do leverage** the contexts of the target domain. For example, our baseline LinUCB-P starts by performing PCA jointly on both source-domain and target-domain contexts and then perform LinUCB. Therefore LinUCB-P does leverage the contexts of the target domain. Similarly, NLinUCB-P starts by performing PCA jointly on both

source-domain and target-domain context and then perform NLinUCB. It therefore also leverages target-domain contexts.

**Seeing Target-Domain Rewards.** Note that in our domain adaptive bandit settings:

- During the **zero-shot** phase, target rewards are **not** visible for all methods (include both baselines and our DABand); it is therefore fair comparison.
- During the **continued-training** phase, all methods (include both baselines and our DABand) start to see the target reward and update their parameters; it is therefore also fair comparison.

#### G.7 WHY MINIMIZE THE PREDICTED REWARD ON THE SOURCE DOMAIN?

While the predicted-reward term is naturally derived from our theoretical analysis, we do find interesting insights when examining this term and its relation to our model. Specifically:

- **Indirect Regularization on Bandit Parameters  $\hat{\theta}$  and the Encoder  $\hat{\phi}(\cdot)$ .** One can see this term as an L1 regularization term. It does not directly regularize the bandit parameter  $\hat{\theta}$ ; however, minimizing the L1 norm of the predicted reward (i.e., a  $K$ -dimensional vector for a  $K$ -arm bandit) does **indirectly** regularize the bandit parameter  $\hat{\theta}$  and the encoder  $\hat{\phi}(\cdot)$ , thereby preventing the L1 norm of the predicted reward from getting to large.
- **Smaller Predicted Rewards for Smaller Variance and Better Stability.** Furthermore, this regularization can help avoid predicting high rewards for all arms in the bandit.
  - Note that for the bandit algorithm to achieve low regret, predicting large rewards are not necessary. This is because one uses the **argmax** operation (i.e., Line 6 in Alg. 1) to select the best arm; an arm  $k$  with a small predicted reward can still be selected as long as all other arms have even smaller predicted rewards.
  - Too higher predicted rewards are not desirable because they increase the model’s sensitivity, leading to higher variance and subsequently increasing the generalization error (i.e., the target regret’s bound).

#### G.8 CLARIFICATION OF DABAND’S CONTRIBUTIONS

Our DABand is the first general method to explore a target domain while only collecting feedback from the source domain, regardless of linear or nonlinear assumptions. No prior methods have explored this setting of exploring a target domain while only collecting feedback from the source domain.

Therefore, DABand’s contribution is two-fold: (1) DABand is the first general method to explore a target domain while only collecting feedback from the source domain; (2) DABand is also the first general method in this domain-adaptive bandit setting that works even under general nonlinear assumptions.

#### G.9 IMPORTANCE OF THE DISCRIMINATOR.

Our bandit setting aims to explore a target domain while only collecting feedback from the source domain. **All reward feedback in the target domain is unknown**, making this approach very **challenging**. Therefore, using a discriminator to align representations for both domains is necessary. If we ignore the discriminator, the method will not work in our settings. This is also evidenced by the whole results for ablation studies in Table 5.

#### G.10 HOW THE TERMS IN THE REGRET DECAY AS $N$ INCREASES.

To see how the terms in the regret decay:

- **Decay in the Source Regret.** Most of the decay occurs in the first term (i.e., Source Regret  $R_S$ ), which has been discussed in the NLinUCB paper (Xu et al., 2020). Specifically the **cumulative** source regret is  $O(\sqrt{N})$ . Then, dividing both sides of Eqn. (5) by  $N$ , we will get the **average** source regret term with  $O(\sqrt{N}/N)$ , which does decay with  $N$ .

- **Decay in the Data Divergence Term.** Decay in the Data Divergence Term. The data divergence term,  $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{S}, \mathcal{T})$ , can actually be further decomposed into
  - an empirical term that estimates the divergence using  $N$  source-domain contexts and  $N$  target-domain contexts, and
  - a term related to  $N$  with complexity  $\mathcal{O}(\log(2N)/N)$  (Ben-David et al., 2010).

Therefore, the second term with  $\mathcal{O}(\log(2N)/N)$  does also decay as  $N$  increases.

#### G.11 UNDERSTANDING THE BOUND IF THE SOURCE AND TARGET DOMAINS ARE EQUIVALENT

Even if the input distributions of the two domains match, it does not imply that the relationship between the input context  $\mathbf{x}$  and the predicted reward is the same for them, since the data divergence term is only responsible for aligning the distribution of  $\mathbf{x}$  (input); it is irrelevant to the output and rewards. This is why other terms are needed to characterize the difference in the relationship between the context  $\mathbf{x}$  and the reward in the source and target domains. The data divergence term alone is not sufficient.

#### G.12 CLARIFICATION ON THE PROBLEM SETTING.

**Simultaneous Observation.** In our setting, source-domain contexts and target-domain contexts are simultaneously observed, with only the reward of the source domain being observed.

**Target-Domain Contexts Available Even Before Running Alg. 1.** Note that in practice, target-domain contexts are usually **available** even before Algorithm 1 starts and **before observing any reward from the source domain**. (This makes it possible to train our DABand by simultaneously using source-domain and target-domain contexts.)

For example, in the case of testing drug reactions on mice (source domain) and humans (target domain), usually one already has the human subjects’ genomics, demographic, and other data as target-domain contexts, even before testing the drug on mice (i.e., the source domain) and collecting rewards.

This is because these human genomics/demographic data are easy to collect at a low cost; in contrast, testing the drug on humans (i.e., the target domain) and collecting rewards involves extremely high costs, due to the risks of fatality, side effects, and the enormous costs of conducting clinical trials.

## H POTENTIAL IMPACT OF DABAND

Our DABand has many potential real-world applications, especially when obtaining responses is costly. For instance, in testing new drugs, we can construct responses from mice for new drug A, and then use DABand to obtain the zero-shot hypothetical regret bound for responses in humans. In another scenario, we can collect data (responses) on humans for another published, similar drug B, and then transfer knowledge by aligning the divergence between drug A and drug B. If both regrets reveal acceptable performance, we might not need excessive costs for back-and-forth testing, which significantly speeds up the process and reduces costs.