# Estimation Considerations in Contextual Bandits

Maria Dimakopoulou[*]    Zhengyuan Zhou [†]    Susan Athey [‡]    Guido Imbens [§]

**Abstract**

Contextual bandit algorithms are sensitive to the estimation method of the outcome model as well as the exploration method used, particularly in the presence of rich heterogeneity or complex outcome models, which can lead to difficult estimation problems along the path of learning. We study a consideration for the exploration vs. exploitation framework that does not arise in multi-armed bandits but is crucial in contextual bandits; the way exploration and exploitation is conducted in the present affects the bias and variance in the potential outcome model estimation in subsequent stages of learning. We develop parametric and non-parametric contextual bandits that integrate balancing methods from the causal inference literature in their estimation to make it less prone to problems of estimation bias. We provide the first regret bound analyses for contextual bandits with balancing in the domain of linear contextual bandits that match the state of the art regret bounds. We demonstrate the strong practical advantage of balanced contextual bandits on a large number of supervised learning datasets and on a synthetic example that simulates model mis-specification and prejudice in the initial training data. Additionally, we develop contextual bandits with simpler assignment policies by leveraging sparse model estimation methods from the econometrics literature and demonstrate empirically that in the early stages they can improve the rate of learning and decrease regret.

---

[*]Stanford University, Management Science and Engineering, madima@stanford.edu
[†]Stanford University, Electrical Engineering, zyzhou@stanford.edu
[‡]Stanford University, Graduate School of Business, and NBER, athey@stanford.edu
[§]Stanford University, Graduate School of Business, and NBER, imbens@stanford.edu

# 1 Introduction

Contextual bandits seek to learn a personalized treatment assignment policy in the presence of treatment effects that vary with observed contextual features. In such settings, there is a need to balance the exploration of treatments[1] for which there is limited knowledge in order to improve performance in the future against the exploitation of existing knowledge in order to attain better performance in the present (see [20] for a survey). Since large amounts of data can be required to learn how the benefits of alternative treatments vary with individual characteristics, contextual bandits can play an important role in making experimentation and learning more efficient. Several successful contextual bandit designs have been proposed [11, 45, 6, 4, 17]. The existing literature has provided regret bounds (e.g., the general bounds of [51], the bounds of [50, 49, 56] in the case of non-parametric function of arm rewards), has demonstrated successful applications (e.g., news article recommendations [45] or mobile health [42]), and has proposed system designs to apply these algorithms in practice [3].

Contextual bandits are poised to play an important role in a wide range of applications: content recommendation in web-services, where the learner wants to personalize recommendations (arm) to the profile of a user (context) to maximize engagement (reward); online education platforms, where the learner wants to select a teaching method (arm) based on the characteristics of a student (context) in order to maximize the student's scores (reward); and survey experiments, where the learner wants to learn what information or persuasion (arm) influences the responses (reward) of subjects as a function of their demographics, political beliefs, or other characteristics (context).

In the contextual setting, one does not expect to see many future observations with the same context as the current observation, and so the value of learning from pulling an arm for this context accrues when that observation is used to estimate the outcome from this arm for a different context in the future. Therefore, the performance of contextual bandit algorithms can be sensitive to the estimation method of the outcome model or the exploration method used. In the initial phases of learning when samples are small, biases are likely to arise in estimating the outcome model using data from previous non-uniform assignments of contexts to arms. The bias issue is aggravated in the case of a mismatch between the generative model and the functional form used for estimation of the outcome model, or similarly, when the heterogeneity in treatment effects is too complex to estimate well with small datasets. In that case methods that proceed under the assumption that the functional form for the outcome model is correct may be overly optimistic about the extent of the learning so far, and

---

[1]A treatment is also referred to as an arm in the literature. In this paper, we use the two terms interchangeably.

emphasize exploitation over exploration. Another case where biases can arise occurs when training observations from certain regions of the context space are scarce (e.g., prejudice in training data if a non-representative set of users arrives in initial batches of data). These problems are common in real-world settings, such as in survey experiments in the domain of social sciences or in applications to health, recommender systems, or education. For example, early adopters of an online course may have different characteristics than later adopters.

We develop parametric and non-parametric contextual bandits that integrate balancing methods from the causal inference literature [36] in their estimation to make it less prone to the aforementioned sources of bias. Our methods aim to balance covariates between treatment groups and achieve contextual bandit designs which are less prone to problems of bias. The balancing will lead to lower estimated precision in the reward functions, and thus will emphasize exploration longer than the conventional linear TS and UCB algorithms, leading to more robust estimates. Balancing can take various forms, ranging from the well-known inverse propensity score weighting to state of the art methods such as approximate residual balancing [8] or the method of [38]. Moreover, balancing can be integrated in contextual bandits with a parametric model estimation such as ridge or LASSO) or non-parametric model estimation such as forests.

We further investigate the effect of balancing via inverse propensity weighting in the domain of linear contextual bandits, by comparing linear Thompson sampling (LinTS) [6] and linear upper confidence bound (LinUCB) [45] – which have strong theoretical guarantees – with our algorithms, *balanced linear Thompson sampling* (BLTS) and *balanced linear upper confidence bound* (BLUCB). Our main contribution here is to provide extensive and convincing empirical evidence for the effectiveness of BLTS and BLUCB (in comparison to LinTS and LinUCB) by considering the problem of multiclass classification with bandit feedback. Specifically, we transform a $K$-class classification task into a $K$-armed contextual bandit [28] and we use 300 public benchmark datasets for our evaluation. Additionally, we provide regret bounds for BLTS and BLUCB, which match the existing state-of-the-art regret bounds for LinTS and LinUCB. It is important to point out that, even though BLTS and LinTS share the same theoretical guarantee, BLTS outperforms LinTS empirically. Similarly, BLUCB has a strong empirical advantage over LinUCB. In bandits, this phenomenon is not uncommon. For instance, it is well-known that even though the existing UCB bounds are often tighter than those of Thompson sampling, Thompson sampling performs better in practice than UCB [22]. Consequently, we take the view that even though regret is a useful theoretical performance metric, it may not always provide clear guidance on which algorithm should be used in practice. We find that this is also the case for balanced linear contextual bandits, as in our evaluation BLTS has a strong empirical advantage over BLUCB. Overall, in this

large-scale evaluation, BLTS outperforms LinUCB, BLUCB and LinTS. In our empirical evaluation, we also consider a synthetic example which simulates in the simplest possible way two issues of bias that often arise in practice; training data with non-representative contexts and model mis-specification. BLTS is much more effective in escaping these biases and, as in the evaluation on supervised learning datasets, it outperforms LinUCB, BLUCB and LinTS by a large margin.

To our knowledge, this is the first work to integrate balancing in the online contextual bandit setting, to perform a large-scale evaluation of it against direct estimation method baselines with theoretical guarantees and to provide a theoretical characterization of balanced contextual bandits that match the regret bound of their direct method counterparts. The balancing technique is well-known in machine learning, especially in domain adaptation and studies in learning-theoretic frameworks [35, 65, 24]. There is a number of recent works which approach contextual bandits through the framework of causality [14, 15, 31, 41]. There is also a significant body of research that leverages balancing for offline evaluation and learning of contextual bandit or reinforcement learning policies from logged data [57, 28, 44, 27, 43, 59, 37, 60, 10, 38, 64, 26, 39, 58, 66]. In the offline setting, the complexity of the historical assignment policy is taken as given, and thus the difficulty of the offline evaluation and learning of optimal policies is taken as given. Therefore, these results lie at the opposite end of the spectrum from our work, which focuses on the online setting, where the complexity of the assignment policy is known and controlled by the algorithm. Methods for reducing the bias due to adaptive data collection have also been studied for non-contextual multi-armed bandits [63, 47], but the nature of the estimation in contextual bandits is qualitatively different. Importance weighted regression in contextual bandits was first mentioned in [4], but without a systematic motivation, analysis and evaluation of this technique. To our knowledge, our paper is the first work to integrate balancing in the online contextual bandit setting, to perform a large-scale evaluation of it against direct estimation method baselines with theoretical guarantees and to provide a theoretical characterization of balanced contextual bandits that match the regret bound of their direct method counterparts. The effect of importance weighted regression is also evaluated in [18], but this paper is a successor to the extended version of our paper.

Reweighting or balancing methods address model misspecification by making the estimation "doubly-robust,", robust against misspecification of the reward function, important here, and robust against the specification of the propensity score (not as important here because in the bandit setting we know the propensity score). The term "doubly-robust" comes from the extensive literature on offline policy evaluation [54]; it means that when comparing two policies using historical data, we get consistent estimates of the average difference in outcomes

for segments of the context whether we have either a well-specified model of rewards or not., because we have good model of the arm assignment policy (i.e., accurate propensity scores). In a contextual bandit, the learner controls the arm assignment policy conditional on the observed context and therefore has access to accurate propensity scores even in small samples. So, even when the reward model is severely misspecified, the learner can obtain more accurate value estimates of the reward function for each value of the context.

In real-world applications, such as health, education, recommender systems, there may be contextual variables that are highly predictive of user outcomes (e.g. previous health metrics, previous test scores, or previous consumer choices), but less important for optimizing arm assignment. We continue by showing that such contextual variables can be particularly problematic in terms of generating bias and variance in the early stages of learning. Even though in principle a variety of methods can be used to consistently estimate outcome models in the presence of complex, non-uniform arm assignment in large samples (Ridge regression as in [6, 45], ordinary least squares as in [32], LASSO as in [16]), in the early stages of learning where samples are small, these methods are only partially effective, and so the extent to which the method exacerbates the estimation problem will affect performance. In the domain of linear contextual bandits, we simplify the assignment policy by simplifying the estimated outcome model. We develop the bootstrap LASSO Thompson sampling and UCB contextual bandits, which use $L_1$ regularization in the estimation of the outcome models and circumvent the lack of closed-form solution by using the bootstrap to form a sampling distribution of the outcomes in order to obtain an approximate posterior for Thompson sampling and an upper confidence bound for UCB. We also propose a method to simplify the assignment policy with a form of smoothing that partitions the context space via a classification tree and defines the assignment rule for each partition rather than for each context distinctly. We show that, all else equal, using assignment policies that are simpler (in terms of how they vary with contextual variables) in the early learning phases of the algorithm can improve the rate of learning and decrease regret. Simple assignment rules may have other advantages as well; for example, [42] highlights the advantages of simplicity for interpretability in health applications of contextual bandits.

## 2 Methodological Designs in Contextual Bandits

In contextual bandits, contexts $x_1, x_2, \ldots, x_T \in \mathcal{X} \subset \mathbb{R}^d$ arrive sequentially. We have a finite set of $K$ arms, $\mathcal{A} = \{1, 2, \ldots, K\}$, which we wish to assign to each context upon its arrival. We posit that the data $(x_t, r_t(1), r_t(2), \ldots, r_t(K))$ are drawn **iid** from a fixed underlying joint distribution on $(x, r(1), r(2), \ldots, r(K))$, where $r(a) \in \mathbb{R}$ denotes the (random) reward under

arm $a$ and context $x$. Note that the marginal distribution on $x$ specifies the distribution from which the contexts are drawn. The observables are $(x_t, a_t, r_t(a_t))$; in particular, only the reward $r_t(a_t)$ for the chosen arm $a_t$ is observed. For each context $x \in \mathbb{R}^d$, the optimal assignment is $a^*(x) = \operatorname{argmax}_a \{\mathbb{E}[r(a)|x]\}$ and we let $a_t^* = a^*(x_t)$, which denotes the optimal assignment for context $x_t$. The objective is to find an assignment rule that sequentially assigns $a_t$ to minimize the cumulative expected regret $\sum_{t=1}^{\top} \mathbb{E}[r(a_t^*) - r(a_t)]$, where the assignment rule is a function of the previous observations $(x_j, a_j, r(a_j))$ for $j = 1, \ldots, t-1$ and of the new context $x_t$. We next discuss three methodological designs in contextual bandits.

## 2.1 Model Estimation

An important component for sequential arm assignment lies in modeling and estimating the conditional expected reward corresponding to each arm $a \in \mathcal{A}$ given context $x$, $\mu_a(x) = \mathbb{E}[r(a)|x]$. In a contextual bandit there are as many models to be estimated as arms. We do this estimation separately for each arm $a \in \mathcal{A}$ on the history of observations corresponding to this arm $\{(x_t, a_t, r_t(a_t)) \mid a_t = a\}$.

### 2.1.1 Parametric Estimation

In parametric estimation, we estimate $\mu_a(x)$ by a parametric form. A commonly used model in contextual bandits is the linear model, where $\mu_a(x) = x^\top \theta_a$, with unknown parameters $\theta_a$. More generally, estimation can be done via generalized linear models[2] $\mu_a(x) = g(x^\top \theta_a)$, where $g : \mathbb{R} \to \mathbb{R}$ is a strictly increasing and known link function. Denote $\mathbf{r}_a$ as the response vector and $\mathbf{X}_a$ as the covariate matrix of the history of observations assigned to $a$. The model parameters can be estimated via $L_1$ (LASSO [62] in linear models) or $L_2$ (ridge [33] in linear models) regularized regression. For a new context $x$, we wish to obtain the conditional mean $\hat{\mu}_a(x)$ of the reward associated with each arm $a \in A$ and its variance $\mathbb{V}(\hat{\mu}_a(x))$. In some cases, the estimates can be computed in closed form (e.g. in the case of a linear model estimated with ridge regression). However, in many cases (e.g. in the case of a linear model estimated with LASSO regression, or generalized linear models) exact computation is intractable and we must perform approximation [53]. When exact computation is intractable, bootstraping provides a viable way to approximate these quantities. More specifically, we can obtain a sampling distribution on $\mu_a(x)$ by training many regularized regression models on bootstrap samples drawn from $(\mathbf{X}_a, \mathbf{r}_a)$. With this sampling distribution, we can then easily to compute the mean estimate $\hat{\mu}_a(x)$ and the variance estimate $\mathbb{V}(\hat{\mu}_a(x))$. Note that as long as one can solve the underlying regression problem efficiently (either in closed form or via some fast

---

[2]GLM bandits are discussed in [46].

numeric scheme), the estimates can be constructed by bootstraping. Consequently, this provides a general-purpose way of computing $\hat{\mu}_a(x)$ and estimate $\mathbb{V}(\hat{\mu}_a(x))$ for contextual bandits.

### 2.1.2 Non-parametric Estimation

Parametric estimation can have high bias when the model is mis-specified (also called unrealizable in the literature). Non-parametric models, on the other hand, are more expressive and comparatively suffer less from the bias problem. In this paper, we consider non-parametric estimation of $\mu_a(x)$ by training a generalized random forest [9] on the history of observations assigned to $a$. Generalized random forest is an ensemble method and hence provides estimates of the conditional mean $\hat{\mu}_a(x)$ and variance $\mathbb{V}(\hat{\mu}_a(x))$. Related approaches based on random forests or decision trees have been proposed in [30, 29]. The generalized random forest is a method that preserves the core elements of random forests [19] including recursive partitioning, sub-sampling, and random split selection, but differs from traditional approaches in several ways, most notably that it integrates "honest" tree estimation [7]: the sample used to select the splits of the tree is independent from the sample used to estimate the improvement in fit yielded by a split. The "honesty" property of generalized random forests reduces bias and overfitting to outliers, which may be of particular concern in early stages of learning. In cases where the outcome functional form is complicated, as is often the case in practice, bandits based on non-parametric estimation tend to perform better. In addition, they can flexibly control for features which are confounders for the estimation of the reward functions (i.e., features that affect outcomes and were also used to determine assignments in earlier contexts).

## 2.2 Treatment Assignment Rules

Thompson sampling [61, 55, 5, 52] and upper confidence bound (UCB) [40, 12] are two different methods for assigning contexts to arms which are highly effective in dealing with the exploration-exploitation trade-off. In both methods, until every arm has been pulled at least once, the first contexts are assigned to arms in $\mathcal{A}$ at random with equal probability. At every time $t$ and for every arm $a$, the two methods use the history of observations $\{(x_t, a_t, r_t(a_t)) \mid a_t = a\}$, to obtain the estimates of the functions $\hat{\mu}_a(x)$ and $\mathbb{V}(\hat{\mu}_a(x))$ with the methods outlined in Section 2.1.

Thompson sampling assumes that the expected reward $\mu_a(x_t)$ associated with arm $a$ conditional on the context $x_t$ is Gaussian $\mathcal{N}(\hat{\mu}_a(x_t), \alpha^2 \mathbb{V}(\hat{\mu}_a(x_t)))$, where $\alpha$ is an appropriately chosen constant. Then, it draws a sample $\tilde{\mu}_a(x_t)$ from the distribution of each arm $a \in$ and

7

context $x_t$ is then assigned to the arm with the highest sample, $a_t = \text{argmax}_a\{\tilde{\mu}_a(x_t)\}$.

On the other hand, UCB computes upper confidence bounds for the expected reward $\mu_a(x_t)$ of context $x_t$ associated with each arm $a \in \mathcal{A}$ and assigns the context to the arm with the highest upper confidence bound, $a_t = \text{argmax}_a \left\{\hat{\mu}_a(x_t) + \alpha\sqrt{\mathbb{V}(\hat{\mu}_a(x_t))}\right\}$, where $\alpha$ is an appropriately chosen constant.

## 2.3 Balancing in Contextual Bandits

In this section, we use balancing methods from the causal inference literature to improve the existing UCB and Thompson sampling algorithms by balancing features between arms to reduce bias. We focus on the method of inverse propensity weighting (IPW) [36]. Denote **r** as the reward vector, **X** as the context matrix and **a** as the arm assignment vector of all previous observations.

For UCB, we train a multi-class logistic regression model of **a** on **X** to estimate the assignment probabilities $p_a(x), a \in \mathcal{A}$, also known as propensity scores. Note that UCB has deterministic assignment rules, so that conditional on the batch, the propensity scores are either zero or one. However, we can consider the ordering of contexts' arrival as random and use $\hat{p}_a(x)$ as a balancing weight to account for non-uniform assignment in previous contexts.

For Thompson sampling, the propensity scores are in principle known because Thompson sampling performs probability matching, i.e., it assigns a context to an arm with the probability that this arm is optimal. Since computing the propensity scores involves high order integration, they can be approximated via Monte-Carlo simulation. Each iteration draws a sample from the posterior reward distribution of every arm $a$ conditional on $x$, where the posterior is the one that the algorithm considered at the end of a randomly selected prior time period. The propensity score $p_a(x)$ is the fraction of the Monte-Carlo iterations in which arm $a$ has the highest sampled reward, where the arrival time of context $x$ is treated as random.

With these propensity estimates, we can modify the model estimation as follows. For parametric estimation, we define the weight of each observation $(x, a, r)$ as the inverse of the estimated propensity score,

$$w_a = 1/\hat{p}_a(x)$$

and we train the weighted counterparts of the regularized regressions discussed in Section 2.1.1. For non-parametric estimation (via generalized random forest), one alternative is to construct an augmented covariate matrix $\tilde{\mathbf{X}}_a$ and reward vector $\tilde{\mathbf{r}}_a$ by replicating $[w_a]$ times each observation $(x, a, r)$ and subsequently to estimate the generalized random forest on $\tilde{\mathbf{X}}_a$ and $\tilde{\mathbf{r}}_a$. Another alternative is to treat the propensity score of each observation as a

contextual variable and estimate the generalized random forest on $[\mathbf{X}_a : \mathbf{p}_a]$ and $\mathbf{r}_a$, where $\mathbf{p}_a$ is the vector of the propensity scores for previous observations for arm $a$.

In both cases, weighting the observations by the inverse propensity scores reduces bias, but even when the propensity scores are known it increases variance, particularly when they are small. Consequently, since eventually assignment probabilities should approach zero or one for all arms and contexts clipping the propensity scores [25, 39] with some threshold $\gamma$, e.g. 0.1 helps control the variance increase.

Finally, note that one could integrate in the contextual bandit estimation other covariate balancing methods, such as the method of approximate residual balancing [8] or the method of [38]. For instance, with approximate residual balancing one would use as weights

$$w_a = \operatorname*{argmin}_w \left\{ (1 - \zeta)\|w\|_2^2 + \zeta\|\bar{x} - \mathbf{X}_a^\top w\|_\infty^2 \text{ s.t. } \sum_{t:a_t=a} w_t = 1 \text{ and } 0 \leq w_t \leq n_a^{-2/3} \right\}$$

where $\zeta \in (0, 1)$ is a tuning parameter, $n_a = \sum_{t=1}^T \mathbf{1}\{a_t = a\}$ and $\bar{x} = \frac{1}{T}\sum_{t=1}^T x_t$ and then use $w_a$ to modify the parametric and non-parametric model estimation as outlined before.

# 3 The Effects of Balancing

In this section, we study empirically and theoretically the effects of balancing as introduced in 2.3. For concreteness, we focus on the method of inverse propensity weighting for balancing and on linear models with $L_2$ regularization (i.e. ridge) and refer to the resulting algorithms as *balanced linear Thompson sampling* (BLTS) and *balanced linear upper confidence bound* (BLUCB), as given in Algorithm 1 and Algorithm 2.

BLTS and BLUCB build on linear contextual bandits LinTS [6] and LinUCB [45] respectively. In LinTS and LinUCB, the expected reward is assumed to be a linear function of the context $x_t$ with some unknown coefficient vector $\theta_a$, $\mathbb{E}[r_t(a)|x_t = x] = x^\top\theta_a$, and the variance is typically assumed to be constant $\mathbb{V}[r_t(a)|x_t = x] = \sigma_a^2$. At time $t$, LinTS and LinUCB apply ridge regression with regularization parameter $\lambda$ to the history of observations $(X_a, r_a)$ for each arm $a \in \mathcal{A}$, in order to obtain an estimate $\hat{\theta}_a$ and its variance $\mathbb{V}_a(\hat{\theta}_a)$. For the new context $x_t$, $\hat{\theta}_a$ and its variance are used by LinTS and LinUCB to obtain the conditional mean $\hat{\mu}_a(x_t) = x_t^\top\hat{\theta}_a$ of the reward associated with each arm $a \in A$, and its variance $\mathbb{V}(\hat{\mu}_a(x_t)) = x_t^\top\mathbb{V}(\hat{\theta}_a)x_t$.

BLTS and BLUCB are linear contextual bandit algorithms that perform balanced estimation of the model of all arms in order to obtain a Gaussian distribution and an upper confidence bound respectively for the reward associated with each arm conditional on the context. The idea is that at every time $t$, the linear contextual bandit weighs each observation

$(x_\tau, a_\tau, r_\tau(a_\tau))$, $\tau = 1, \ldots, t$ in the history up to time $t$ by the inverse propensity score, $p_{a_\tau}(x_\tau)$. Then, for each arm $a \in \mathcal{A}$, the linear contextual bandit weighs each observation $(x_\tau, a, r_\tau(a))$ in the history of arm $a$ by $w_a = 1/p_a(x_\tau)$ and uses weighted regression to obtain the estimate $\hat{\theta}_a^{\text{balanced}}$ with variance $\mathbb{V}(\hat{\theta}_a^{\text{balanced}})$.

In BLTS (Algorithm 1), the observations are weighted by the known propensity scores. For every arm $a$, the history $(X_a, r_a, p_a)$ is used to obtain a balanced estimate $\hat{\theta}_a^{\text{BLTS}}$ of $\theta_a$ and its variance $\mathbb{V}(\hat{\theta}_a^{\text{BLTS}})$ which produce a normally distributed estimate of $\tilde{\mu}_a \sim \mathcal{N}\left(x_t^\top \hat{\theta}_a^{\text{BLTS}}, \alpha^2 x_t^\top \mathbb{V}(\hat{\theta}_a^{\text{BLTS}}) x_t\right)$ of the reward of arm $a$ for context $x_t$, where $\alpha$ is a parameter of the algorithm.

In BLUCB (Algorithm 2), the observations are weighted by the estimated propensity scores. For every arm $a$, the history $(X_a, r_a, \hat{p}_a)$ is used to obtain a balanced estimate $\hat{\theta}_a^{\text{BLUCB}}$ of $\theta_a$ and its variance $\mathbb{V}(\hat{\theta}_a^{\text{BLUCB}})$. These are used to construct the upper confidence bound, $x_t^\top \hat{\theta}_a + \alpha \sqrt{x_t^\top \mathbb{V}(\hat{\theta}_a^{\text{BLUCB}}) x_t}$, for the reward of arm $a$ for context $x_t$, where $\alpha$ is a constant. (For some results, e.g., [13], $\alpha$ needs to be slowly increasing in $t$.)

We present computational results that compare the performance of BLTS and BLUCB, with the existing state-of-the-art linear contextual bandits algorithms LinTS [6] and LinUCB [45]. More specifically, we first present a simple synthetic example that simulates bias in the training data by under-representation or over-representation of certain regions of the context space and investigates the performance of the considered linear contextual bandits both when the outcome model of the arms matches the true reward generative process and when it does not match the true reward generative process. Second, we conduct an experiment by leveraging 300 public, supervised cost-sensitive classification datasets to obtain contextual bandit problems, treating the features as the context, the labels as the arms and revealing only the reward for the chosen label. We show that BLTS performs better than LinTS and that BLUCB performs better than LinUCB. The randomized assignment nature of Thompson sampling facilitates the estimation of the arms' outcomes models compared to UCB, and as a result LinTS outperforms LinUCB and BLTS outperforms BLUCB. Overall, BLTS has the best performance. In the supplemental material, we include experiments against the policy-based contextual bandit from [4] which is statistically optimal but it is also outperformed by BLTS. Finally, we give a theoretical guarantee that matches the existing regret performance of BLTS and BLUCB.

## 3.1    A Synthetic Example

This simulated example aims to reflect in a simple way two issues that often arise in practice. The first issue is the presence of bias in the training data by under-representation or over-

**Algorithm 1** Balanced Linear Thompson Sampling

---

1: **Input:** Regularization parameter $\lambda > 0$, propensity score threshold $\gamma \in (0, 1)$, constant $\alpha$ (deafult is 1)
2: Set $\hat{\theta}_a^{\mathrm{BLTS}} \leftarrow \mathbf{null}, B_a \leftarrow \mathbf{null}, \forall a \in \mathcal{A}$
3: Set $X_a \leftarrow$ empty matrix, $r_a \leftarrow$ empty vector $\forall a \in \mathcal{A}$
4: **for** $t = 1, 2, \ldots, T$ **do**
5:     **if** $\exists a \in \mathcal{A}$ s.t. $\hat{\theta}_a^{\mathrm{BLTS}} = \mathbf{null}$ or $B_a = \mathbf{null}$ **then**
6:         Select $a \sim \mathrm{Uniform}(\mathcal{A})$
7:     **else**
8:         Draw $\tilde{\theta}_a$ from $\mathcal{N}\left(\hat{\theta}_a^{\mathrm{BLTS}}, \alpha^2 \mathbb{V}(\hat{\theta}_a^{\mathrm{BLTS}})\right)$ for all $a \in \mathcal{A}$
9:         Select $a = \arg\max_{a \in \mathcal{A}} x_t^\top \tilde{\theta}_a$
10:     **end if**
11:     Observe reward $r_t(a)$.
12:     Set $W_a \leftarrow$ empty matrix
13:     **for** $\tau = 1, \ldots, t$ **do**
14:         **if** $a_\tau = a$ **then**
15:             Compute $p_a(x_\tau)$ and set $w = \frac{1}{\max(\gamma, p_a(x_\tau))}$
16:             $W_a \leftarrow \mathrm{diag}(W_a, w)$
17:         **end if**
18:     **end for**
19:     $X_a \leftarrow [X_a : x_t^\top]$
20:     $B_a \leftarrow X_a^\top W_a X_a + \lambda \mathbf{I}$
21:     $r_a \leftarrow [r_a : r_t(a)]$
22:     $\hat{\theta}_a^{\mathrm{BLTS}} \leftarrow B_a^{-1} X_a^\top W_a r_a$
23:     $\mathbb{V}(\hat{\theta}_a^{\mathrm{BLTS}}) \leftarrow B_a^{-1}\left((r_a - X_a^\top \hat{\theta}_a^{\mathrm{BLTS}})^\top W_a(r_a - X_a^\top \hat{\theta}_a^{\mathrm{BLTS}})\right)$
24: **end for**

---

representation of certain regions. A personalized policy that is trained based on such data and is applied to the entire context space will result in biased decisions for certain contexts. The second issue is the problem of mismatch between the true reward generative process and the functional form used for estimation of the outcome model of the arms, which is common in applications with complex generative models. Model mis-specification aggravates the presence of bias in the learned policies.

We use this simple example to present in an intuitive manner why balancing and randomized assignment rule help with these issues, before moving on to a large-scale evaluation of the algorithms in real datasets in the next section.

Consider a simulation design where there is a warm-start batch of training observations, but it consists of contexts focused on one region of the context space. There are three arms $\mathcal{A} = \{0, 1, 2\}$ and the contexts $x_t = (x_{t,0}, x_{t,1})$ are two-dimensional with $x_{t,j} \sim \mathcal{N}(0, 1)$, $j \in \{0, 1\}$. The rewards corresponding to each arm $a \in \mathcal{A}$ are generated as follows; $r_t(0) =$

**Algorithm 2** Balanced Linear UCB

---

1: **Input:** Regularization parameter $\lambda > 0$, propensity score threshold $\gamma \in (0, 1)$, constant $\alpha$.

2: Set $\hat{\theta}_a^{\mathrm{BLUCB}} \leftarrow \mathbf{null}, B_a \leftarrow \mathbf{null}, \forall a \in \mathcal{A}$

3: Set $X_a \leftarrow$ empty matrix, $r_a \leftarrow$ empty vector $\forall a \in \mathcal{A}$

4: **for** $t = 1, 2, \ldots, T$ **do**

5:     **if** $\exists a \in \mathcal{A}$ s.t. $\hat{\theta}_a^{\mathrm{BLUCB}} = \mathbf{null}$ or $B_a = \mathbf{null}$ **then**

6:         Select $a \sim \mathrm{Uniform}(\mathcal{A})$

7:     **else**

8:         Select $a = \arg\max_{a \in \mathcal{A}} \left( x_t^\top \hat{\theta}_a^{\mathrm{BLUCB}} + \alpha \sqrt{x_t^\top \mathbb{V}(\hat{\theta}_a^{\mathrm{BLUCB}}) x_t} \right)$

9:     **end if**

10:     Observe reward $r_t(a)$.

11:     Set $W_a \leftarrow$ empty matrix

12:     **for** $\tau = 1, \ldots, t$ **do**

13:         **if** $a_\tau = a$ **then**

14:             Estimate $\hat{p}_a(x_\tau)$ and set $w = \frac{1}{\max(\gamma, \hat{p}_a(x_\tau))}$

15:             $W_a \leftarrow \mathrm{diag}(W_a, w)$

16:         **end if**

17:     **end for**

18:     $X_a \leftarrow [X_a : x_t^\top]$

19:     $B_a \leftarrow X_a^\top W_a X_a + \lambda \mathbf{I}$

20:     $r_a \leftarrow [r_a : r_t(a)]$

21:     $\hat{\theta}_a \leftarrow B_a^{-1} X_a^\top W_a r_a$

22:     $\mathbb{V}(\hat{\theta}_a^{\mathrm{BLUCB}}) \leftarrow B_a^{-1} \left( (r_a - X_a^\top \hat{\theta}_a^{\mathrm{BLUCB}})^\top W_a (r_a - X_a^\top \hat{\theta}_a^{\mathrm{BLUCB}}) \right)$

23: **end for**

---

$0.5(x_{t,0} + 1)^2 + 0.5(x_{t,1} + 1)^2 + \epsilon_t$, $r_t(1) = 1 + \epsilon_t$, and $r_t(2) = 2 - 0.5(x_{t,0} + 1)^2 - 0.5(x_{t,1} + 1)^2 + \epsilon_t$, where $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$, $\sigma^2 = 0.01$. The expected values of the three arms' rewards are shown in Figure 1.

In the warm-start data, $x_{t,0}$ and $x_{t,1}$ are generated from a truncated normal distribution $\mathcal{N}(0, 1)$ on the interval $(-1.15, -0.85)$, while in subsequent data $x_{t,0}$ and $x_{t,1}$ are drawn from $\mathcal{N}(0, 1)$ without the truncation. Each one of the 50 warm-start contexts is assigned to one of the three arms at random with equal probability. Note that the warm-start contexts belong to a region of the context space where the reward surfaces do not change much with the context. Therefore, when training the reward model for the first time, the estimated reward of arm $a = 2$ (blue) is the highest, the one of arm $a = 1$ (yellow) is the second highest and the one of arm $a = 0$ (red) is the lowest across the context space.

We run our experiment with a learning horizon $T = 10000$. The regularization parameter $\lambda$, which is present in all algorithms, is chosen via cross-validation every time the model is updated. The constant $\alpha$, which is present in all algorithms, is optimized among values
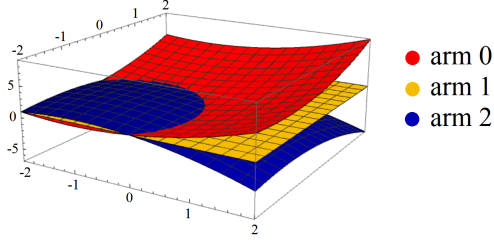
Figure 1: Expectation of each arm's reward, $\mathbb{E}[r_t(0)] = 0.5(x_{t,0}+1)^2 + 0.5(x_{t,1}+1)^2$ (red), $\mathbb{E}[r_t(1)] = 1$ (yellow), $\mathbb{E}[r_t(2)] = 2 - 0.5(x_{t,0}+1)^2 - 0.5(x_{t,1}+1)^2$ (blue).

$0.25, 0.5, 1$ in the Thompson sampling bandits (the value $\alpha = 1$ corresponds to standard Thompson sampling, [22] suggest that smaller values may lower regret) and among values $1, 2, 4$ in the UCB bandits [22]. The propensity threshold $\gamma$ for BLTS and BLUCB is optimized among the values $0.01, 0.05, 0.1, 0.2$.

### 3.1.1 Well-Specified Outcome Models

In this section, we compare the behavior of LinTS, LinUCB, BLTS and BLUCB when the outcome model of the contextual bandits is well-specified, i.e., it includes both linear and quadratic terms. Note that this is still in the domain of linear contextual bandits, if we treat the quadratic terms as part of the context.

First, we compare LinTS and LinUCB. Figure 2a shows that the uncertainty and the stochastic nature of LinTS leads to a "dispersed" assignment of arms $a = 1$ and $a = 2$ and to the crucial assignment of a few contexts to arm $a = 0$. This allows LinTS to start decreasing the bias in the estimation of all three arms in the subsequent time periods. Within the first few learning observations, LinTS estimates the outcome models of all three arms correctly and finds the optimal assignment. On the other hand, Figure 2b, shows that the deterministic nature of LinUCB assigns entire regions of the context space to the same arm. As a result not enough contexts are assigned to $a = 0$ and LinUCB delays the correction of bias in the estimation of this arm. Another way to understand the problem is that the outcome model in the LinUCB bandit has biased coefficients combined with estimated uncertainty that is too low to incentivize the exploration of arm $a = 0$ initially. LinUCB finds the correct assignment after 240 observations.

Second, we study the performance of BLTS and BLUCB. In Figure 2d, we observe that balancing has a significant impact on the performance of UCB, since BLUCB finds the optimal assignment after 110 observations, much faster than LinUCB. This is because the few observations of arm $a = 0$ outside of the context region of the warm-start batch are weighted more heavily by BLUCB. As a result, BLUCB, despite its deterministic nature

(a) Well-specified LinTS



(b) Well-specified LinUCB

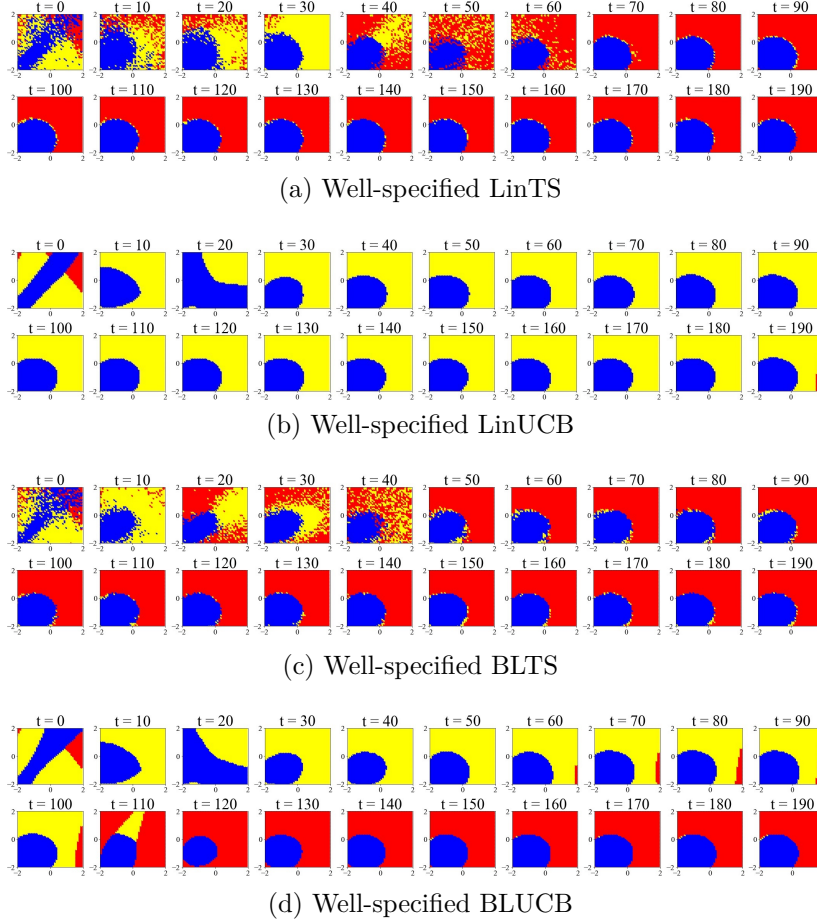

(c) Well-specified BLTS



(d) Well-specified BLUCB

Figure 2: Evolution of the arm assignment in the context space for well-specified LinTS, LinUCB, BLTS, BLUCB.

which complicates estimation, is able to reduce its bias more quickly via balancing Figure 2c shows that BLTS is also able to find the optimal assignment a few observations earlier than LinTS. Figure 3 shows the evolution of the estimation bias for all three arms for the well-specified LinTS, LinUCB, BLTS and BLUCB.



Figure 3: Evolution of the potential outcomes estimation bias in the $(x_0, x_1)$ context space for well-specified LinTS, LinUCB, BLTS and BLUCB. Blue indicates that the actual estimate is lower than the predicted, whereas red indicates that the actual estimate is higher than the predicted.

(a) Mis-specified LinTS



(b) Mis-specified LinUCB



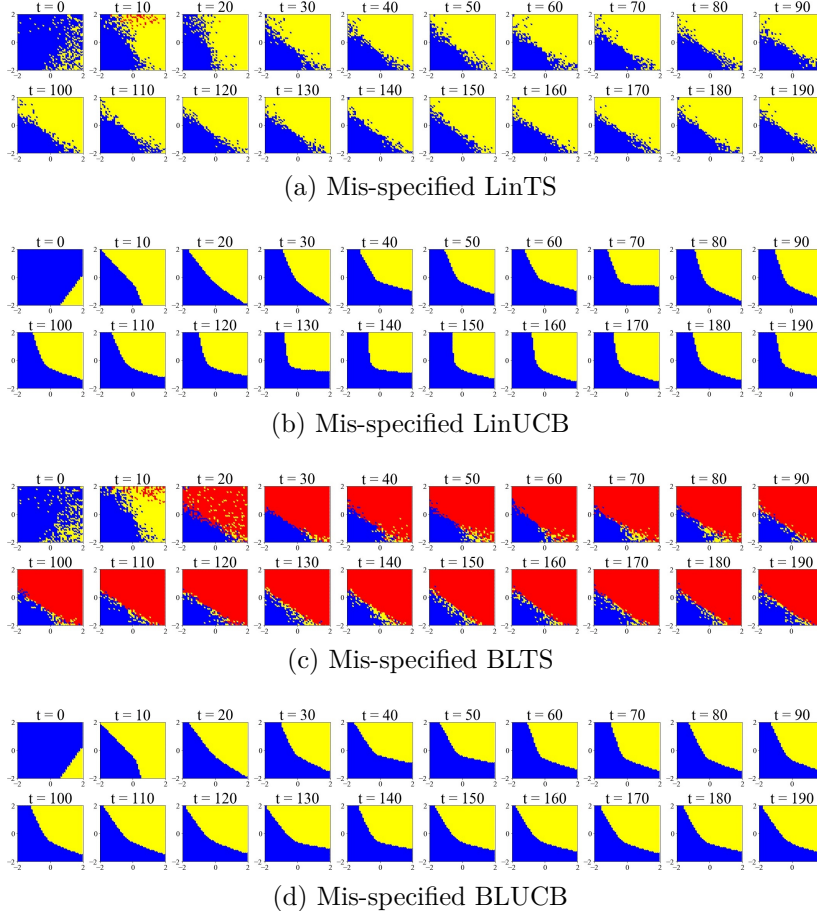(c) Mis-specified BLTS



(d) Mis-specified BLUCB

Figure 4: Evolution of the arm assignment in the context space for mis-specified LinTS, LinUCB, BLTS, BLUCB.

The first column of Table 1 shows the percentage of simulations in which LinTS, LinUCB, BLTS and BLUCB find the optimal assignment within $T = 10000$ contexts for the well-specified case. BLTS outperforms all other algorithms by a large margin.

### 3.1.2 Mis-Specified Outcome Models

We now study the behavior of LinTS, LinUCB, BLTS and BLUCB when the outcome models include only linear terms of the context and therefore are mis-specified. In real-world domains, the true data generative process is complex and very difficult to capture by the simpler outcome models assumed by the learning algorithms. Hence, model mismatch is very likely.

We first compare LinTS and LinUCB. In Figures 4a, 4b, we see that during the first time periods, both bandits assign most contexts to arm $a = 2$ and a few contexts to arm $a = 1$. LinTS finds faster than LinUCB the linearly approximated area in which arm $a = 2$ is suboptimal. However, both LinTS and LinUCB have trouble identifying that the optimal
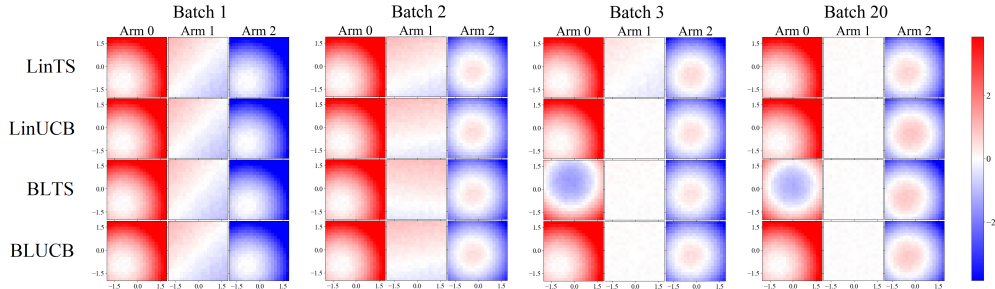
15

Figure 5: Evolution of the potential outcomes estimation bias in the $(x_0, x_1)$ context space for mis-specified LinTS, LinUCB, BLTS and BLUCB. Blue indicates that the actual estimate is lower than the predicted, whereas red indicates that the actual estimate is higher than the predicted.

arm is $a = 0$. Due to the low estimate of $a = 0$ from the mis-representative warm-start observations, LinUCB does not assign contexts to arm $a = 0$ for a long time and therefore, delays to estimate the model of $a = 0$ correctly. LinTS does assign a few contexts to arm $a = 0$, but they are not enough to quickly correct the estimation bias of arm $a = 0$ either. On the other hand, BLTS is able to harness the advantages of the stochastic assignment rule of Thompson sampling. The few contexts assigned to arm $a = 0$ are weighted more heavily by BLTS. Therefore, as shown in Figure 4c, BLTS corrects the estimation error of arm $a = 0$ and finds the (constrained) optimal assignment already after 20 observations. On the other hand, BLUCB does not handle better than LinUCB the estimation problem created by the deterministic nature of the assignment in the mis-specified case, as shown in Figure 4d. Figure 5 shows the evolution of the estimation bias for all three arms for the mis-specified LinTS, LinUCB, BLTS and BLUCB.

The second column of table 1 shows the percentage of simulations in which LinTS, LinUCB, BLTS and BLUCB find the optimal assignment within $T = 10000$ contexts for the mis-specified case. Again, BLTS has a strong advantage.

This simple synthetic example allowed us to explain transparently where the benefits of balancing in linear bandits stem from. Balancing helps escape biases in the training data and can be more robust in the case of model mis-specification. While, as we proved, balanced linear contextual bandits share the same strong theoretical guarantees, this indicates towards their better performance in practice compared to other contextual bandits with linear realizability assumption. We investigate this further in the next section with an extensive evaluation on real cost-sensitive classification datasets.

16

|        | Well-Specified | Mis-Specified |
|--------|----------------|---------------|
| LinTS  | 84%            | 39%           |
| LinUCB | 51%            | 29%           |
| BLTS   | 92%            | 58%           |
| BLUCB  | 79%            | 30%           |

Table 1: Percentage of simulations in which LinTS, LinUCB, BLTS and BLUCB find the optimal assignment within learning horizon of 10000 contexts

## 3.2 Multiclass Classification with Bandit Feedback

Adapting a classification task to a bandit problem is a common method for comparing contextual bandit algorithms [28], [4], [18]. In a classification task, we assume data are drawn IID from a fixed distribution: $(x, c) \sim D$, where $x \in \mathcal{X}$ is the context and $c \in 1, 2, \ldots, K$ is the class. The goal is to find a classifier $\pi : \mathcal{X} \to \{1, 2, \ldots, K\}$ that minimizes the classification error $\mathbb{E}_{(x,c)\sim D}\mathbf{1}\{\pi(x) \neq c\}$. The classifier can be seen as an arm-selection policy and the classification error is the policy's expected regret. Further, if only the loss associated with the policy's chosen arm is revealed, this becomes a contextual bandit setting. So, at time $t$, context $x_t$ is sampled from the dataset, the contextual bandit selects arm $a_t \in \{1, 2, \ldots, K\}$ and observes reward $r_t(a_t) = \mathbf{1}\{a_t = c_t\}$, where $c_t$ is the unknown, true class of $x_t$. The performance of a contextual bandit algorithm on a dataset with $n$ observations is measured with respect to the normalized cumulative regret, $\frac{1}{n}\sum_{t=1}^{n}(1 - r_t(a_t))$.

We use 300 multiclass datasets from the Open Media Library (OpenML). The datasets vary in number of observations, number of classes and number of features. Table 2 summarizes the characteristics of these benchmark datasets. Each dataset is randomly shuffled.

| Observations | Datasets |
|--------------|----------|
| $\leq 100$ | 58 |
| $> 100$ and $\leq 1000$ | 152 |
| $> 1000$ and $\leq 10000$ | 57 |
| $> 10000$ | 33 |

| Classes | Count |
|---------|-------|
| 2 | 243 |
| $> 2$ and 10 | 48 |
| $> 10$ | 9 |

| Features | Count |
|----------|-------|
| $\leq 10$ | 154 |
| $> 10$ and $\leq 100$ | 106 |
| $> 100$ | 40 |

Table 2: Characteristics of the 300 datasets used for the experiments of multiclass classification with bandit feedback.

We evaluate LinTS, BLTS, LinUCB and BLUCB on these 300 benchmark datasets. We run each contextual bandit on every dataset for different choices of input parameters. The regularization parameter $\lambda$, which is present in all algorithms, is chosen via cross-validation every time the model is updated. The constant $\alpha$, which is present in all algorithms, is optimized among values $0.25, 0.5, 1$ in the Thompson sampling bandits [22] and among

values $1, 2, 4$ in the UCB bandits [22]. The propensity threshold $\gamma$ for BLTS and BLUCB is optimized among the values $0.01, 0.05, 0.1, 0.2$. Apart from baselines that belong in the family of contextual bandits with linear realizability assumption and have strong theoretical guarantees, we also evaluate the policy-based ILOVETOCONBANDITS (ILTCB) from [4] that does not estimate a model, but instead it assumes access to an oracle for solving fully supervised cost-sensitive classification problems and achieves the statistically optimal regret.



Figure 6: Pairwise comparison of LinTS, BLTS, LinUCB, BLUCB on 300 classification datasets. BLUCB outperforms LinUCB. BLTS outperforms LinTS, LinUCB, BLUCB.

Figure 6 shows the pairwise comparison of LinTS, BLTS, LinUCB, BLUCB and ILTCB on the 300 classification datasets. Each point corresponds to a dataset. The $x$ coordinate is

the normalized cumulative regret of the column bandit and the $y$ coordinate is the normalized cumulative regret of the row bandit. The point is blue when the row bandit has smaller normalized cumulative regret and wins over the column bandit. The point is red when the row bandit loses from the column bandit. The point's size grows with the significance of the win or loss.

The first important observation is that the improved model estimation achieved via balancing leads to better practical performance across a large number of contextual bandit instances. Specifically, BLTS outperforms LinTS and BLUCB outperforms LinUCB. The second important observation is that deterministic assignment rule bandits are at a disadvantage compared to randomized assignment rule bandits. The improvement in estimation via balancing is not enough to outweigh the fact that estimation is more difficult when the assignment is deterministic and BLUCB is outperformed by LinTS. Overall, BLTS which has both balancing and a randomized assignment rule, outperforms all other linear contextual bandits with strong theoretical guarantees. BLTS also outperforms the model-agnostic ILTCB algorithm.

We refer the reader to Appendix B of the supplemental material for details on the datasets.

## 3.3   Theoretical Guarantee

Here we establish theoretical guarantees of BLTS and BLUCB that are comparable to LinTS and LinUCB. We start with a few technical assumptions that are standard in the contextual bandits literature.

**Assumption 1.** *__Linear Realizability:__ There exist parameters $\{\theta_a\}_{a \in \mathcal{A}}$ such that given any context $x$, $\mathbb{E}[r_t(a)|x] = x^\top \theta_a, \forall a \in \mathcal{A}, \forall t \geq 0$.*

We use the standard (frequentist) regret criterion to measure performance as defined next:

**Definition 1.** *The instantaneous regret at iteration $t$ is $x_t^\top \theta_{a_t^*} - x_t^\top \theta_{a_t}$, where $a_t^*$ is the optimal arm at iteration $t$ and $a_t$ is the arm taken at iteration $t$. The cumulative regret $R(T)$ with horizon $T$ is the defined as $R(T) = \sum_{t=1}^{T} \left( x_t^\top \theta_{a_t^*} - x_t^\top \theta_{a_t} \right)$.*

**Definition 2.** *We denote the canonical filtration of the underlying contextual bandits problem by $\{\mathcal{F}_t\}_{t=1}^{\infty}$, where $\mathcal{F}_t = \sigma(\{x_s\}_{s=1}^{t}, \{a_s\}_{s=1}^{t}, \{r_s(a_s)\}_{s=1}^{t}, x_{t+1})$: the sigma algebra[3] generated by all the random variables up to and including iteration $t$, plus $x_{t+1}$. In other words, $\mathcal{F}_t$ contains all the information that is available before making the decision for iteration $t + 1$.*

---

[3]All the random variables $x_t, a_t, r_t$ are defined on some common underlying probability space, which we do not write out explicitly here.

Next, we make the standard assumptions on the regularity of the distributions:

**Assumption 2.** *For each $a \in \mathcal{A}$ and every $t \geq 1$:*

1. **Sub-Gaussian Noise:** *$r_t(a) - x_t^\top \theta_a$ is conditionally sub-Gaussian: there exists some $L_a > 0$, such that $\mathbb{E}[e^{s(r_t(a) - x_t^\top \theta_a)} \mid \mathcal{F}_{t-1}] \leq \exp(\frac{s^2 L_a^2}{2}), \forall s, \forall x_t$.*

2. **Bounded Contexts and Parameters:** *The contexts $x_t$ and parameters $\theta_a$ are assumed to be bounded. Consequently, without loss of generality, we can rescale them such that $\|x_t\|_2 \leq 1, \|\theta_a\|_2 \leq 1, \forall a, t$.*

**Remark 1.** *Note that we make no assumption of the underlying $\{x_t\}_{t=1}^\infty$ process: the contexts $\{x_t\}_{t=1}^\infty$ need not to be fixed beforehand or come from some stationary process. Further, they can even be adapted to $\sigma(\{x_s\}_{s=1}^t, \{a_s\}_{s=1}^t, \{r_s(a_s)\}_{s=1}^t)$, in which case they are called adversarial contexts in the literature as the contexts can be chosen by an adversary who chooses a context after observing the arms played and the corresponding rewards. If $\{x_t\}_{t=1}^\infty$ is an IID process, then the problem is known as stochastic contextual bandits. From this viewpoint, adversarial contextual bandits are more general, but the regret bounds tend to be worse. Both are studied in the literature.*

**Theorem 1.** *Under Assumption 1 and Assumption 2:*

1. *If BLTS is run with $\alpha = \sqrt{\frac{\log \frac{1}{\delta}}{\epsilon}}$ in Algorithm 1, then with probability at least $1 - \delta$, $R(T) = \tilde{O}\left(d\sqrt{\frac{KT^{1+\epsilon}}{\epsilon}}\right)$.*

2. *If BLUCB is run with $\alpha = \sqrt{\log \frac{TK}{\delta}}$ in Algorithm 2, then with probability at least $1 - \delta$, $R(T) = \tilde{O}\left(\sqrt{TdK}\right)$.*

We refer the reader to Appendix A of the supplemental material for the regret bound proofs.

**Remark 2.** *The above bound essentially matches the existing state-of-the art regret bounds for linear Thompson sampling with direct model estimation (e.g. [6]). Note that in [6], an infinite number of arms is also allowed, but all arms share the same parameter $\theta$. The final regret bound is $\tilde{O}\left(d^2 \frac{\sqrt{T^{1+\epsilon}}}{\epsilon}\right)$. Note that even though no explicit dependence on $K$ is present in the regret bound (and hence our regret bound appears as a factor of $\sqrt{K}$ worse), this is to be expected, as we have $K$ parameters to estimate, one for each arm. Note that here we do not assume any structure on the $K$ arms; they are just $K$ stand-alone parameters, each of which needs to be independently estimated. Similarly, for BLUCB, our regret bound is $\tilde{O}\left(\sqrt{TdK}\right)$, which is a factor of $\sqrt{K}$ worse than that of [23], which establishes a $\tilde{O}\left(\sqrt{Td}\right)$ regret bound.*

*Again, this is because a single true $\theta^*$ is assumed in [23], rather than $K$ arm-dependent parameters.*

*Of course, we also point out that our regret bounds are not tight, nor do they achieve state-of-the-art regret bounds in contextual bandits algorithms in general. The lower bound $\Omega(\sqrt{dT})$ is established in [23] for linear contextual bandits (again in the context of a single parameter $\theta$ for all $K$ arms). In general, UCB based algorithms ([11, 23, 21, 1]) tend to have better (and sometimes near-optimal) theoretical regret bounds. In particular, the state-of-the-art bound of $O(\sqrt{dT \log K})$ for linear contextual bandits is given in [21] (optimal up to a $O(\log K)$ factor). However, as mentioned in the introduction, Thompson sampling based algorithms tend to perform much better in practice (even though their regret bounds tend not to match UCB based algorithms, as is also the case here). Hence, our objective here is not to provide state-of-the-art regret guarantees. Rather, we are motivated to design algorithms that have better empirical performance (compared to both the existing UCB style algorithms and Thompson sampling style algorithms), which also enjoy the baseline theoretical guarantee.*

*Finally, we give some quick intuition for the proof. For BLTS, we first show that estimated means concentrate around true mean (i.e. $x_t^\top \hat\theta_a$ concentrates around $x_t^\top \theta_a$). Then, we establish that sampled means concentrate around the estimated means (i.e. $x_t^\top \tilde\theta_a$ concentrates around $x_t^\top \hat\theta_a$). These two steps together indicate that the sampled mean is close to the true mean. A further consequence of that is we can then bound the instantaneous regret (regret at each time step $t$) in terms of the sum of two standard deviations: one corresponds to the optimal arm at time $t$, the other corresponds to the actual selected arm at $t$. The rest of the proof then follows by giving tight characterizations of these two standard deviations. For BLUCB, the proof again utilizes the first concentration mentioned above: the estimated means concentrate around true mean (note that there is no sampled means in BLUCB). The rest of the proof adopts a similar structure as in [23].*

# 4    Additional Estimation Considerations

In this section, we investigate three additional estimation considerations and present further evidence from simulations that demonstrate our design outlined in Section 2 can be more effective in a variety of different contexts.

## 4.1    The Effect of Simpler Outcome Models: LASSO v.s. Ridge

In contextual bandits, the assignment to an arm is a function of the contexts and not all contexts have the same assignment probabilities. This creates an environment with

confounding, often in combination with small samples, in which the choice of the estimation method plays a very significant role. In this section, we study the significance maintaining a simple outcome model in contextual bandits. The potential outcome models corresponding to each arm are estimated in every batch and are used by the bandit to determine the assignment of contexts to arms, e.g. via the construction of upper confidence bounds. Therefore, maintaining a simple outcome model results in a simple assignment model.

We compare a contextual bandit that maintains a simpler outcome model, such as LASSO, with a contextual bandit that maintains a more complex outcome model, such as ridge. As mentioned in 2, exact computation of the reward mean and variance in the case of LASSO is intractable and we must perform approximation [53]. Bootstraping provides a viable way to approximate the quantities. More specifically, we can obtain a sampling distribution on $\mu_a(x)$ by training many regularized regression models on bootstrap samples drawn from $(\mathbf{X}_a, \mathbf{r}_a)$. With this sampling distribution, we can then easily to compute the mean estimate $\hat{\mu}_a(x)$ and the variance estimate $\hat{\sigma}_a^2(x)$. This section compares bootstrap LASSO Thompson sampling and bootstrap ridge Thompson sampling. In Appendix C, we provide a Bayesian way of using LASSO estimation in a Thompson sampling or UCB contextual bandit motivated by [48]. The theoretical analysis of Bayesian LASSO contextual bandit may be proven more straightforward than the theoretical analysis of bootstrap LASSO contextual bandit, though we leave this analysis for future work.

Consider a simulation setting where the contexts $x_t$ are $d$-dimensional and $x_t \sim \mathcal{N}(0_d, I_d)$. There are three arm $\mathcal{A} = \{0, 1, 2\}$ and the potential outcomes are generated as $r_t(0) = x_{t,0} + f(x_t) + \epsilon_t$, $r_t(1) = 1 - x_{t,0} + f(x_t) + \epsilon_t$, and $r_t(2) = x_{t,1} + f(x_t) + \epsilon_t$, where $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ and $f(x_t)$ is a function that shifts the potential outcomes of all arms. Therefore, in this design, only contextual variables $x_{t,0}$ and $x_{t,1}$ are relevant to the assignment model. Among the remaining $d - 2$ contextual variables, there are $q$ nuisance contextual variables that appear in the outcome model, but should not play a role in the assignment model, as their effect to the potential outcomes is the same for all arms. These nuisance contextual variables shift the potential outcomes by $f(x_t) = 2\sigma \sum_{j=2}^{q+1} x_{t,j}$. The remaining $d - q - 2$ contextual variables are noise contextual variables and do not play a role in either the assignment or the outcome model.

We choose $d = 102$ and $q = 51$. So, among the 100 non-relevant contextual variables, 51 are nuisance contextual variables and 49 are noise contextual variables. First, we compare LASSO and ridge in terms of fit on 1000 observations, when the assignment of contexts to arms is purely randomized. The number of nuisance contextual variables, $q$, is chosen so that LASSO and ridge perform equivalently in terms of fit on purely randomized data, in order to evaluate the choice of the estimation method in the adaptive setting all else being equal.

22

Indeed, as demonstrated in Figure 7, the LASSO and ridge models are almost identical in terms of mean squared error (MSE) when trained on batches of randomized data.
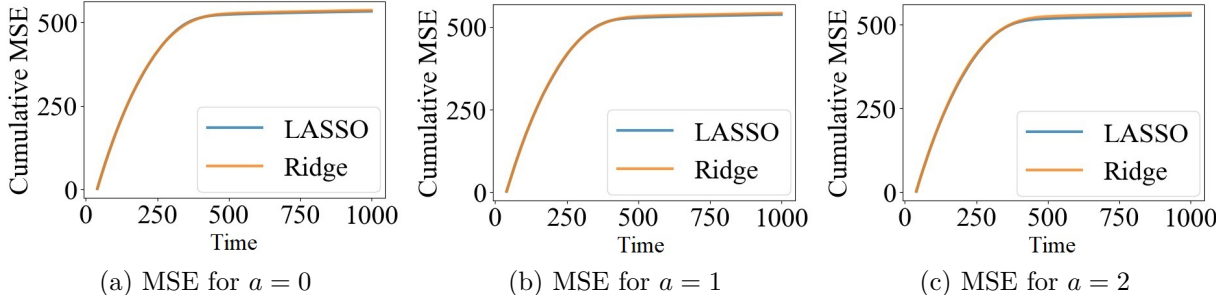


(a) MSE for $a = 0$      (b) MSE for $a = 1$      (c) MSE for $a = 2$

Figure 7: Cumulative MSE averaged over hundreds of simulations for LASSO and ridge of all arms' outcome model estimation on 1000 randomized observations.

Subsequently, we compare the performance of LASSO and ridge in the bandit setting with learning horizon 1000 context. Despite the initial equivalence of LASSO and ridge in the training setting, LASSO clearly outperforms ridge in terms of regret in the adaptive learning setting, as shown in Figure 8.



Figure 8: Cumulative regret averaged over hundreds of simulations for LASSO Thompson sampling and ridge Thompson sampling on 1000 learning observations.

The contributing factor to the bandit performance dissimilarity of these seemingly equivalent models on randomized training data is the presence of confounding. In the initial learning observations, a ridge bandit, due to $L_2$ regularization, brings in all of the nuisance and noise contextual variables, as shown in Figures 9a, 9c. The nuisance contextual variables affect assignment (possibly in non-linear ways) and act as confounders. There is insufficient data in the early observations to accurately control for all of them in the outcome model estimation. Both the nuisance and the noise contextual variables create more extreme and variable assignment probabilities, increasing the variance of estimation. A LASSO bandit, due to the $L_1$ regularization, excludes from the outcome model most of the noise contextual

(a) Ridge Weak Coefficients

(b) LASSO Weak Coefficients

(c) Ridge Noise Coefficients
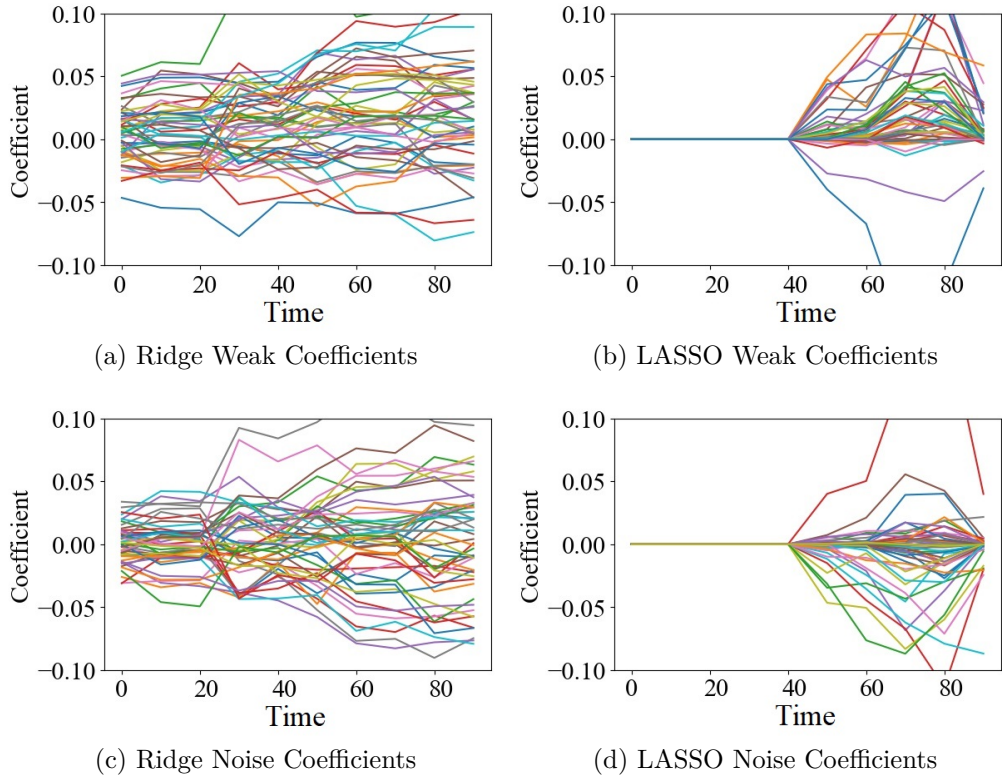
(d) LASSO Noise Coefficients

Figure 9: Coefficient paths of nuisance and noise contextual variables of the first arm's outcome model in the first 100 learning observations for ridge Thompson sampling and LASSO Thompson sampling.

variables and initially, the nuisance contextual variables, as shown in Figures 9b, 9d. As a result, in the early observations, there are fewer confounders compared to a ridge bandit. Therefore, in the subsequent stages of learning, there is less bias as well as less noise in the assignment process for the LASSO bandit than for the ridge bandit.

## 4.2 Significance of Non-Parametric Bandits under Non-Linearities

Real-world applications, such as recommendation systems, may have inherently difficult and complex outcome models. In such cases, contextual bandits based on non-parametric model estimation may be more agile. In this section, we compare non-parametric contextual bandits (generalized random forest Thompson sampling) with parametric contextual bandits that make linearity assumptions on the potential outcome models (bootstrap LASSO Thompson sampling and bootstrap ridge Thompson sampling).

Consider a simple non-linear simulation design with 10-dimensional contexts $x_t$ such that $x_{t,j} \sim \mathcal{N}(0,1)$, $j = 0, \ldots, 9$. There are three arms $\mathcal{A} = \{0,1,2\}$ and the potential outcomes are generated as $r_t(0) = x_{t,0} + \epsilon_t$, $r_t(1) = -x_{t,0} + \epsilon_t$, and $r_t(2) = \mathbf{1}\{x_{t,1} <$

$0\}\left(\min\{x_{t,0}, -x_{t,0}\} - 1\right) + \mathbf{1}\{x_{t,1} > 0\}\left(\max\{x_{t,0}, -x_{t,0}\} + 1\right) + \epsilon_t$, where $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ with $\sigma = 0.1$. Therefore, only contextual variables $x_{t,0}$ and $x_{t,1}$ are relevant to the arm assignment model. In this design, the correct assignment is $a = 2$ in the first and second quadrants, $a = 1$ in the third quadrant and $a = 0$ in the fourth quadrant. We run the bandits on 5000 observations. The models of LASSO and ridge include quadratic and second order interarm terms, there are $B = 100$ bootstrap samples, and the regularization parameter is chosen via cross-validation. The number of trees in the generalized random forest is $m = 200$ and the tuning parameters are the default, specified in [9] and in the `grf` R package.

Figure 10 shows that the generalized random forest bandit outperforms the LASSO and the ridge bandits. In cases where the outcome functional form is complicated, which is expected in real-world settings, bandits based on non-parametric model estimation may be proven useful and perform better. However, one needs to bear in mind that similarly to the ridge bandit, the generalized random forest bandit creates a complex assignment model and is subject to the disadvantages discussed in Section 4.1. In this simulation design, the presence of noise contextual variables leads the LASSO bandit to outperform the ridge bandit. But the inability of the LASSO bandit and the ability of the generalized random forest bandit to fit the potential outcome model of the third arm, results in a strong performance edge of the latter. Figure 11 shows the assignment evolution in the $(x_{\cdot,0}, x_{\cdot,1})$ for the ridge, the LASSO and the generalized random forest bandit. The generalized random forest bandit has the advantage that the outcome model is non-parametric, and thus is able to account for non-linear functions of the contextual variables, in principle reducing problems that might arise if the assignment model is a non-linear function of contextual variables. Additionally, the "honest" estimation property featured by the generalized random forest bandit reduces bias.
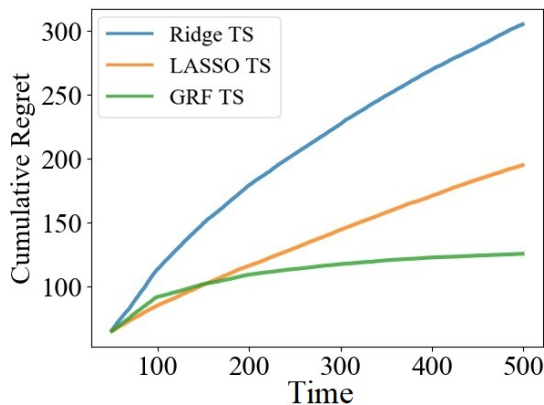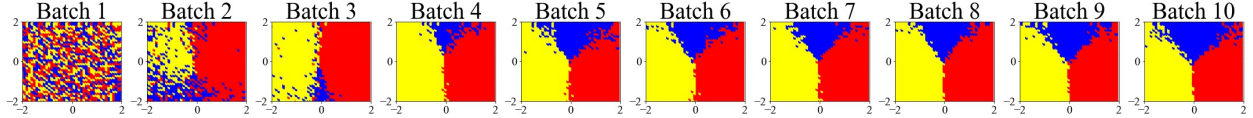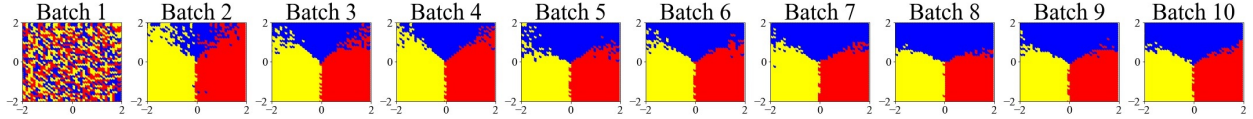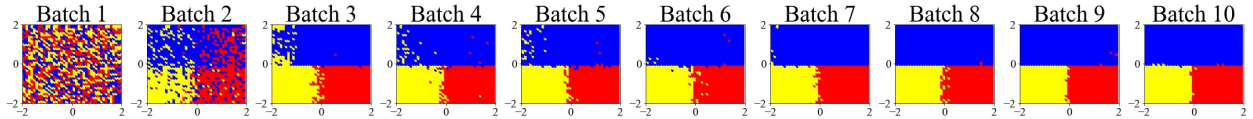


Figure 10: Cumulative regret averaged over hundreds of simulations of ridge, LASSO and generalized random forest Thompson sampling over 5000.

Figure 11: Evolution of the arm assignment in the $(x_{.,0}, x_{.,1})$ context space for ridge, LASSO and generalized random forest Thompson sampling. The optimal assignment is $a = 2$ (blue) in the 1st and 2nd quadrants, $a = 1$ (yellow) in the 3rd quadrant and $a = 0$ (red) in the 4th quadrant.

## 4.3 Smoothing the Assignment Policy

So far, we have considered the assignment rule for each context distinctly; the rule depends on the mean and variance of estimates at each context $x$. A literature on optimal policy evaluation in the offline world (e.g. [10]) derives efficient methods for offline policy estimation when the policy is constrained to be of limited complexity. One example uses trees of limited depth as the relevant policy class. The method first constructs the efficient score $\hat{\mu}_a(x) + \frac{r - \hat{\mu}_a(x)}{\hat{p}_a(x)}$ for each observation $(x, a, r)$. Subsequently, it estimates a classification tree on the history of contexts $\mathbf{X}$ and assignments $\mathbf{a}$, weighted by the absolute value of the efficient scores in order to determine the optimal choice of arm in each leaf, where leaves are regions of the context space.

Here, we propose to follow their method to estimate a policy assignment tree. However, we use the output differently: rather than deterministically assigning each context to the estimated optimal policy, instead we use estimates of the mean and variance of each arm within each leaf together with Thompson sampling or UCB to determine assignments. One complication that potentially arises in the online setting is that in the early stages of learning the estimation of the "nuisance parameters" in the efficient score, $\hat{\mu}_a(x)$ and $\hat{p}_a(x)$, may be noisy due to the small number of observations.

To understand why using simpler assignment rules through a form of smoothing can be beneficial, it is useful to contrast two cases, one where the probability that an arm is

26

best is estimated very precisely, and the second where we have a noisy estimate of that probability. Suppose that sampling according to the true probability balances exploration and exploitation in an ideal way (e.g. that the Thompson sampling heuristic is ideal in a given setting). Then, when shifting to the second case where the probabilities are unknown, adding a small amount of smoothing to the assignment rule will have little effect on the exploitation side of the bandit trade-off (given that the estimates were noisy, a little smoothing does not introduce first-order mistakes in assignment). From a practical perspective, this suggests using simpler assignment rules can improve performance; further, the Thompson sampling heuristic, which incorporates its own form of smoothing by randomizing assignment, may also improve performance over UCB. However, smoothing improves the exploration side of the trade-off, since it enables lower-variance estimation in future batches. Note that simple assignment rules may have other advantages; for example, [42] highlights the advantages of simplicity for interpretability in health applications of contextual bandits.

# 5  Closing Remarks

Contextual bandits are poised to play an important role in a wide range of applications: In these settings, there are many potential sources of bias in estimation of outcome models, not only due to the inherent adaptive data collection, but also due to mismatch between the true data generating process and the outcome model assumptions, and due to prejudice in the training data in form of under-representation or over-representation of certain regions of the context space. content recommendation in web-services, where the learner wants to personalize recommendations (arm) to the profile of a user (context) to maximize engagement (reward); online education platforms, where the learner wants to select a teaching method (arm) based on the characteristics of a student (context) in order to maximize the student's scores (reward); and survey experiments, where the learner wants to learn what information or persuasion (arm) influences the responses (reward) of subjects as a function of their demographics, political beliefs, or other characteristics (context). In these settings, there are many potential sources of bias in estimation of outcome models, not only due to the inherent adaptive data collection, but also due to mismatch between the true data generating process and the outcome model assumptions, and prejudice in the training data in form of under-representation or over-representation of certain regions of the context space in the cold-start training data. To reduce bias, we proposed new contextual bandit designs which integrate balancing methods from the causal inference literature and parametric and non-parametric model estimation methods with optimism or randomization based exploration methods. We provided the first regret bound analysis for linear contextual bandits with

balancing that matches the theoretical guarantees of the linear contextual bandits with direct model estimation We showed that contextual bandit designs with randomization based exploration and balancing in model estimation are more robust to the aforementioned sources of bias and have improved learning rates. We also showed that using simpler, less variable assignment policies, reduces variance in estimation and bias due to confounding and decreases regret. Through a range of simulations and experiments on real-world datasets, we aimed to highlight key tradeoffs and thus provide methodological guidelines to anyone who wishes to efficiently use contextual bandits for learning personalized policies in practice.

# 6    Acknowledgments

# References

[1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *NIPS*, 2011.

[2] Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*, volume 55. Courier Corporation, 1964.

[3] A. Agarwal, S. Bird, M. Cozowicz, L. Hoang, J. Langford, S. Lee, J. Li, D. Melamed, G. Oshri, O. Ribas, S. Sen, and A. Slivkins. Making contextual decisions with low technical debt. 2017.

[4] A. Agarwal, D. Hsu, S. Kale, J. Langford, L. Li, and R. Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. *ICML*, 2014.

[5] S. Agrawal and N. Goyal. Analysis of thompson sampling for the multi-armed bandit problem. *Journal of Machine Learning Research Workshop and Conference Proceedings*, 2012.

[6] S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. *ICML*, 2013.

[7] S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 2016.

[8] S. Athey, G. Imbens, and S. Wager. Approximate residual balancing. *arXiv*, 2017.

[9] S. Athey, J. Tibshirani, and S. Wager. Generalized random forests. *arXiv*, 2017.

[10] S. Athey and S. Wager. Efficient policy learning. *arXiv*, 2017.

[11] P. Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 2003.

[12] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 2002.

[13] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov), 2002.

[14] E. Bareinboim, A. Forney, and J. Pearl. Bandits with unobserved confounders: A causal approach. *ICML*, 2015.

[15] E. Bareinboim and J. Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 2015.

[16] H. Bastani and M. Bayati. Online decision-making with high-dimensional covariates. 2015.

[17] Hamsa Bastani and Mohsen Bayati. Online decision-making with high-dimensional covariates. 2015.

[18] Alberto Bietti, Alekh Agarwal, and John Langford. A contextual bandit bake-off. *arXiv preprint arXiv:1802.04064*, 2018.

[19] L. Breiman. Random forests. *Machine Learning*, 2001.

[20] S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 2012.

[21] Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and non-stochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 2012.

[22] O. Chapelle and L. Li. An empirical evaluation of thompson sampling. *NIPS*, 2011.

[23] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *AISTATS*, 2011.

[24] C. Cortes, Y. Mansour, and M. Mohri. Learning bounds for importance eeighting. *NIPS*, 2010.

[25] Richard K Crump, V Joseph Hotz, Guido W Imbens, and Oscar A Mitnik. Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, 96(1):187–199, 2009.

[26] Y. Deshpande, L. Mackey, V. Syrgkanis, and M. Taddy. Accurate inference in adaptive linear models. *arXiv*, 2017.

[27] M. Dudik, D. Erhan, J. Langford, and L. Li. Doubly robust policy evaluation and optimization. *Statistical Science*, 2014.

[28] M. Dudik, J. Langford, and L. Li. Doubly robust policy evaluation and learning. *ICML*, 2011.

[29] A. Elmachtoub, R. McNellis, S. Oh, and M. Petrik. A practical method for solving contextual bandit problems using decision trees. *arXiv*, 2017.

[30] R. Feraud, R. Allesiardo, T. Urvoy, and F. Clerot. Random forest for the contextual bandit problem. *AISTATS*, 2016.

[31] A. Forney, J. Pearl, and E. Bareinboim. Counterfactual data-fusion for online reinforcement learners. *ICML*, 2017.

[32] A. Goldenshluger and A. Zeevi. A linear response bandit problem. *Stochastic Systems*, 2013.

[33] A. Hoerl and R. Kennard. Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 1970.

[34] Roger A Horn, Roger A Horn, and Charles R Johnson. *Matrix analysis*. Cambridge university press, 1990.

[35] J. Huang, A. Gretton, K. M. Borgwardt, B. Scholkopf, and A. J. Smola. Correcting sample selection bias by unlabeled data. *NIPS*, 2007.

[36] G. W. Imbens and D. B. Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*. 2015.

[37] N. Jiang and L. Li. Doubly robust off-policy value evaluation for reinforcement learning. *ICML*, 2016.

[38] N. Kallus. Balanced policy evaluation and learning. *arXiv*, 2017.

[39] Nathan Kallus and Angela Zhou. Policy evaluation and optimization with continuous treatments. *AISTATS*, 2018.

[40] T. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 1985.

[41] F. Lattimore, T. Lattimore, and M. D. Reid. Causal bandits: Learning good interventions via causal inference. *NIPS*, 2016.

[42] H. Lei, A. Tewari, and S. Murphy. An actor-critic contextual bandit algorithm for personalized mobile health interventions. *arXiv*, 2017.

[43] L. Li, S. Chen, J. Kleban, and A. Gupta. Counterfactual estimation and optimization of click metrics for search engines. *CoRR*, 2014.

[44] L. Li, W. Chu, J. Langford, T. Moon, and X. Wang. An unbiased offline evaluation of contextual bandit algorithms with generalized linear models. *Journal of Machine Learning Research Workshop and Conference Proceedings*, 2012.

[45] L. Li, W. Chu, J. Langford, and R. Schapire. A contextual-bandit approach to personalized news article recommendation. *WWW*, 2010.

[46] L. Li, Y. Lu, and D. Zhou. Provably optimal algorithms for generalized linear contextual bandits. *ICML*, 2017.

[47] X. Nie, X. Tian, J. Taylor, and J. Zou. Why adaptively collected data have negative bias and how to correct for it. *AISTATS*, 2018.

[48] Trevor Park and George Casella. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.

[49] V. Perchet and P. Rigollet. The multi-armed bandit problem with covariates. *The Annals of Statistics*, 2013.

[50] P. Rigollet and R. Zeevi. Nonparametric bandits with covariates. *COLT*, 2010.

[51] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 2014.

[52] D. Russo, B. Van Roy, A. Kazerouni, I. Osband, and Z. Wen. A tutorial on Thompson sampling. *arXiv*, 2017.

[53] Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.

[54] Daniel O Scharfstein, Andrea Rotnitzky, and James M Robins. Adjusting for nonignorable drop-out using semiparametric nonresponse models. *Journal of the American Statistical Association*, 94(448):1096–1120, 1999.

[55] S. Scott. A modern bayesian look at the multi-armed bandit. *Applied Stochastic Models in Business and Industry*, 2010.

[56] A. Slivkins. Contextual bandits with similarity information. *Journal of Machine Learning Research*, 2014.

[57] A. Strehl, J. Langford, L. Li, and S. Kakade. Learning from logged implicit exploration data. *NIPS*, 2010.

[58] Alex Strehl, John Langford, Lihong Li, and Sham M Kakade. Learning from logged implicit exploration data. In *NIPS*, 2010.

[59] A. Swaminathan and T. Joachims. Batch learning from logged bandit feedback through counterfactual risk minimization. *Journal of Machine Learning Research*, 2015.

[60] P. Thomas and E. Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. *ICML*, 2016.

[61] W. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 1933.

[62] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society*, 1996.

[63] S. Villar, J. Bowden, and J. Wason. Multi-armed bandit models for the optimal design of clinical trials: benefits and challenges. *Statistical Science*, 2015.

[64] Y. X. Wang, A. Agarwal, and M. Dudik. Optimal and adaptive off-policy evaluation in contextual bandits. *ICML*, 2017.

[65] B. Zadrozny. Learning and evaluating classifiers under sample selection bias. *ICML*, 2004.

[66] Zhengyuan Zhou, Susan Athey, and Stefan Wager. Offline multi-action policy learning: Generalization and optimization. *arXiv preprint arXiv:1810.04778*, 2018.

# A    Appendix A: Regret Bound Proofs

## A.1    Auxiliary Results

We collect in this section all the existing results in the literature for later use.

The first one is the self-normalized bound for vector-valued martingales in [1].

**Lemma 1.** *Let $\{F_t\}_{t=0}^{\infty}$ be a filtration. Let $\{\eta_t\}_{t=1}^{\infty}$ be a real-valued stochastic process such that $\eta_t$ is $F_t$ measurable and $\eta_t$ is conditionally $R$-sub-Guassian for some $R \geq 0$: $\mathbf{E}[e^{\lambda \eta_t} \mid F_{t-1}] \leq \exp(\frac{\lambda^2 R^2}{2})$. Let $\{Z_t\}_{t=1}^{\infty}$ be an $\mathbb{R}^d$-valued stochastic process such that $Z_t$ is $F_{t-1}$ measurable. Assume that $V$ is a $d \times d$ positive definite matrix, and for any integer $t \geq 0$, define:*

$$V_t = V + \sum_{s=1}^{t} Z_s Z_s^T, S_t = \sum_{s=1}^{t} \eta_s Z_s.$$

*Then, for any $\delta > 0$, with probability at least $1 - \delta$, we have:*

$$(1) \qquad S_t^T V_t^{-1} S_t \leq 2R^2 \log\left(\frac{(\det(V_t))^{\frac{1}{2}}(\det(V))^{\frac{1}{2}}}{\delta}\right).$$

The second one, taken from [2], gives large deviation bounds for Guassian random variables.

**Lemma 2.** *Let $Z$ be a Guassian random variable with mean $m$ and variance $\sigma^2$. Then for any real number $r \geq 1$,*

$$(2) \qquad \frac{1}{2\sqrt{\pi}r}e^{-\frac{r^2}{2}} \leq \mathbb{P}(|Z - m| > r\sigma) \leq \frac{1}{\sqrt{\pi}r}e^{-\frac{r^2}{2}}.$$

The third one is taken from [13].

**Lemma 3.** *Let $Z'$ and $Z$ be two symmetric and positive semidefinite $d \times d$ matrices where the corresponding eigenvalues are $\lambda'_1, \ldots, \lambda'_d$ and $\lambda_1, \ldots, \lambda_d$ respectively. If $Z' = Z + xx^T$, then the eigenvalues can be arranged in such a way that $\lambda_j \leq \lambda'_j$ and $x^T Z^{-1} x \leq 10 \sum_{j=1}^{d} \frac{\lambda'_j - \lambda_j}{\lambda_j}$.*

The next one is a basic fact from matrix algebra ( [34]).

**Lemma 4.** *Let $M$ be a symmetric, positive definite matrix, and $x, y$ be vectors (all with appropriate dimensions). Then the weighted inner product $\langle \cdot, \cdot \rangle_M$ is defined as $\langle x, y \rangle_M =$*

$x^T M y$. *Furthermore,* $\sqrt{\langle x, x \rangle_M}$ *is a norm on* $x$, *which we denote by* $\|x\|_M$ *and we have* $|\langle x, y \rangle_M| \leq \|x\|_M \|y\|_M$.

## A.2  Proof of Theorem 1

In this section, we provide the proof to the regret bound given in Theorem 1. We start with BLTS, which is more involved. For ease of exposition, we break the proof into several steps, each of which will be explained. We start by setting up some notation. Let $\hat{\theta}_a(t)$ and $\tilde{\theta}_a(t)$ denote the estimated mean and the sampled mean for arm $a$ in the BLTS algorithm at time $t$, respectively. By some algebra, one can show that $(r_a(t) - X_a^T \hat{\theta}_a(t))^T W_a (r_a(t) - X_a^T \hat{\theta}_a(t)) = \sigma_a^2 \frac{\sum_{i=1}^t w_i}{t} \in [\sigma_a^2 \gamma, \sigma_a^2 \frac{1}{\gamma}]$. Consequently, $B_a(t)$ is only a constant factor of the variance $\mathbf{V}(\hat{\theta}_a(t))$ term and it suffices to focus on $B_a(t)$. Let $w_t$ denote the thresholded inverse propensity score computed at $t$. For each arm $a$, an equivalent way of writing the updates in the BLTS algorithm is:

$$B_a(t+1) = \begin{cases} B_a(t) + w_t x_{t+1} x_{t+1}^T, & \textbf{if } a \text{ is selected in } t \\ B_a(t), & \textbf{otherwise}, \end{cases}$$

$$\hat{\theta}_a(t+1) = \begin{cases} B_a(t+1)^{-1} \sum_{s \in \mathcal{S}_a(t)} w_s x_s r_a(s), & \textbf{if } a \text{ is selected in } t \\ \hat{\theta}_a(t), & \textbf{otherwise}, \end{cases}$$

where $\mathcal{S}_a(t) = \{1 \leq s \leq t \mid a(s) = a\}$ keeps track of all the iterations where action $a$ is taken. Note that an equivalent way of expressing $B_a(t)$ that will be used later is $B_a(t) = \lambda \mathbf{I} + \sum_{s \in \mathcal{S}_a(t)} \sqrt{w_s} x_s (\sqrt{w_s} x_s)^T$. We shall freely use and switch between the update written in this incremental format and the update format given in the main text.

*Step 1: High concentration bound of estimated means.*

Here we show that $x_t^T \hat{\theta}_a(t)$ concentrates around $x_t^T \theta_a$. Furthermore, such concentration is uniform across time $t$ and different arms. Specifically, define the event $C = \bigcap_{t=1}^T \bigcap_{a \in \mathcal{A}} C_a(t)$, where $C_a(t)$ is the following event:

$$x_t^T \hat{\theta}_a(t) \in [x_t^T \theta_a - O\left(\sqrt{d \log(\frac{KT}{\delta})}\right) \sqrt{x_t^T (B_a(t))^{-1} x_t}, \ x_t^T \theta_a + O\left(\sqrt{d \log(\frac{KT}{\delta})}\right) \sqrt{x_t^T (B_a(t))^{-1} x_t}].$$

The event $C_a(t)$ essentially says that $x_t^T \hat{\theta}_a(t)$ is within some constant multiple of standard deviations from $x_t^T \theta_a$ (the true mean reward). Note that since $\tilde{\theta}$ is drawn from a Guassian with mean $\hat{\theta}$ and variance $\frac{\log \frac{1}{\delta}}{\epsilon}(B_a(t))^{-1}$, the quantity $\sqrt{x_t^T (B_a(t))^{-1} x_t}$ is the standard deviation (up to some constant multiple) of $x_t^T \tilde{\theta}_a$, which is the sample reward associated with arm $a$ estimated by the algorithm BLTS. Our goal in this step is to show that the event $C$ occurs

with high probability. Intuitively, this means that the estimated parameter $\hat{\theta}_a$ is concentrated around the true parameter $\theta_a$: this is an important ingredient for low regret, because if the estimated mean parameters were wrong, then the model maintained by the algorithm would be incorrect, in which case it is impossible to attain any good performance. Specifically, we establish that $\mathbf{P}(C) \geq 1 - \frac{\delta}{T}$, $x_t^T \hat{\theta}_a(t) \in [x_t^T \theta_a - O\left(\sqrt{d \log(\frac{KT}{\delta})}\right)\sqrt{x_t^T (B_a(t))^{-1} x_t}, x_t^T \theta_a +$
$O\left(\sqrt{d \log(\frac{KT}{\delta})}\right)\sqrt{x_t^T (B_a(t))^{-1} x_t}], \forall t \in \{1, \dots, T\}, \forall a \in \mathcal{A}$. To show this, let

$$
\begin{aligned}
\mathcal{X}_s &= \sqrt{w_s} x_s, \quad \eta_s^a = \sqrt{w_s}(r_a(s) - x_s^T \theta_a) \\
V_t^a &= \lambda \mathbf{I} + \sum_{s \in \mathcal{S}_a(t)} \sqrt{w_s} x_s (\sqrt{w_s} x_s)^T = \lambda \mathbf{I} + \sum_{s \in \mathcal{S}_a(t)} \mathcal{X}_s \mathcal{X}_s^T \\
S_t^a &= \sum_{s \in \mathcal{S}_a(t)} \sqrt{w_s} x_s \eta_s^a = \sum_{s \in \mathcal{S}_a(t)} \eta_s^a \mathcal{X}_s.
\end{aligned}
$$

Since $r_a(s) - x_t^T \theta_a$ is conditionally sub-Gaussian (conditioned on $\mathcal{F}_{t-1}$), and $w_s$ is bounded, $\sqrt{w_s}(r_a(s) - x_t^T \theta_a)$ is also conditionally sub-Gaussian. Let $R^*$ be the universal sub-Gaussian constant for the entire proof. It is easy to check that under the canonical filtration $\mathcal{F}_t$, all the adaptability assumptions in Lemma 1 hold: specifically, $\mathcal{X}_t$ is adapted to $\mathcal{F}_{t-1}$ and $\eta_t$ is adapted to $\mathcal{F}_t$. Consequently, by Lemma 1, we have for each $a \in \mathcal{A}$, with probability at least $1 - \delta$:

$$
(3) \qquad S_t^T V_t^{-1} S_t \leq 2R^2 \log(\frac{(\det(V_t))^{\frac{1}{2}}(\det(\lambda \mathbf{I}))^{-\frac{1}{2}}}{\delta}) = 2R^2 \log(\frac{(\det(V_t))^{\frac{1}{2}}(\lambda)^{-\frac{p}{2}}}{\delta})
$$

$$
(4) \qquad \leq 2R^2 \log(\frac{(\frac{|\mathcal{S}_a(t)|}{\gamma} + \lambda)^{\frac{p}{2}} \lambda^{-\frac{p}{2}}}{\delta}) \leq 2R^2 \log(\frac{(\frac{t}{\gamma} + \lambda)^{\frac{p}{2}} \lambda^{-\frac{p}{2}}}{\delta}) = pR^2 \log(\frac{1 + \frac{t}{\gamma\lambda}}{\delta}),
$$

where the last inequality follows because the largest eigenvalue of the matrix $x_t^T x_t$ can be at most one and hence by a simple calculation, the largest eigenvalue of $V_t^a$ can be at most $\frac{|\mathcal{S}_a(t)|}{\gamma} + \lambda$, which in turn is upper bounded by $\frac{t}{\gamma} + \lambda$. Since the determinant equals the product of all the eigenvalues, and since all the matrices involved are positive-semidefinite, the result follows.

With the above bound in place, we are now ready to bound $x_t^T \hat{\theta}_a(t) - x_t^T \theta_a$. First note that $|x_t^T \hat{\theta}_a(t) - x_t^T \theta_a| = |x_t^T (V_t^a)^{-1}(S_t^a - \theta_a)|$. Since $V_t^a$ is positive definite, its inverse is positive definite as well. Consequently, by Lemma 4, we have:

$$
(5) \qquad |x_t^T \hat{\theta}_a(t) - x_t^T \theta_a| = |x_t^T (V_t^a)^{-1}(S_t^a - \lambda \theta_a)| \leq \|x_t^T\|_{(V_t^a)^{-1}} \|S_t^a - \lambda \theta_a\|_{(V_t^a)^{-1}}
$$

$$
(6) \qquad = \|x_t^T\|_{(B_t^a)^{-1}} \|S_t^a - \lambda \theta_a\|_{(V_t^a)^{-1}} = \sqrt{x_t^T (B_a(t))^{-1} x_t} \|S_t^a - \lambda \theta_a\|_{(V_t^a)^{-1}},
$$

where the first equality follows from the following algebraic calculation:

$$(7) \quad (V_t^a)^{-1}(S_t^a - \theta_a) = (\lambda \mathbf{I} + \sum_{s \in \mathcal{S}_a(t)} \mathcal{X}_s \mathcal{X}_s^T)^{-1}(S_t^a - \lambda \theta_a)$$

$$(8) \quad = (X_a^T W X_a + \lambda \mathbf{I})(X_a^T W_a r_a - X_a^T W_a X_a \theta_a - \lambda \mathbf{I} \theta_a)$$

$$(9) \quad = (X_a^T W X_a + \lambda \mathbf{I})^{-1}(X_a^T W_a r_a) - (X_a^T W X_a + \lambda \mathbf{I})^{-1}(X_a^T W_a X_a \theta_a + \lambda \mathbf{I} \theta_a) = \hat{\theta}_a(t) - \theta_a,$$

and the second inequality follows from Lemma 4. Note that since

$$(10) \quad \|\theta_a\|_{(V_t^a)^{-1}} = \sqrt{\theta_a^T (V_t^a)^{-1} \theta_a} \le \sqrt{\lambda_{\max}[(V_t^a)^{-1}]}\|\theta_a\| \le \sqrt{\lambda_{\max}[(V_t^a)^{-1}]}$$

$$(11) \quad \le \frac{1}{\sqrt{\lambda_{\min}[V_t^a]}} \le \frac{1}{\sqrt{\lambda_{\min}[\lambda \mathbf{I}]}} = \frac{1}{\sqrt{\lambda}},$$

where $\lambda_{\max}$ and $\lambda_{\min}$ denote the maximum eigenvalue and the minimum eigenvalue of a matrix respectively. Consequently, we have that with probability at least $1 - \delta$:

$$(12)$$

$$|x_t^T \hat{\theta}_a(t) - x_t^T \theta_a| \le \sqrt{x_t^T (B_a(t))^{-1} x_t} \|S_t^a - \lambda \theta_a\|_{(V_t^a)^{-1}}$$

$$(13)$$

$$\le \sqrt{x_t^T (B_a(t))^{-1} x_t}\left(\|S_t^a\|_{(V_t^a)^{-1}} + \lambda\|\theta_a\|_{(V_t^a)^{-1}}\right)$$

$$(14)$$

$$\le \sqrt{x_t^T (B_a(t))^{-1} x_t}\left(\|S_t^a\|_{(V_t^a)^{-1}} + \lambda\frac{1}{\sqrt{\lambda}}\right) \le \sqrt{x_t^T (B_a(t))^{-1} x_t}\left(\sqrt{d}R\sqrt{\log(\frac{1 + \frac{t}{\gamma\lambda}}{\delta})} + \sqrt{\lambda}\right)$$

$$(15)$$

$$\le \sqrt{x_t^T (B_a(t))^{-1} x_t}\left(\sqrt{d}R\sqrt{\log(\frac{1 + \frac{T}{\gamma\lambda}}{\delta})} + \sqrt{\lambda}\right)$$

$$(16)$$

where the second inequality follows from the triangle inequality of a norm, the third inequality follows from Equation (10) and the last inequality follows from Equation (3). Now take $\delta$ to be $\frac{\delta}{KT^2}$, and absorbing all the constants into the big-$O$, we have, with probability at least

37

$1 - \frac{\delta}{KT^2}$,

$$\mathbf{P}(C_a(t)) \leq \sqrt{x_t^T (B_a(t))^{-1} x_t} \left( \sqrt{d} R \sqrt{\log(\frac{KT^2 + K\frac{T^3}{\gamma\lambda}}{\delta})} + \sqrt{\lambda} \right) = \sqrt{x_t^T (B_a(t))^{-1} x_t} O\left( \sqrt{d \log(\frac{KT}{\delta})} \right).$$

By a union bound, we have: $\mathcal{P}(C) = 1 - \mathbf{P}(\bigcup_{t=1}^T \bigcup_{a \in \mathcal{A}} C_a^c(t)) \leq \sum_{t=1}^T \sum_{a \in \mathcal{A}} \mathbf{P}(C_a^c(t)) \leq \frac{\delta}{T}$. Consequently, we have with probability at least $1 - \frac{\delta}{T}$, $x_t^T \hat{\theta}_a(t) \in [x_t^T \theta_a - O\left( \sqrt{d \log(\frac{KT}{\delta})} \right) \sqrt{x_t^T (B_a(t))^{-1} x_t}$, $x_t^T \theta_a + O\left( \sqrt{d \log(\frac{KT}{\delta})} \right) \sqrt{x_t^T (B_a(t))^{-1} x_t}], \forall t \in \{1, \ldots, T\}, \forall a \in \mathcal{A}$.

*Step 2: High concentration bound of sampled means.*

The above step establishes that $\hat{\theta}_a$ concentrates around the true $\theta_a$. We now establish that $\tilde{\theta}_a$ is concentrates around $\hat{\theta}_a$. This concentration is very much to be expected, because $\tilde{\theta}_a$ is exactly drawn from a Guassian with mean $\hat{\theta}_a$. Consequently, since Guassian random variables concentrate around its mean (Lemma 2), by choosing an appropriate multiple of its standard deviation, we can get the desired concentration probability relatively easily. Specifically, here we have that with probablity at least $1 - \frac{1}{T}$, for all $t = 1, \ldots, T$, and for all $a \in \mathcal{A}$:

$$(17) \qquad x_t^T \tilde{\theta}_a(t) \in [x_t^T \hat{\theta}_a - O\left( \sqrt{d \frac{\log \frac{1}{\delta}}{\epsilon} \log(dKT)} \right) \sqrt{x_t^T (B_a(t))^{-1} x_t},$$

$$(18) \qquad x_t^T \hat{\theta}_a + O\left( \sqrt{d \frac{\log \frac{1}{\delta}}{\epsilon} \log(dKT)} \right) \sqrt{x_t^T (B_a(t))^{-1} x_t}].$$

To see this, denote $E_a(t)$ to be the above event. Then by a straightfoward application of Lemma 2, one can easily check that

$$\mathcal{P}\left( \|B_a(t)^{0.5}\left( \tilde{\theta}_a(t) - \hat{\theta}_a(t) \right)\| \geq \sqrt{\frac{\log \frac{1}{\delta}}{\epsilon}} \sqrt{4d \log(dNT)} \right) \leq \frac{1}{KT^2}.$$

Consequently, by a further union bound across all $a$ and across all $t$, we have with probability at least $1 - \frac{1}{T}$, $\|B_a(t)^{0.5}\left( \tilde{\theta}_a(t) - \hat{\theta}_a(t) \right)\| \leq \sqrt{\frac{\log \frac{1}{\delta}}{\epsilon}} \sqrt{4d \log(dKT)}$. Note that when this holds, we can easily bound $|x_t^T \tilde{\theta}_a(t) - x_t^T \hat{\theta}_a|$ as follows:

$$(19) \qquad |x_t^T \tilde{\theta}_a(t) - x_t^T \hat{\theta}_a| = |x_t^T B_a(t)^{-0.5} B_a(t)^{0.5} (\tilde{\theta}_a(t) - \hat{\theta}_a)|$$

$$(20) \qquad \leq \|x_t^T B_a(t)^{-0.5}\| \|B_a(t)^{0.5} (\tilde{\theta}_a(t) - \hat{\theta}_a)\| \leq \sqrt{x_t^T B_a(t)^{-1} x_t} \sqrt{\frac{\log \frac{1}{\delta}}{\epsilon}} \sqrt{4d \log(dKT)}.$$

Putting the above two pieces together and using $O(\cdot)$ to simplify all the constants, yields that with probablity at least $1-\frac{1}{T}$, $x_t^T \tilde{\theta}_a(t) \in [x_t^T \hat{\theta}_a - O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}, x_t^T\hat{\theta}_a + O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}]$.

*Step 3: Bounding regret in terms of standard deviations.*

The previous two steps combined together establish that the samples $\tilde{\theta}_a(t)$ produced by BLTS are close to the true $\theta_a$. An immediate consequence of this is that we can then bound the instantaneous regret $ir(t)$ at time $t$ in terms of the standard deviations. More specifically, under the above two concentration events, for each arm $a$, $x_t^T \tilde{\theta}_a(t)$ differs from $x_t^T \hat{\theta}_a$ by at most $O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}$ and $x_t^T\hat{\theta}_a(t)$ differs from $x_t^T\theta_a$ by at most $O\left(\sqrt{d\log(\frac{KT}{\delta})}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}]$. Consequently, $x_t^T\tilde{\theta}_a(t)$ differs from $x_t^T\theta_a$ by at most $O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))} + \sqrt{d\log(\frac{KT}{\delta})}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}$. Next, since action $a$ is chosen at time $t$, it must be that $x_t^T\tilde{\theta}_a(t)$ yields the largest value, and in particualr, $x_t^T\tilde{\theta}_a(t) \geq x_t^T\theta_{a^*(t)}$. Putting the above discussion together, we can bound the instantaneous regret as follows:

$$(21) \quad ir(t) = x_t^T(\theta_{a^*(t)} - \theta_{a(t)})$$

$$(22) \quad \leq O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))} + \sqrt{d\log(\frac{KT}{\delta})}\right)\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right),$$

where $a(t)$ is the arm chosen at $t$ and $a^*(t)$ is the optimal arm at $t$. Consequently,

$$(23)$$
$$R(T) = \sum_{t=1}^{T}ir(t) \leq O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))} + \sqrt{d\log(\frac{KT}{\delta})}\right)\sum_{t=1}^{T}\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right)$$

$$(24)$$
$$= O\left(\sqrt{d\frac{\log\frac{1}{\delta}}{\epsilon}\log(dKT))}\right)\sum_{t=1}^{T}\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right)$$

$$(25)$$
$$= \tilde{O}\left(\sqrt{\frac{d}{\epsilon}}\right)\sum_{t=1}^{T}\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right).$$

The rest of the proof can then be completed by bounding the sum in the right-hand side of the above equation. First, following a similar analysis (with differences only in constants) as in [6], one can use a martingale based approach to bound $\sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}$ in terms

of $\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t}$. This is done by dividing arms into two different categories and do a careful analysis of the bound in each case. The final bound on the sum, after dropping all the lower order terms is that

$$(26) \quad \sum_{t=1}^{T}\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right) \leq \sum_{t=1}^{T}(1+O(\sqrt{T^\epsilon}))\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t}.$$

Next, in [23], it is shown that for each arm $a$, define $\mathcal{T}_a = \{1 \leq t \leq T \mid a(t) = a\}$, then the following holds:

$$\sum_{t\in\mathcal{T}_a}\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} \leq 5\sqrt{dN_a(T)\log(N_a(T))},$$

where $N_a(T) = |\mathcal{T}_a|$ is the total number of times arm $a$ is selected (note in particular $\sum_{a\in\mathcal{A}} N_a(T) = T$). Consequently, summing over all $a$, we have:

(27)
$$\sum_{t=1}^{T}\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} = \sum_{a\in\mathcal{A}}\sum_{t\in\mathcal{T}_a}\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} \leq \sum_{a\in\mathcal{A}} 5\sqrt{dN_a(T)\log(N_a(T))} = O(\sqrt{dKT\log T}).$$

Consequently, combining the above inequality with Equation (26), we have:

$$(28) \quad \sum_{t=1}^{T}\left(\sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} + \sqrt{x_t^T(B_{a^*(t)}(t))^{-1}x_t}\right) \leq (1+O(\sqrt{T^\epsilon}))O(\sqrt{pNT\log T})$$

$$(29) \quad = O(\sqrt{dKT^{1+\epsilon}\log T}) = \tilde{O}(\sqrt{dKT^{1+\epsilon}}).$$

Finally, combining the preceding inequality with Equation (2.23), we obtain the final regret bound (note that all the inequalities hold with probability $1 - 2\delta$, since we need the concentration events to hold true):

$$R(T) = \tilde{O}\left(\sqrt{\frac{d}{\epsilon}}\right)\tilde{O}(\sqrt{dKT^{1+\epsilon}}) = \tilde{O}(d\frac{\sqrt{KT^{1+\epsilon}}}{\epsilon}).$$

Next, the proof of regret bound for BLUCB follows closely to that of LinUCB in [23] and LinRel in [13], where one divides the BLUCB into two parts for analysis, the first part is BaseBLUCB (which corresponds to BaseLinUCB) and the second part is SuperBLUCB (which is the same as SuperLinUCB in [23]). Hence, the analysis only needs to be adjusted for BaseBLUCB (and we omit the discussion of SuperBLUCB as it is the same as SuperLinUCB in [23]). We start by noting that parameters in BLUCB follow the same update as BLST,

---

**Algorithm 3** BaseBLUCB at Step $t$

---

1: Inputs: $\alpha > 0$ and $\Phi_t \subset \{1, 2, \ldots, t-1\}$
2: **for** $\tau = 1, \ldots, t-1$ **do**
3:     Estimate $\hat{p}_a(x_\tau)$ and set $w_\tau = \frac{1}{\max(\gamma, \hat{p}_a(x_\tau))}$.
4: **end for**
5: $B_a(t) \leftarrow \lambda\mathbf{I} + \sum_{s \in \mathcal{S}_a(t)} \mathbf{1}_{s \in \Phi_t} \sqrt{w_s}x_s(\sqrt{w_s}x_s)^T$, for each $a$.
6: $\hat{\theta}_a(t) \leftarrow B_a(t)^{-1} \sum_{s \in \mathcal{S}_a(t)} \mathbf{1}_{s \in \Phi_t} w_s x_s r_a(s)$, for each $a$.
7: Observe feature $x_t$.
8: $s_a(t) \leftarrow \alpha\sqrt{x_t^T B_a^{-1}(t) x_t}$, for each $a$.
9: $r_a(t) \rightarrow \hat{\theta}_a^T x_t$, for each $a$.

---

which are written as follows:

$$B_a(t+1) = \begin{cases} B_a(t) + w_t x_{t+1} x_{t+1}^T, & \textbf{if } a \text{ is selected in } t \\ B_a(t), & \textbf{otherwise}, \end{cases}$$

$$\hat{\theta}_a(t+1) = \begin{cases} B_a(t+1)^{-1} \sum_{s \in \mathcal{S}_a(t)} w_s x_s r_a(s), & \textbf{if } a \text{ is selected in } t \\ \hat{\theta}_a(t), & \textbf{otherwise}, \end{cases}$$

where $\mathcal{S}_a(t) = \{1 \leq s \leq t \mid a(s) = a\}$ keeps track of all the iterations where action $a$ is taken.

Consequently, from Step 1 above, and with $\alpha = \sqrt{\log \frac{TK}{\delta}}$, we know that with probability at least $1 - \frac{\delta}{T}$, for each $a \in \mathcal{A}$ and each $t = 1, 2, \ldots, T$:

$$|x_t^T \hat{\theta}_a(t) - x_t^T \theta_a| \leq O\left(\sqrt{\log(\frac{KT}{\delta})}\right)\sqrt{x_t^T(B_a(t))^{-1}x_t}.$$

Furthermore, by Equation 27 in Step 3, we have:

$$(30) \qquad \sum_{t = \in \Phi_{T+1}} \sqrt{x_t^T(B_{a(t)}(t))^{-1}x_t} = O(\sqrt{dK|\Phi_{T+1}|\log|\Phi_{T+1}|}).$$

With these two main ingredients in place, the rest of the proof follows the same steps as in [23] (which in turn follows [13]), giving a $O(\sqrt{TdK\log^3(\frac{\log KT\log(T)}{\delta})})$ regret bound.

# B    Experiments on Multi-Class Classification Datasets

We use 300 multiclass datasets from the Open Media Library (OpenML). The full list of datasets in alphabetical order is:

2dplanes, abalone (183), abalone (720), acute-inflammations, Agrawal1, aids, ailerons, airlines, analcatdata_apnea2, analcatdata_apnea3, analcatdata_asbestos, analcatdata_authorship, analcatdata_challenger, analcatdata_creditscore, analcatdata_dmft, analcatdata_germangss, analcatdata_japansolvent, analcatdata_michiganacc, analcatdata_olympic2000, analcatdata_seropositive, analcatdata_vehicle, analcatdata_vineyard, analcatdata_wildcat, AP_Breast_Kidney, AP_Breast_Omentum, AP_Breast_Ovary, AP_Colon_Kidney, AP_Colon_Lung, AP_Endometrium_Breast, AP_Endometrium_Kidney, AP_Endometrium_Ovary, AP_Endometrium_Prostate, AP_Lung_Kidney, AP_Lung_Uterus, AP_Omentum_Lung, AP_Omentum_Prostate, AP_Omentum_Uterus, AP_Ovary_Kidney, AP_Ovary_Lung, AP_Prostate_Ovary, AP_Uterus_Kidney, ar3, ar4, ar5, ar6, arsenic-male-bladder, arsenic-male-lung, artificial-characters, Australian, autoPrice, balance-scale, banana, baskball, bodyfat, bolts, boston (853), boston (872), boston_corrected, BurkittLymphoma, cal_housing, car, chatfield_4, chscase_adopt, chscase_census2, chscase_census3, chscase_vine2, cmc, codrna, codrnaNorm, collins (478), collins (987), confidence (468), confidence (1015), covertype (180), covertype (293), cpu, cpu_small, delta_ailerons, delta_elevators, desharnais, diabetes, diabetes_numeric, diggle_table_a2, disclosure_x_bias, disclosure_x_noise, dresses-sales, eating, ecoli, electricity, elevators, elusage, energy-efficiency, first-order-theorem-proving, fl2000, flags (285), flags (1012), fri_c0_1000_25, fri_c0_1000_5, fri_c0_100_10, fri_c0_100_50, fri_c0_250_10, fri_c0_250_25, fri_c0_250_5, fri_c0_500_25, fri_c0_500_50, fri_c1_1000_5, fri_c1_1000_50, fri_c1_100_10, fri_c1_100_5, fri_c1_100_50, fri_c1_250_10, fri_c1_250_5, fri_c1_250_50, fri_c1_500_10, fri_c1_500_5, fri_c1_500_50, fri_c2_1000_10, fri_c2_1000_25, fri_c2_1000_5, fri_c2_1000_50, fri_c2_100_10, fri_c2_100_25, fri_c2_100_5, fri_c2_100_50, fri_c2_250_10, fri_c2_250_25, fri_c2_250_5, fri_c2_250_50, fri_c2_500_10, fri_c2_500_25, fri_c2_500_5, fri_c3_1000_10, fri_c3_1000_25, fri_c3_1000_5, fri_c3_1000_50, fri_c3_100_10, fri_c3_100_5, fri_c3_250_10, fri_c3_250_5, fri_c3_500_10, fri_c3_500_25, fri_c3_500_5, fri_c4_1000_10, fri_c4_1000_50, fri_c4_100_10, fri_c4_100_100, fri_c4_100_50, fri_c4_250_25, fri_c4_500_10, fri_c4_500_100, fri_c4_500_50, gina_agnostic, gina_prior2, glass, grub-damage (338), grub-damage (1026), hayes-roth, heart-statlog, houses, humandevel, hutsof99_child_witness, hutsof99_logis, Hyperplane_10_1E-4, iris, JapaneseVowels, jEdit_4.0_4.2, jEdit_4.2_4.3, kc1, kc1-binary, kc2, kin8nm, kr-vs-k, kr-vs-kp, kropt, leaf, letter (6), letter (977), leukemia, lowbwt, lupus, machine_cpu, MagicTelescope, mammography, mc2, meta_all.arff,

meta_ensembles.arff, mfeat-factors, mfeat-fourier, mfeat-karhunen (16), mfeat-karhunen (1020), mfeat-morphological (18), mfeat-morphological (962), mfeat-pixel (20), mfeat-pixel (1022), mfeat-zernike (22), mfeat-zernike (995), monks-problems-3, mu284, musk, mv, mw1, MyIris, newton_hema, no2, nursery (26), nursery (959), optdigits, OVA_Endometrium, OVA_Ovary, OVA_Prostate, page-blocks, pasture (339), pasture (964), pc1, pc2, pc3, pc4, pendigits, PieChart1, PieChart3, PieChart4, PizzaCutter3, plasma_retinol, pm10, pollen, pollution, prnn_cushings, prnn_fglass (952), prnn_fglass (996), puma32H, puma8NH, pwLinear, pyrim, quake (209), quake (772), quake (948), rabe_131, rabe_148, rabe_265, rabe_266, rabe_97, rmftsa_ladata, rmftsa_sleepdata (679), rmftsa_sleepdata (741), rsctc2010_1, rsctc2010_2, satellite_image, satimage, scene, schlvote, SEA(50), SEA(50000), segment, sensory, sleuth_case1102, sleuth_case1201, sleuth_case1202, sleuth_case2002, sleuth_ex1605, sleuth_ex2016 (682), sleuth_ex2016 (862), socmob, sonar, space_ga, spambase, spectrometer (313), spectrometer (754), Stagger1, Stagger2, Stagger3, stock, strikes, sylva_agnostic, sylva_prior, synthetic_control (377), synthetic_control (1004), tae, teachingAssistant, tecator, tic-tac-toe, tr21.wc, tr23.wc, vehicle (54), vehicle (994), vehicle_sensIT, vehicleNorm, vinnie, visualizing_environmental, visualizing_galaxy, visualizing_hamster, visualizing_soil, vowel, waveform-5000, white-clover, wind, wind_correlations, wine, wine_quality, witmer_census_1980, zoo

The datasets vary in number of observations, number of classes and number of features. Table 3 summarizes the characteristics of these benchmark datasets.

| Observations | Datasets |
|---|---|
| $\leq 100$ | 58 |
| $> 100$ and $\leq 1000$ | 152 |
| $> 1000$ and $\leq 10000$ | 57 |
| $> 10000$ | 33 |

| Classes | Count |  | Features | Count |
|---|---|---|---|---|
| 2 | 243 |  | $\leq 10$ | 154 |
| $> 2$ and 10 | 48 |  | $> 10$ and $\leq 100$ | 106 |
| $> 10$ | 9 |  | $> 100$ | 40 |

Table 3: Characteristics of the 300 datasets used for the multiclass classification with bandit feedback experiment.

In the main paper, we focused on the family of contextual bandits with linear realizability assumption and on how to improve model estimation in linear contextual bandits using balancing. Hence, we compared our algorithms, balanced linear Thompson sampling (BLTS)

and balanced linear UCB (BLUCB), with LinTS [6] and LinUCB [45], which are baselines that belong in the family of contextual bandits with linear realizability assumption and have strong theoretical guarantees. Here, apart from BLTS, BLUCB, LinTS and LinUCB, we also evaluate the policy-based ILOVETOCONBANDITS (ILTCB) from [4] that does not estimate a model, but instead it assumes access to an oracle for solving fully supervised cost-sensitive classification problems and achieves the statistically optimal regret guarantee.
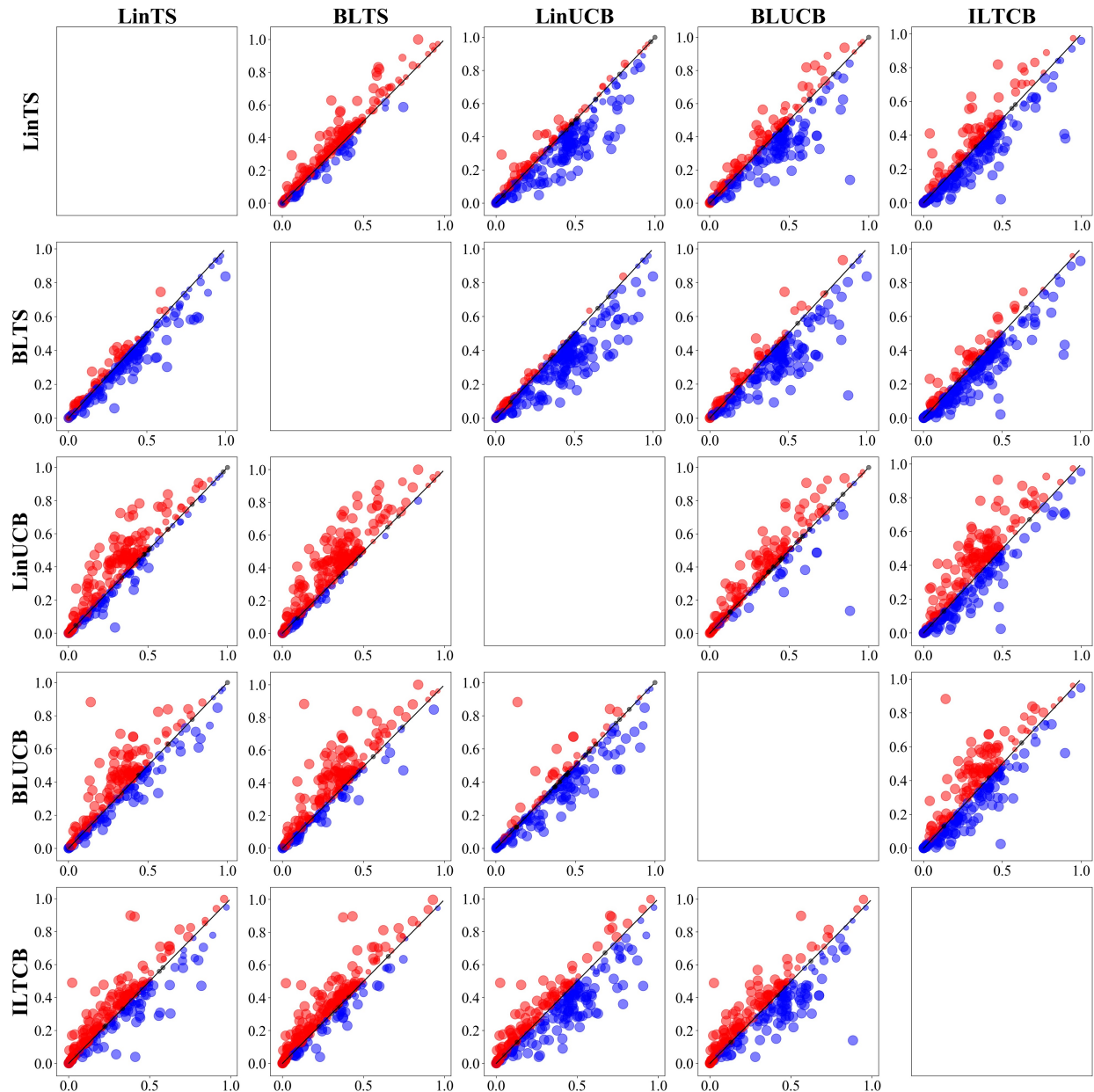


Figure 12: Pairwise comparison of LinTS, BLTS, LinUCB, BLUCB and ILTCB on the 300 classification datasets. BLTS outperforms LinTS, LinUCB, BLUCB and ILTCB.

Figure 12 shows the pairwise comparison of LinTS, BLTS, LinUCB, BLUCB and ILTCB on the 300 classification datasets. Each point corresponds to a dataset. The $x$ coordinate is the normalized cumulative regret of the column bandit and the $y$ coordinate is the normalized cumulative regret of the row bandit. The point is blue when the row bandit has smaller normalized cumulative regret and wins over the column bandit indicated. The point is red when the row bandit loses from the column bandit. The point's size grows with the significance of the win or loss. Table 4 presents the pairwise comparison of LinTS, BLTS, LinUCB, BLUCB and ILTCB on the 300 classification datasets in table-form.

In LinTS, LinUCB, BLTS and BLUCB, the ridge regularization parameter $\lambda$ is chosen via cross-validation every time the model is updated. In LinTS and BLTS, the constant $\alpha$ is optimized among values $0.25, 0.5, 1$, while in LinUCB and BLUCB the constant $\alpha$ is optimized among values $1, 2, 4$. In BLTS and BLUCB, the propensity threshold $\gamma$ is optimized among the values $0.01, 0.05, 0.1, 0.2$. In ILTCB, the parameter $\mu$ of algorithm 1 in [4] is optimized among the values $0.01, 0.1, 1$.

| (down vs. right) | LinTS | BLTS | LinUCB | BLUCB | ILTCB |
|---|---|---|---|---|---|
| LinTS | W=0, L=0 | W=58, L=233 | W=180, L=96 | W=149, L=130 | W=191, L=95 |
| BLTS | W=233, L=58 | W=0, L=0 | W=232, L=55 | W=191, L=92 | W=222, L=60 |
| LinUCB | W=96, L=180 | W=55, L=232 | W=0, L=0 | W=56, L=185 | W=135, L=153 |
| BLUCB | W=130, L=149 | W=92, L=191 | W=185, L=56 | W=0, L=0 | W=153, L=129 |
| ILTCB | W=95, L=191 | W=60, L=222 | W=153, L=135 | W=129, L=153 | W=0, L=0 |

Table 4: Number of the 300 classification datasets in which the contextual bandit algorithm of the row name wins over (W) or loses from (L) the contextual bandit algorithm of the column name. The $300 - W - L$ remaining datasets are ties.

# C    Bayesian LASSO Contextual Bandit

We now provide a Bayesian way of using LASSO estimation in a Thompson sampling contextual bandit inspired by [48]. The theoretical analysis of Bayesian LASSO contextual bandit may be proven more straightforward than the theoretical analysis of bootstrap LASSO contextual bandit, though we leave this analysis for future work.

We model the reward as a linear function of the context $r = x^T\theta + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$, where $\theta$ is $p$-dimensional and sparse. LASSO estimates can be interpreted as posterior mode estimates when the coefficients have **iid** Laplace priors [62]. [48] propose an expanded hierarchy with conjugate normal priors for the regression parameters and independent exponential priors on their variances to perform Gibbs sampling on this posterior. Algorithm 4 proposes a contextual bandit that uses the Gibbs sampler hierarchy of [48].

---

**Algorithm 4** Bayesian LASSO Thompson sampling

1: **Input:** Regularization parameter $\lambda > 0$
2: Set $\mathbf{X}_a \leftarrow$ empty matrix, $\mathbf{r}_a \leftarrow$ empty vector $\forall a \in \mathcal{A}$
3: Sample **iid** $\tau^2_{a,1}, \ldots, \tau^2_{a,p} \sim \prod^p_{j=1} \frac{\lambda^2}{2} e^{-\frac{\lambda^2 \tau^2_j}{2}} d\tau^2_j \ \forall a \in \mathcal{A}$
4: Set $\mathbf{D}_{\tau_a} \leftarrow \text{diag}(\tau^2_{a,1}, \ldots, \tau^2_{a,p}) \ \forall a \in \mathcal{A}$
5: Sample $\theta_{a,1} \sim \mathcal{N}(0, \sigma^2 \mathbf{D}_{\tau_a})$
6: **for** $t = 1, 2, \ldots, T$ **do**
7:      Observe $x_t$
8:      Select $a \leftarrow \arg\max_{a \in \mathcal{A}} x_t^T \theta_{a,t}$
9:      Observe reward $r_t(a)$.
10:      $\mathbf{X}_a \leftarrow [\mathbf{X}_a : x_t^T]$
11:      $\mathbf{r}_a \leftarrow [\mathbf{r}_a : r_t(a)]$
12:      **for** $k = 1, \ldots, K$ Gibbs sampling iterations **do**
13:          **if** $k = 1$ **then**
14:              $\theta^k_{a,t+1} \leftarrow \theta_{a,t}$
15:          **end if**
16:          Sample $\left(\tau^k_{a,j}\right)^2 \sim \text{InverseGaussian}\left(\mu' = \sqrt{\frac{\lambda^2 \sigma^2}{\left(\theta^k_{a,t+1,j}\right)^2}}, \lambda' = \lambda^2\right) \ \forall j = 1, \ldots, p$
17:          Set $\mathbf{D}_{\tau^k_a} \leftarrow \text{diag}\left(\left(\tau^k_{a,1}\right)^2, \ldots, \left(\tau^k_{a,p}\right)^2\right)$ and $\mathbf{A}_a = \mathbf{X}_a^T \mathbf{X}_a + \mathbf{D}_{\tau^k_a}$
18:          Sample $\theta^k_{a,t+1} \sim \mathcal{N}(\mathbf{A}_a^{-1}\mathbf{X}_a^T\mathbf{r}_a, \sigma^2 \mathbf{A}_a^{-1})$
19:      **end for**
20:      Sample $\theta_{a,t+1} \sim \left(\theta^1_{a,t+1}, \ldots, \theta^K_{a,t+1}\right)$.
21:      Set $\theta_{a',t+1} \leftarrow \theta_{a',t} \ \forall a' \in \mathcal{A}\backslash a$
22: **end for**

---