# Stochastic Conservative Contextual Linear Bandits

Jiabin Lin, Xian Yeow Lee, Talukder Jubery, Shana Moothedath,
Soumik Sarkar, and Baskar Ganapathysubramanian

*Abstract*— Many physical systems have underlying safety considerations that require that the strategy deployed ensures the satisfaction of a set of constraints. Further, often we have only partial information on the state of the system. We study the problem of safe real-time decision making under uncertainty. In this paper, we formulate a conservative stochastic contextual bandit formulation for real-time decision making when an adversary chooses a distribution on the set of possible contexts and the learner is subject to certain safety/performance constraints. The learner observes only the context distribution and the exact context is unknown, and the goal is to develop an algorithm that selects a sequence of optimal actions to maximize the cumulative reward without violating the safety constraints at any time step. By leveraging the UCB algorithm for this setting, we propose a conservative linear UCB algorithm for stochastic bandits with context distribution. We prove an upper bound on the regret of the algorithm and show that it can be decomposed into three terms: (i) an upper bound for the regret of the standard linear UCB algorithm, (ii) a constant term (independent of time horizon) that accounts for the loss of being conservative in order to satisfy the safety constraint, and (ii) a constant term (independent of time horizon) that accounts for the loss for the contexts being unknown and only the distrbution being known. To validate the performance of our approach we perform extensive simulations on synthetic data and on real-world maize data collected through the Genomes to Fields (G2F) initiative.

## I. INTRODUCTION

Decision making under critical and uncertain situations is a common problem in a wide range of domains including online marketing, finance, health sciences, and robotics. There exist learning algorithms that can learn good policies/strategies for optimal decision making. Contextual bandits is one such framework that models the sequential decision making process by utilizing the side information which is referred to as context [1]. One real-world example of a contextual multi-armed bandit problem is when a news website has to make a decision about which articles to display to a visitor when some information about the visitor is known [2]. In the contextual bandit model a learner interacts with the environment in several rounds. In each round the environment presents a context to the learner and the goal of the learner is to choose an action. Upon selecting an action the learner is presented with a reward associated with the chosen action and the goal of the learner is to maximize the cumulative reward.

Most of the existing work on contextual bandit model assumes that the contexts are known and there are no additional constraints on the learner. However, in many applications there exists scenarios where the contexts are noisy or are forecasting measurements (e.g., weather forecasting or stock market prediction) so that the actual context is unknown, rather a distribution on the context is only available. In such cases, the exact context is a sample from this distribution. Such a model has been studied in [3] and an Upper Confidence Bound (UCB)-based algorithm has been proposed with regret bound guarantee. Additionally, safety/performance is a major concern while making decisions and it is crucial to develop learning algorithms that can perform decision making while ensuring that certain safety/performance conditions are satisfied at each round. Contextual bandits with safety constraints have been studied in [4], [5], [6] and algorithms with guarantees were proposed.

Our goal in this paper is to develop a framework for solving sequential decision making when the contexts are unknown and there are safety/performance constraints imposed on the learner. We motivate our problem setting through a scenario. Consider a scenario where the goal is to develop a recommendation system for smart farming such that based on the details of the farm and the farming conditions, including information on the weather and soil properties, the system presents recommendations on the choice of the crop/seed in order to maximize the overall net profit of the farmer. In this setting, the contexts are not observable, rather a distribution of the contexts are known as the weather and soil conditions are forecasting rather than accurate measurements. Additionally, often farmers impose performance constraints such as the net profit must be at least a certain value. Thus for a given farmland and set of soil properties and climate indices, our goal is to provide recommendations for the crop/seed type such that the annual net profit of the farmer is maximized and the associated constraints are satisfied.

This paper makes the following contributions.

- We formulate a conservative stochastic contextual bandit formulation for real-time decision making when an adversary chooses a distribution on the set of possible contexts and the learner is subject to certain safety/performance constraints.
- We present a UCB-based algorithm, conservative (safe) linear UCB algorithm for stochastic bandits with context distribution and unknown contexts.
- We prove an upper bound on the regret of the algorithm and show that it can be decomposed into three terms: (i) an upper bound for the regret of the standard linear UCB algorithm, (ii) a constant term (independent of time horizon) that accounts for the loss of being conservative in order to satisfy the safety constraint, and (ii) a constant term (independent of time horizon) that

J. Lin and S. Moothedath are with the Department of Electrical and Computer Engineering, Iowa State University, USA. Email: {jiabin, mshana}@iastate.edu.

X. Y. Lee, T. Jubery, S. Sarkar, and B. Ganapathysubramanian are with the Department of Mechanical Engineering, Iowa State University, USA. Email: {xylee, znjubery,soumiks, baskarg}@iastate.edu.

accounts for the loss for the contexts being unknown and only the distrbution being known.

- We validated the performance of our approach via extensive simulations on synthetic data and on real-world maize data collected through the Genomes to Fields (G2F) initiative.

The rest of the paper is organized as follows. In Section II we present the notations and the problem formulation. In Section III, we present the related work. In Section IV we present the solution approach for the conservative stochastic bandit problem. In Section V, we present the regret analysis and prove an upper bound on the regret of our proposed algorithm. In Section VII, we present the conclusion and future work.

## II. NOTATIONS AND PROBLEM FORMULATION

In this section, we first specify the standard linear bandit problem below and then explain the stochastic constrained bandit setting. Let $\mathcal{X}$ denote the action set and $\mathcal{C}$ denote the context set. The environment is defined by a fixed and unknown function $y : \mathcal{X} \times \mathcal{C} \to \mathbb{R}$. In linear bandit setting, at any time $t \in \mathbb{N}$, the agent observes a context $c_t \in \mathcal{C}$ and has to choose an action $x_t \in \mathcal{X}$. Each context-action pair $(x,c)$, $x \in \mathcal{X}$ and $c \in \mathcal{C}$, is associated with a feature vector $\phi_{x,c} \in \mathbb{R}^d$, i.e., $\phi_{x_t,c_t} = \phi(x_t, c_t)$. Upon selection of an action $x_t$, the agent observes a reward $y_t \in \mathbb{R}$

$$y_t := \langle \theta^\star, \phi_{x_t,c_t} \rangle + \eta_t, \qquad (1)$$

where $\theta^\star \in \mathbb{R}^d$ is the unknown reward parameter, $\langle \theta^\star, \phi_{x_t,c_t} \rangle = r(x_t, c_t)$ is the expected reward for action $x_t$ at time $t$, i.e., $r(x_t, c_t) = \mathbb{E}[y_t]$, and $\eta_t$ is $\sigma-$Gaussian, additive noise. The goal is to choose optimal actions $x_t^\star$ for all $t \in T$ such that the cumulative reward, $\sum_{t=1}^T y_t$, is maximized. This is equivalent to minimizing the cumulative (pseudo)-regret denoted as

$$\mathcal{R}_T = \sum_{t=1}^T \langle \theta^\star, \phi_{x_t^\star,c_t}^t \rangle - \sum_{t=1}^T \langle \theta^\star, \phi_{x_t,c_t}^t \rangle. \qquad (2)$$

Here $x_t^\star$ is the optimal/best action for context $c_t$ and $x_t$ is the action chosen by the agent for context $c_t$. We make the standard assumptions on the additive noise $\eta_t$ and the unknown parameter $\theta^\star$ [3], [4].

**Assumption 1.** *Each element $\eta_t$ of the noise sequence $\{\eta_t\}_{t=1}^\infty$ is conditionally $\sigma-$subGaussian, i.e.,*

$$\textit{For all } \zeta \in \mathcal{R}, \mathbb{E}[e^{\zeta \eta_t} | x_{1:t}, \varepsilon_{1:t-1}] \geqslant exp(\frac{\zeta^2 \sigma^2}{2}).$$

**Assumption 2.** *There exists constant $A, D \geqslant 0$ such that $\|\theta^\star\|_2 \leqslant A$, $\|\phi_{x,c_t}\|_2 \leqslant D$, and $\phi_{x,c_t}^\top \theta \in [0,1]$, for all $t$ and all $x \in \mathcal{X}$.*

In this work, we consider a *conservative and stochastic* linear bandit setting with context distribution and unknown contexts, i.e., a bandit problem with performance constraints and unknown contexts. We assume that the context at time $t$, $c_t$ is unobservable rather only a distribution of the context denoted as $\mu_t$ is observed by the agent. At round $t$, the environment chooses a distribution $\mu_t \in \mathcal{P}(\mathcal{C})$ over the context

set and samples a context realization $c_t \sim \mu_t$. The learner observes only $\mu_t$ and not $c_t$ and chooses an action, say $x_t$. In addition, there exists a baseline policy (farmer's strategy) $\pi_b$ that at each round $t$, selects action $b_t \in \mathcal{X}$ and incurs the expected reward $r(b_t, c_t) = \langle \theta^\star, \phi_{b_t,c_t} \rangle$. We assume that the expected rewards of the actions taken by the baseline policy, $r(b_t, c_t)$, are known. This assumption is often reasonable as we typically have access to a large amount of data generated using the baseline policy (i.e., the farmer's strategy) and hence can obtain a good estimate of the baseline reward function [4].

Based on the baseline policy, a conservative linear bandit imposes performance constraints. The constraints are such that at round $t$, the difference between the performances of the baseline and the learner's policies should remain above a pre-defined fraction $\alpha \in (0,1)$ of the baseline performance. Our aim is to learn an optimal mapping/policy $g : \mathcal{C} \to \mathcal{X}$ of contexts to actions such that the cumulative reward, $\sum_{t=1}^T y_t$ is maximized while simultaneously satisfying the performance constraints. Formally, our aim is to minimize the cumulative regret

$$\mathcal{R}_T = \sum_{t=1}^T \langle \theta^\star, \phi_{x_t^\star,c_t} \rangle - \sum_{t=1}^T \langle \theta^\star, \phi_{x_t,c_t} \rangle. \qquad (3)$$

such that

$$\sum_{i=1}^t r(b_t, c_t) - \sum_{i=1}^t r(x_t, c_t) \leqslant \alpha \sum_{i=1}^t r(b_t, c_t), \text{ for all } t \in T. \quad (4)$$

Here, $x_t^\star = \arg\max_{x \in \mathcal{X}} \mathbb{E}_{c \sim \mu_t}[r(x,c)]$ is the best action provided we know $\mu_t$, but not $c_t$, $T$ is the number of rounds, and $\alpha \in (0,1)$ is the maximum decrease in the performance the decision maker is willing to accept. Eq. (4) is equivalent to $\sum_{i=1}^t r(x_t, c_t) \geqslant (1-\alpha) \sum_{i=1}^t r(b_t, c_t)$.

## III. RELATED WORK

Bandit algorithms are well studied in the literature, for a survey see [1] and [7]. Recently, contextual bandits have attracted increased attention. Related to our work is stochastic contextual bandits, where the learner chooses actions after observing the contexts and the goal is to learn an optimal mapping from contexts to actions. While stochastic contextual bandits have similarities to Reinforcement Learning (RL) [8], the key difference is that the sequence of contexts can be arbitrary and even chosen by an adversary unlike in an RL setting which has a specific transition structure. Linear contextual bandits is a popular variant of the contextual bandits and it has been studied in [9], [10], [11], [2], [12], [13], [14] and strong theoretical guarantees are established using different solution approaches. The most popular solution approach is the Upper Confidence Bound (UCB) algorithm [2], [15]. The UCB was later improved in [9], [11], [16] with stronger guarantees. Another solution approach is using Thompson sampling and algorithms with theoretical guarantees are provided in [13]. In the linear contextual bandit setting, there are no constraints that need to be satisfied by the learner and the context in round $t$ is known and hence it is a special case of the bandit setting considered in this paper with no constraints and the choice

of the distribution $\mu_t$ as a Dirac delta distribution denoted as $\mu_t = \delta_{c_t}$ for all $t \in T$.

Our work is more closely related to two settings of the contextual bandit problem, the stochastic bandit framework and the constrained contextual bandit framework. A linear contextual bandit setting with uncertainty in the context is studied in [3], [17], [18]. While [18] considered a setting with perturbed contexts, [3] considered a setting in which the context itself is not observable rather a distribution on the context is available and is more closely related to this work. We note that, there are no safety constraints in [3]. There are two different settings where constraints have been applied to the stochastic MAB problem [4], [5], [19], [6], [20], [21]. The first line of work considers the MAB problem with global budget constraints where each arm is associated with a random resource consumption and the objective is to maximize the total reward before the learner exhausts all of its resources [19], [22]. Constrained linear bandit with linear budget constraints is studied in [19] and a primal-dual algorithm is presented. A generalized version of the problem studied in [19], where the objective is concave and constraints are convex is studied in [22] and a UCB-based algorithm was proposed. We note that, the constraints in [19], [22] are modeled as budget constraints unlike in this paper which consider a performance constraint. The second line of work considers safety/performance constraints for bandit problems by ensuring that the performance of the learning algorithm should remain above a pre-defined fraction of the performance of a baseline policy [4], [5], [6]. Among these our work is closely related to [4], but the key difference is that the contextx are unknown in our setting. In this paper we built on the works in [3], [4] to address the conservative and stochastic contextual bandit problem in which the contexts are uncertain and the learner is subject to performance constraints imposed by some baseline policy.

## IV. SOLUTION APPROACH: STOCHASTIC CONSERVATIVE CONTEXTUAL BANDIT

In this section, we present the algorithm for solving the stochastic conservative contextual bandit problem. Our solution approach is built on the works of [3] and [4]. Given the distribution $\mu_t$, we construct the expected feature vector, $\Psi_t = \{\bar{\psi}_{x,\mu_t} : x \in \mathcal{X}\}$ where $\{\bar{\psi}_{x,\mu_t} := \mathbb{E}_{c \sim \mu_t}[\phi_{x,c}]\}$ (step: 5). We note that, each feature $\bar{\psi}_{x,\mu_t}$ corresponds to exactly one action $x \in \mathcal{X}$ and we use $\Psi_t$ as the feature context set at time $t$. The proposed algorithm is based on the *optimism in the face of uncertainty* principle, where the algorithm maintains a confidence set $\mathcal{B}_t \subset \mathcal{R}^d$ that contains the unknown parameter vector $\theta^\star$ with high probability [9]. The algorithm then chooses an optimistic estimate $\tilde{\theta}_t = \arg\max_{\hat{\theta} \in \mathcal{B}_t} (\max_{x \in \mathcal{X}} \bar{\psi}_{x,\mu_t}^\top \hat{\theta})$ and chooses an action $x'_t = \arg\max_{x \in \mathcal{X}} \bar{\psi}_{x,\mu_t}^\top \tilde{\theta}_t$. Equivalently the algorithm chooses the pair $(x'_t, \tilde{\theta}_t) \in \arg\max_{(x,\hat{\theta}) \in \mathcal{X} \times \mathcal{B}_t} \bar{\psi}_{x,\mu_t}^\top \hat{\theta}$ which jointly maximizes the reward.

To ensure that the action chosen by the algorithm guarantees satisfaction of the constraints, the algorithm plays the action $x'_t$ only if it satisfies the constraint for the worst choice of the parameter $\hat{\theta} \in \mathcal{B}_t$ [4]. We formally define

this by introducing two sets $S^b_{t-1}$ and $S_{t-1}$. Let $S_{t-1}$ be the set of rounds $i$ before round $t$ at which the algorithm has played the optimistic action, i.e., $x_i = x'_i$. Then $S^b_{t-1} = \{1, 2, \ldots, t-1\} - S_{t-1}$ is the set of rounds $j$ before round $t$ at which the algorithm has followed the baseline policy, i.e., $x_j = b_j$. To ensure that constraint in Eq. (4) is satisfied the algorithm plays optimal action $x_t = x'_t$ at round $t$ if it satisfies

$$\min_{\hat{\theta} \in \mathcal{B}_t} \Big[ \sum_{i \in S^b_{t-1}} r(b_t, c_t) + (\sum_{i \in S_{t-1}} \bar{\psi}_{x_i,\mu_i})^\top \hat{\theta} + \bar{\psi}_{x'_t,\mu_t}^\top \hat{\theta} \Big] \geqslant (1-\alpha) \sum_{i=1}^t r(b_i, c_i),$$

and plays the action chosen by the farmer, i.e., $x_t = b_t$ otherwise.

---

**Algorithm IV.1** Pseudocode for conservative stochastic contextual bandit with context distribution

**Input:** $\alpha, \mathcal{B} = \mathbb{R}^d$
1: **Initialize:** $S_0 = \emptyset, \ell_0 = 0 \in \mathbb{R}^d$, $\mathcal{B}_1 = \mathcal{B}$
2: **for** $t = 1, 2, \ldots, T$ **do**
3:     Nature chooses $\mu_t \in \mathcal{P}(\mathcal{C})$
4:     Learner observes $\mu_t$
5:     Set $\Psi_t = \{\bar{\psi}_{x,\mu_t} : x \in \mathcal{X}\}$ where $\{\bar{\psi}_{x,\mu_t} := \mathbb{E}_{c \sim \mu_t}[\phi_{x,c}]\}$
6:     Query baseline strategy $b_t \leftarrow \pi(\Psi_t)$
7:     Find $(x'_t, \tilde{\theta}_t) \in \arg\max_{(x,\hat{\theta}) \in \mathcal{X} \times \mathcal{B}_t} \bar{\psi}_{x,\mu_t}^\top \hat{\theta}$
8:     Compute $L_t = \min_{\hat{\theta} \in \mathcal{B}_t} \langle \ell_{t-1} + \bar{\psi}_{x'_t,\mu_t}, \hat{\theta} \rangle$
9:     **if** $L_t + \sum_{i \in S^b_{t-1}} r(b_t, c_t) \geqslant (1-\alpha) \sum_{i=1}^t r(b_t, c_t)$ **then**
10:         Play $x_t = x'_t$ and observe reward $y_t$ in Eq. (1)
11:         Set $\ell_t = \ell_{t-1} + \bar{\psi}_{x_t,\mu_t}$, $S_t = S_{t-1} \cup t$, $S^b_t = S^b_{t-1}$
12:         Given $x_t, y_t$ construct $\mathcal{B}_{t+1}$ using Eq. (6)
13:     **else**
14:         Play $x_t = b_t$ and observe reward $y_t$ in Eq. (1)
15:         Set $\ell_t = \ell_{t-1}$, $S_t = S_{t-1}$, $S^b_t = S^b_{t-1} \cup t$, $\mathcal{B}_{t+1} = \mathcal{B}_t$
16:     **end if**
17: **end for**

---

**Construction of the Confidence Set $\mathcal{B}_t$:** We denote the confidence set in round $t$ as $\mathcal{B}_t$. The proposed algorithm starts by the most general confidence set i.e., $\mathcal{B}_1 = \mathcal{B} = \mathbb{R}^d$, and updates the confidence set only when the optimistic action proposed by the learner is played. This is because that unless the learner's action is played, no additional information is gained about the unknown parameter $\theta$. Let $S_t = \{i_1, i_2, \ldots, i_{m_t}\}$ be the set of rounds up to and including $t$ during which the the algorithm played the optimistic action. Here $m_t = |S_t|$. For a fixed value $\lambda > 0$, the regularized least square estimate of $\hat{\theta}$ at round $t$ is given by

$$\bar{\theta}_t = \left( \Phi_t \Phi_t^\top + \lambda I \right)^{-1} \Phi_t Y_t, \tag{5}$$

where $\Phi_t = [\bar{\psi}_{x_{i_1},\mu_{i_1}}, \bar{\psi}_{x_{i_2},\mu_{i_2}}, \ldots, \bar{\psi}_{x_{m_t},\mu_{m_t}}]$ and $Y_t = [y_{i_1}, y_{i_2}, \ldots, y_{m_t}]^\top$. For a given confidence parameter $\delta \in (0,1)$, we construct the confidence set for the next round $t+1$ as

$$\mathcal{B}_{t+1} = \{\hat{\theta} \in \mathcal{R}^d : \|\hat{\theta} - \bar{\theta}_t\|_{V_t} \leq \beta_{t+1}\}, \tag{6}$$

where $\beta_{t+1} = \sigma \sqrt{d \log(\frac{1 + (m_t+1)D^2/\lambda}{\delta})} + \sqrt{\lambda} A$, $V_t = \lambda I + \Phi_t \Phi_t^\top$, and the weighted norm is defined as $\|u\|_V = \sqrt{u^\top V u}$ for any $u \in \mathcal{R}^d$ and positive definite $V \in \mathcal{R}^{d \times d}$.

**Proposition 1.** *For any $\delta > 0$ and the confidence set $\mathcal{B}_t$ defined by Eq. (6), we have*

$$\mathbb{P}[\theta^\star \in \mathcal{B}_t, \forall t \in \mathbb{N}] \geqslant 1 - \delta.$$

At each round $t$, Algorithm IV.1 ensures that Eq. (4) holds for all $\theta \in \mathcal{B}_t$. From Proposition 1, $\mathbb{P}[\theta^\star \in \mathcal{B}_t] \geqslant 1 - \delta$ for all $t \in \mathbb{N}$. Thus, Proposition 1 ensures that at each round $t$, Algorithm IV.1 satisfies the baseline criteria in Eq. (4) with probability at least $1 - \delta$.

## V. REGRET ANALYSIS

In this section, we prove the regret bound for Algorithm IV.1. Let $\Delta_{b_t}^t = r(x_t^\star, c_t) - r(b_t, c_t)$ be the baseline gap at round $t$, i.e., the difference between the expected rewards of optimal action and baseline action at round $t$.

**Assumption 3.** *There exists $0 \leq \Delta_\ell \leq \Delta_h$ and $0 < r_\ell < r_h$ such that, at each round $t$,*

$$\Delta_\ell \leq \Delta_{b_t}^t \leq \Delta_h \text{ and } r_\ell \leq r(b^t, c_t) \leq r_h.$$

Since the rewards belong to $[0,1]$ (Assumption 2), we set $\Delta_h = r_h = 1$, and $\Delta_\ell = 0$. The reward lower bound $r_l$ ensures that the baseline policy satisfies a minimum level of performance guarantee at each round of the algorithm.

**Proposition 2** ([3], Lemma 3). *The regret of the UCB algorithm for linear stochastic bandits with expected feature set $\Psi_t$ is bounded in time $T$ with probability at least $1 - \delta$,*

$$\mathcal{R}_T \leq \mathcal{R}_T^{UCB} + 4\sqrt{2T \log \frac{1}{\delta}}.$$

**Lemma 1.** *The regret of Algorithm IV.1 with expected feature set $\Psi_t$ is bounded in time $T$ with probability at least $1 - \delta$,*

$$\mathcal{R}_T \leq \mathcal{R}_{S_T}^{UCB} + 4\sqrt{2m_T \log \frac{1}{\delta}} + n_T \Delta_h,$$

*where $\mathcal{R}_{S_T}^{UCB}$ is the cumulative (pseudo)-regret of linear UCB algorithm at rounds $t \in S_T$, $m_T = |S_T|$ is the number of times Algorithm IV.1 played the learner's action, and $n_T = |S_T^b| = T - |S_T| = T - m_T$ is the number of times Algorithm IV.1 played the baseline action.*

*Proof.* From the definition of regret in Eq. (2)

$$
\begin{aligned}
\mathcal{R}_T &= \sum_{t=1}^T r(x_t^\star, c_t) - \sum_{t=1}^T r(x_t, c_t), \\
&= \sum_{t \in S_T}(r(x_t^\star, c_t) - r(x_t, c_t)) + \sum_{t \in S_T^b}(r(x_t^\star, c_t) - r(x_t, c_t)), \\
&= \sum_{t \in S_T}(r(x_t^\star, c_t) - r(x_t, c_t)) + \sum_{t \in S_T^b}\Delta_{b_t}^t, \\
&\leq \sum_{t \in S_T}(r(x_t^\star, c_t) - r(x_t, c_t)) + n_T \Delta_h, \\
&\leq \mathcal{R}_{S_T}^{UCB} + 4\sqrt{2m_T \log \frac{1}{\delta}} + n_T \Delta_h. \quad (7)
\end{aligned}
$$

Inequality in Eq. (7) follows from Proposition 2 since for $t \in S_T$, Algorithm IV.1 plays the same actions as the UCB algorithm in [3] and this completes the proof. $\square$

The regret bound for linear UCB algorithm for the confidence set given in Eq. (6) is given in [9]. Let $\varepsilon$ be the event that $\theta^\star \in \mathcal{B}_t$ for all $t \in \mathbb{N}$. By Proposition 1 the probability of $\varepsilon$ is at least $1 - \delta$. The result below from [9] presents the bound for $\mathcal{R}_{S_T}^{UCB}$.

**Proposition 3** ([4], Proposition 3). *On event $\varepsilon$, for any $T \in \mathbb{N}$, we have*

$$
\begin{aligned}
\mathcal{R}_{S_T}^{UCB} &\leq 4\sqrt{m_T d \log\left(\lambda + \frac{m_T D}{d}\right)} \\
&\quad \times \left[A\sqrt{\lambda} + \sigma\sqrt{2\log(1/\delta) + d\log\left(1 + \frac{m_T D}{\lambda d}\right)}\right] \\
&= O\left(d\log(\frac{D}{\lambda \delta}T)\sqrt{T}\right). \quad (8)
\end{aligned}
$$

We note that, to bound the regret of Algorithm IV.1, we only need to find upper bounds on $n_T$, the number of times Algorithm IV.1 deviates from the UCB algorithm for linear stochastic bandits and plays the baseline, and $m_T$, the number of times Algorithm IV.1 plays the action suggested by the UCB algorithm for linear stochastic bandits. Since $m_T = T - n_T$, it also suffices to find an upper and lower bounds for $n_T$. An upper bound for $n_T$ is given in [4] which is presented in the proposition below.

**Proposition 4** ([4], Theorem 5). *Assume that $\lambda \geqslant \max\{1, D^2\}$. On event $\varepsilon$, for any horizon $T \in \mathbb{N}$, we have*

$$n_T \leqslant 1 + 114 d^2 \frac{(A\sqrt{\lambda} + \sigma)^2}{\alpha r_\ell (\Delta_\ell + \alpha r_\ell)}\left[\log\left(\frac{62 d(A\sqrt{\lambda} + \sigma)}{\sqrt{\delta}(\Delta_\ell + \alpha r_\ell)}\right)\right]^2.$$

Thus the only thing remaining to prove is a lower bound on $n_T$. To prove a lower bounds on $n_T$, we use Proposition 5 from [4] and Lemma 2.

**Proposition 5** ([4], Lemma 4). *For given $k \in \mathbb{N}$, $\lambda > 0$, and any sequence $Y_1, Y_2, \ldots, Y_k$ in $\mathbb{R}^d$ such that for all $i : \|Y_i\|_2 \leqslant D$, let $V_0 = \lambda I$ and $V_i = \lambda I + \sum_{j=1}^i Y_j Y_i^\top$ for $1 \leqslant i \leq k$. Then, we have*

$$\sum_{i=1}^k \min\left(1, \|Y_i\|_{V_{i-1}^{-1}}^2\right) \leqslant 2d\log\left(1 + \frac{kD^2}{\lambda d}\right). \quad (9)$$

**Lemma 2.** *For any $m \geq 2$ and $c_1, c_2, c_3 > 0$, $-c_3 m - c_1\sqrt{m}\log(c_2 m) \geq \frac{16c_1^2}{25c_3}\left[\log(\frac{2c_1\sqrt{c_2}e}{c_3})\right]^2.$*

*Proof.* Let $g(m) = -c_3 m - c_1\sqrt{m}\log(c_2 m)$. Then, $g'(m) = -c_3 - \frac{c_1(2 + \log(c_2 m))}{2\sqrt{m}}$ and $g''(m) = \frac{c_1 \log(c_2 m)}{4m\sqrt{m}}$. Since $c_2 > 1$, $g$ is a convex function over its domain $[2, \infty)$, and thus a global optimum $m^\star$ exists for $g$. By the first order condition, we get $g'(m^\star) = 0$. This gives

$$2 + \log(c_2 m^\star) = \frac{-2c_3}{c_1}\sqrt{m^\star}. \quad (10)$$

Thus $g^\star = g(m^\star) = c_3 m^\star + 2c_1\sqrt{m^\star}$. Using change of variables $z = \frac{c_3}{2c_1}\sqrt{m^\star}$, we get

$$g^\star = \frac{4c_1^2}{c_3}(z^2 + z). \quad (11)$$

Eq. (10) becomes

$$2 + \log(\frac{4c_2 c_1^2}{c_3^2}) + 2\log(z) = -4z.$$

After taking exponential on both sides,

$$\frac{e^{-4z}}{z^2} = \frac{4c_1^2 c_2 e^2}{c_3^2}.$$

Using $e^z > z^2$,

$$\frac{4c_1^2 c_2 e^2}{c_3^2} = \frac{e^{-4z}}{z^2} > \frac{e^{-4z}}{e^z} = e^{-5z}.$$

Thus

$$z \geq \frac{-1}{5}\log(\frac{4c_1^2 c_2 e^2}{c_3^2}).$$

Substituting in Eq. (11), we get

$$g^\star \geq \frac{4c_1^2}{c_3}z^2 \geq \frac{4c_1^2}{25c_3}\left[\log(\frac{4c_1^2 c_2 e^2}{c_3^2})\right]^2 = \frac{16c_1^2}{25c_3}\left[\log(\frac{2c_1\sqrt{c_2}e}{c_3})\right]^2.$$

$\square$

**Theorem 1.** *Assume that $\lambda \geq D^2$. On event $\varepsilon$, for any horizon $T \in \mathbb{N}$, the following holds*

$$n_T \geq \frac{d^2(A\sqrt{\lambda}+\sigma)^2}{\alpha r_h(\Delta_h + \alpha r_h)}\left[\log\left(\frac{10d(A\sqrt{\lambda}+\sigma)}{\sqrt{\delta}(\Delta_h + \alpha r_h)}\right)\right]^2.$$

*Proof.* Let $\tau$ be the last round in which Algorithm IV.1 plays the learner's action, $\tau = \max\{1 \leq t \leq T | x_t = x'_t\}$.

$$\min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t \rangle + \sum_{t \in S_{\tau-1}^b}r(b_t, c_t) \geq (1-\alpha)\sum_{t=1}^\tau r(b_t, c_t),$$

$$\alpha\sum_{t=1}^\tau r(b_t, c_t) \geq \sum_{t \in S_{\tau-1}}r(b_t, c_t) + r(b_\tau, c_\tau) - \min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle,$$

$$\geq \sum_{t \in S_{\tau-1}}(r(b_t, c_t) - \langle \theta^\star, \bar{\psi}_{x_t}^t\rangle) + (r(b_\tau, c_\tau) - \langle \theta^\star, \bar{\psi}_{x'_\tau}^\tau\rangle)$$

$$+ \langle \theta^\star, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle - \min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle$$

$$\geq \sum_{t \in S_{\tau-1}}(-\Delta_{b_t}^t) - \Delta_{b_\tau}^\tau - \min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle$$

$$\geq \sum_{t \in S_{\tau-1}}(-\Delta_h) - \Delta_h - \min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle$$

$$= -(m_{\tau-1}+1)\Delta_h - \min_{\theta \in \mathcal{B}_\tau}\langle \theta, \bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\rangle \quad (12)$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \|\theta\|_{V_\tau}\left\|\bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\right\|_{V_\tau^{-1}}$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \beta_\tau\left\|\bar{\psi}_{x'_\tau}^\tau + \sum_{t \in S_{\tau-1}}\bar{\psi}_{x_t}^t\right\|_{V_\tau^{-1}}$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \beta_\tau(\left\|\bar{\psi}_{x'_\tau}^\tau\right\|_{V_\tau^{-1}} + \sum_{t \in S_{\tau-1}}\left\|\bar{\psi}_{x_t}^t\right\|_{V_t^{-1}}) \quad (13)$$

From Eq. (12), we get

$$\alpha\sum_{t=1}^\tau r(b_t, c_t) \geq -(m_{\tau-1}+1)\Delta_h - (m_{\tau-1}+1) \quad (14)$$

We note that, $\beta_\tau$ is non decreasing and is greater than 1. Hence from Eqs. (13) and (14),

$$\alpha\sum_{t=1}^\tau r(b_t, c_t) \geq -(m_{\tau-1}+1)\Delta_h - \beta_\tau[\min(\left\|\bar{\psi}_{x'_\tau}^\tau\right\|_{V_\tau^{-1}}, 1)$$

$$+ \sum_{t \in S_{\tau-1}}\min(\left\|\bar{\psi}_{x_t}^t\right\|_{V_t^{-1}}, 1)] \quad (15)$$

In order to simplify the equation, we introduce $\Gamma$ as

$$\Gamma := \left[\min(\left\|\bar{\psi}_{x'_\tau}^\tau\right\|_{V_\tau^{-1}}^2, 1) + \sum_{t \in S_{\tau-1}}\min(\left\|\bar{\psi}_{x_t}^t\right\|_{V_t^{-1}}^2, 1)]\right].$$

By Cauchy-Schwarz inequality and using Proposition 5, and $\Gamma$, Eq. (15) can be written as

$$\alpha\sum_{t=1}^\tau r(b_t, c_t) \geq -(m_{\tau-1}+1)\Delta_h - \beta_\tau\sqrt{(m_{\tau-1}+1)\Gamma}$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \beta_\tau\sqrt{2(m_{\tau-1}+1)d\log(1+\frac{(m_{\tau-1}+1)D^2}{\lambda d})}$$

$$= -(m_{\tau-1}+1)\Delta_h - \sqrt{2(m_{\tau-1}+1)d\log(1+\frac{(m_{\tau-1}+1)D^2}{\lambda d})}$$

$$\times \left(\sqrt{\lambda}A + \sigma\sqrt{d\log(\frac{1+(m_{\tau-1}+1)D^2/\lambda}{\delta})}\right)$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \sqrt{2(m_{\tau-1}+1)d\log(1+\frac{(m_{\tau-1}+1)}{d})}$$

$$\times \left(\sqrt{\lambda}A + \sigma\sqrt{d\log(\frac{1+(m_{\tau-1}+1)}{\delta})}\right) \quad (16)$$

$$\geq -(m_{\tau-1}+1)\Delta_h - \left(\sqrt{2d}\sqrt{m_{\tau-1}+1}(A\sqrt{\lambda}+\sigma)\right)$$

$$\times \log(\frac{2(m_{\tau-1}+1)}{\delta}) \quad (17)$$

The inequality in Eq. (16) follows from $D^2 \geq \lambda$ and the inequality in Eq. (17) holds since

$$\left(\sqrt{\lambda}A + \sigma\sqrt{d\log(\frac{1+(m_{\tau-1}+1)}{\delta})}\right)$$

$$\leq \left(A\sqrt{\lambda}+\sigma\right)\sqrt{d\log(\frac{2(m_{\tau-1}+1)}{\delta})}$$

Eq. (17) can be rewritten as

$$\alpha\sum_{t=1}^\tau r_h \geq -(m_{\tau-1}+1)\Delta_h - \left(\sqrt{2d}\sqrt{m_{\tau-1}+1}(A\sqrt{\lambda}+\sigma)\right)$$

$$\times \log(\frac{2(m_{\tau-1}+1)}{\delta})$$

$$\alpha n_{\tau-1}r_h \geq -(m_{\tau-1}+1)(\Delta_h + \alpha r_h) - \left(\sqrt{2d}\sqrt{m_{\tau-1}+1}\right)$$

$$\times (A\sqrt{\lambda}+\sigma)\log(\frac{2(m_{\tau-1}+1)}{\delta}) \quad (18)$$

Eq. (18) follows after substituting $\alpha\sum_{t=1}^\tau r_h = \alpha(m_{\tau-1}+n_{\tau-1}+1)r_h$. To prove a lower bound on the LHS of Eq. (18), we first present a lower bound for the RHS of Eq. (18). Let $m = (m_{\tau-1}+1)$, $c_1 = \sqrt{2}d(A\sqrt{\lambda}+\sigma)$, $c_2 = \frac{2}{\delta}$, $c_3 = (\Delta_h + \alpha r_h)$. By Lemma 2,

$$\alpha n_{\tau-1}r_h \geq \frac{d^2(A\sqrt{\lambda}+\sigma)^2}{(\Delta_h + \alpha r_h)}\left[\log\left(\frac{10d(A\sqrt{\lambda}+\sigma)}{\sqrt{\delta}(\Delta_h + \alpha r_h)}\right)\right]^2$$

The result follows as $n_T \geq n_\tau = n_{\tau-1}$. $\square$

We now present the regret bound on Algorithm IV.1 in Theorem 2. Proof of Theorem 2 follows from Lemma 1, Proposition 4, Theorem 1, and Proposition 3.

**Theorem 2.** *With probability at least $1-\delta$, Algorithm IV.1 satisfies the performace constraint in Eq. (4) for all $t \in \mathbb{N}$, and*
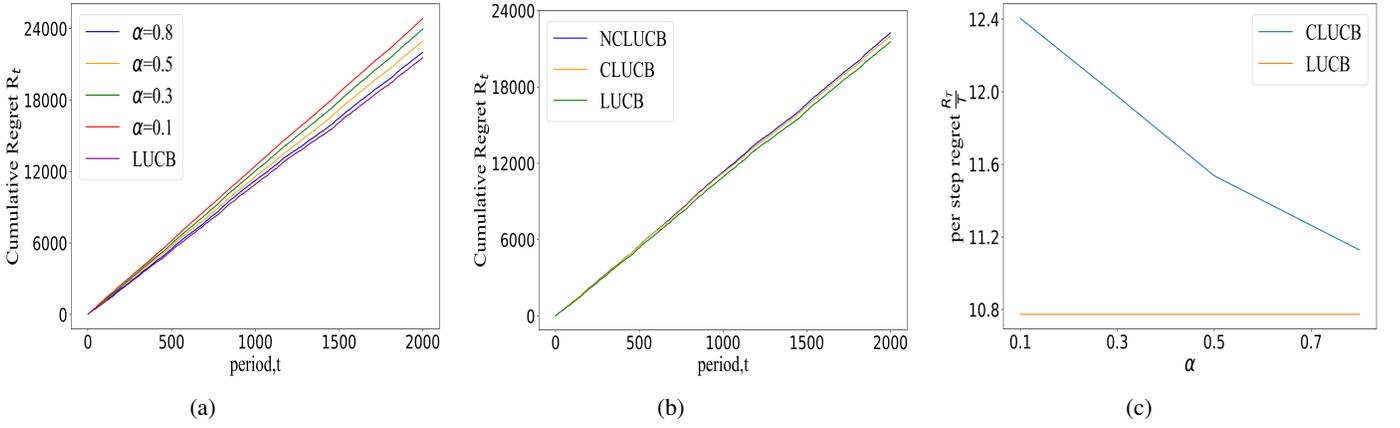
Fig. 1: Plots for synthetic data. (a) Cumulative regret of the standard linear UCB algorithm (LUCB [9]) and conservative stochastic bandit algorithm with context distribution (Algorithm IV.1) with $\alpha = 0.1, 0.3, 0.5, 0.8$, (b) Cumulative regret for three settings: (i) when the learner observes the context and there are no safety constraints (LUCB [9]), (ii) when the learner observes the context and there are safety constraints (conservative linear UCB, CLUCB [4]), and (iii) when the learner observes only the context distribution and there are safety constraints (Algorithm IV.1) (c) Comparison of per step regret $\mathcal{R}_T/T$ at $T = 2000$ for different values of $\alpha$ for (i), (ii), and (iii).
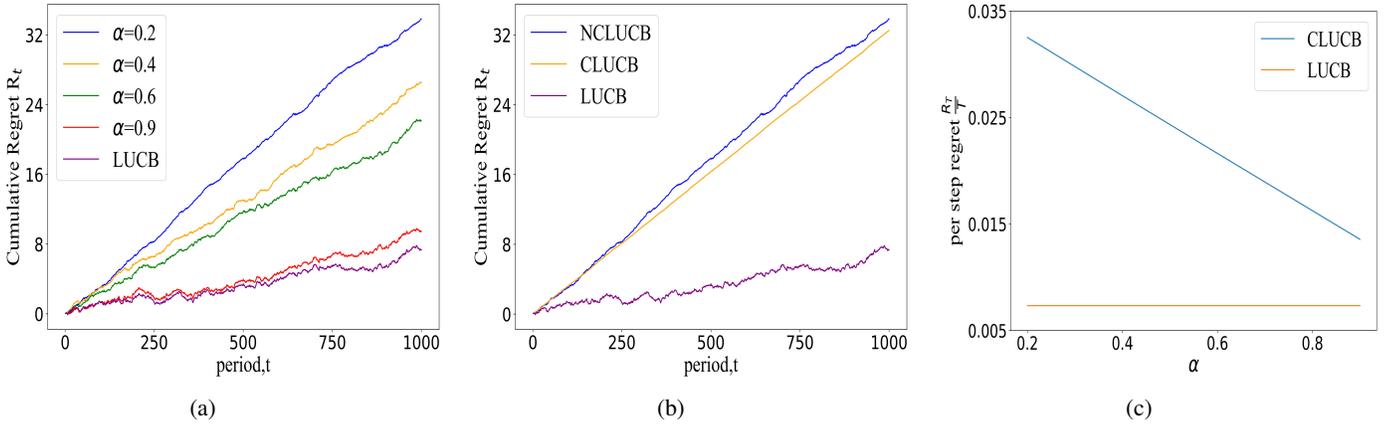


Fig. 2: Plots for maize yield data. (a) Cumulative regret of the standard linear UCB algorithm (LUCB [9]) and conservative stochastic bandit algorithm with context distribution (Algorithm IV.1) with $\alpha = 0.2, 0.4, 0.6, 0.9$, (b) Cumulative regret for three settings: (i) when the learner observes the context and there are no safety constraints (LUCB [9]), (ii) when the learner observes the context and there are safety constraints (conservative linear UCB, CLUCB [4]), and (iii) when the learner observes only the context distribution and there are safety constraints (Algorithm IV.1) (c) Comparison of per step regret $\mathcal{R}_T/T$ at $T = 1000$ for different values of $\alpha$ for (i), (ii), and (iii).

*satisfies the following regret bound*

$$\mathcal{R}_T \leq O\left(d\log(\frac{D}{\lambda\delta}T)\sqrt{T} + \frac{K_h\Delta_h}{\alpha r_\ell} + \frac{K_\ell\sqrt{\log(1/\delta)}}{\sqrt{\alpha r_h(\Delta_h + \alpha r_h)}}\right),$$

*where $K_h$ and $K_\ell$ are constants that depend only on the parameters of the problem as $K_h = 1 + 114d^2 \frac{(A\sqrt{\lambda}+\sigma)^2}{\Delta_\ell + \alpha r_\ell}\left[\log\left(\frac{62d(A\sqrt{\lambda}+\sigma)}{\sqrt{\delta}(\Delta_\ell + \alpha r_\ell)}\right)\right]^2$ and $K_\ell = d(A\sqrt{\lambda}+\sigma)\left[\log\left(\frac{10d(A\sqrt{\lambda}+\sigma)}{\sqrt{\delta}(\Delta_h + \alpha r_h)}\right)\right].$*

*Proof.* The proof follows from Lemma 1, proposition 3, proposition 4, and Theorem 1. $\square$

## VI. Experimental Analysis

In this section, we present the experimental validation of our approach on two data sets (i) synthetic data and (ii) real-world maize data.

*A. Synthetic Data:* We considered a context set $\mathcal{C}$ and an action set $\mathcal{X}$ with 5-dimensional contexts and actions, i.e., $c \in \mathbb{R}^5$ and $x \in \mathbb{R}^5$,. Further we set the reward function $r(x_i, c_i) = \sum_{i=1}^{5}(x_i - c_i)^2$. Thus the parameterized vector $\phi(x,c)$ is given by $\phi() = [x_1^2, \ldots, x_5^2, c_1^2, \ldots, c_5^2, x_1 c_1, \ldots, x_5 c_5]$ and the reward parameter $\theta^\star = [1,1,1,1,1,1,1,1,1,1,-2,-2,-2,-2,-2]$. The action set consists of 20 actions that we sample from a standard Gaussian distribution. At each round $t \in T$, we sample the context $c_t$ from a multi-variate normal distribution and set the context distribution as $\mu_t = \mathcal{N}(c_t, \mathbb{1}_5)$. The observation noise $\eta_t$ is set as Gaussian with zero mean and standard deviation 0.1 and the mean reward of the baseline policy at any time is taken to be the reward associated with the 10th best action.

*B. Maize Yield Data:* We use a maize yield data set acquired over four years by the maize Genomes to Fields (G2F) initiative [23], a multi-institutional effort in North America over 68 unique locations. The data set includes yields, planting dates, flowering times, and harvest dates, as well as hourly weather data from

in-field weather stations, such as temperature, humidity, solar radiation, rainfall, and soil wind speed, as well as soil characteristics such as soil texture, organic matter, texture, and nitrogen, phosphorous, potassium, sulfur, and sodium levels (in parts per million). There are 2158 yield measurements (rewards) for 24 crops (action set) collected from 22 different locations in this data set. The weather data of the whole growing season was summarized by crop growth stages as in [24]. These are average daily solar radiation [MJ/m2], average daily minimum temperature below 0°C in absolute values [°C] as a measure of frost impacts, average daily mean temperature [°C] as a measure of temperature determining plant growth, average daily maximum temperature above 35 C [°C] as a measure of heat stress, and average photoperiod.

We first constructed a data set $\mathcal{D} = \{(c_i, x_i, y_i)\}$, where for each data point $i = 1, 2, \ldots, 2158$, the context $c_i \in \mathbb{R}^{28}$ is a $28-$dimensional vector that includes 6-dimensional weather and soil data information (% of sand, % of silt, % of clay in the soil, daily average temperature, radiation, and photosynthesis) and a $22-$dimensional one-hot encoding that captures the field ID, and $x_i, y_i$ are the seed/crop identifier, yield, respectively. We first fit a bilinear model [25] such that $y_i \approx c_i^\top W V_{x_i}$, where $V_{x_i} \in \mathbb{R}^{10}$ is the feature vector for crop type $x_i$ [3]. Our data set consists of 24 varieties of crops and hence there are 24 feature vectors, $V_1, V_2, \ldots, V_{24}$. The bilinear model captures the correlation between site features $c_i^\top W$ and $V_{x_i}$ for each data point and serves as the interactive setting that provides the rewards (yield) for our bandit setting.

We fitted a bilinear model on the historical maize data, collected through the G2F initiative [23], via stochastic gradient descent using the loss function

$$L(V, W) = \sum_{i=1}^{n} (y_i - c_i^T W V_{x_i})^2 + \lambda_v ||V_{x_i}||^2 + \lambda_w ||W||^2,$$

where $\lambda_v$ and $\lambda_w$ denotes the regularization terms. Training this model for 300 iterations resulted in a mean square error loss of 0.002 using a learning rate of 0.015, $\lambda_v = \lambda_w = 0.001$ and a latent dimension of 10, i.e., $V_{x_i} \in \mathbb{R}^{10}$ for all $i$. The observation noise $\eta_t$ is set as Gaussian with zero mean and standard deviation 0.1 and the mean reward of the baseline policy at any time is taken to be the reward associated with the 16th best action.

*C. Experiments and Analysis:* We performed two experiments on the synthetic and real data and all the points are averaged over 100 independent trials. In the first experiment we varied the value of the constraint parameter $\alpha$ and we plot the cumulative regret $\mathcal{R}_t$ at each round $t$. Figure 1a shows the comparison of the cumulative regret of the standard linear UCB algorithm (LUCB) in [9], when contexts are known and there are no constraints, and the conservative stochastic bandit with context distribution and $\alpha = 0.1, 0.3, 0.5$, and 0.8 for the synthetic data. Figure 2a shows the comparison of the cumulative regret of the LUCB algorithm in [9] and the conservative stochastic bandit with context distribution and $\alpha = 0.2, 0.4, 0.6$, and 0.9 for the maize yield data data. From Figures 1a and 2a we observe that as the value of $\alpha$ increases the cumulative regret decreases which is expected as larger value of $\alpha$ means weaker constraint.

In the second experiment, we implemented three cases of our bandit setting: (i) when the learner observes the context and there are no safety/performance constraints (LUCB [9]), (ii) when the learner observes the context and there are safety constraints (conservative linear UCB, CLUCB [4]), and (iii) when the learner

observes only the context distribution and there are safety constraints (conservative stochastic UCB algorithm with context distribution, Algorithm IV.1). Figures 1b and 2b presents the plots for the cumulative regret for the three settings for synthetic data and maize yield data, respectively. We note that, in (i) the decisions are based on the observed context, in (ii) the decisions are based on the observed context and the safety constraint and in (iii) the decisions are based on the context distribution and the safety constraint. From the plots we notice that (i) outperforms (ii) and (ii) outperforms (iii) which is expected. We also present the per step cumulative regret $\mathcal{R}_T/T$ for the synthetic data with $T = 2000$ and fro the maize yield data with $T = 1000$ while varying the value of $\alpha$. From Figures 1c and 2c we notice that the gap between the standard linear UCB algorithm and the conservative stochastic bandit algorithm decreases as the $\alpha$ value increases, which is expected as larger $\alpha$ means weaker constraint.

## VII. CONCLUSION

In this paper, we presented a conservative stochastic contextual bandit framework for sequential decision making when an adversary chooses a distribution on the set of possible contexts and the learner is subject to certain safety/performance constraints. Our bandit formulation is conservative in the sense that we incorporate constraints on the learned policy such that the learned policy need to satisfy certain baseline performance criteria while maximizing the reward. Furthermore, our bandit formulation is stochastic in the sense that the contexts are not observable, rather a distribution of the contexts are known. We proposed a conservative linear UCB algorithm for stochastic bandits with context distribution. We proved an upper bound on the regret of the algorithm and showed that it can be decomposed into three terms: (i) an upper bound for the regret of the standard linear UCB algorithm, (ii) a constant term (independent of time horizon) that accounts for the loss of being conservative in order to satisfy the safety constraint, and (ii) a constant term (independent of time horizon) that accounts for the loss of the contexts being unknown and only the distrbution being known. We validated the performance of our approach through extensive simulations on synthetic data and on real-world maize data collected through the Genomes to Fields (G2F) initiative.

## REFERENCES

[1] S. Bubeck and N. Cesa-Bianchi, "Regret analysis of stochastic and nonstochastic multi-armed bandit problems," *arXiv preprint arXiv:1204.5721*, 2012.

[2] L. Li, W. Chu, J. Langford, and R. E. Schapire, "A contextual-bandit approach to personalized news article recommendation," in *Proceedings of the 19th international conference on World wide web*, 2010, pp. 661–670.

[3] J. Kirschner and A. Krause, "Stochastic bandits with context distributions," *Advances in Neural Information Processing Systems*, vol. 32, pp. 14 113–14 122, 2019.

[4] A. Kazerouni, M. Ghavamzadeh, Y. Abbasi-Yadkori, and B. Van Roy, "Conservative contextual linear bandits," *Advances in Neural Information Processing Systems*, 2017.

[5] Y. Wu, R. Shariff, T. Lattimore, and C. Szepesvári, "Conservative bandits," in *International Conference on Machine Learning*, 2016, pp. 1254–1262.

[6] S. Amani, M. Alizadeh, and C. Thrampoulidis, "Linear stochastic bandits under safety constraints," *arXiv:1908.05814*, 2019.

[7] T. Lattimore and C. Szepesvári, *Bandit algorithms.* Cambridge University Press, 2020.

[8] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction.* MIT press, 2018.

[9] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári, "Improved algorithms for linear stochastic bandits," *Advances in neural information processing systems*, vol. 24, pp. 2312–2320, 2011.

[10] P. Auer, "Using confidence bounds for exploitation-exploration trade-offs," *Journal of Machine Learning Research*, vol. 3, pp. 397–422, 2002.

[11] V. Dani, T. P. Hayes, and S. M. Kakade, "Stochastic linear optimization under bandit feedback," *Annual Conference on Learning Theory (COLT)*, 2008.

[12] W. Chu, L. Li, L. Reyzin, and R. Schapire, "Contextual bandits with linear payoff functions," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 208–214.

[13] S. Agrawal and N. Goyal, "Thompson sampling for contextual bandits with linear payoffs," in *International Conference on Machine Learning*, 2013, pp. 127–135.

[14] R. Allesiardo, R. Féraud, and D. Bouneffouf, "A neural networks committee for the contextual bandit problem," in *International Conference on Neural Information Processing*, 2014, pp. 374–381.

[15] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Machine learning*, vol. 47, no. 2, pp. 235–256, 2002.

[16] Y. Li, Y. Wang, X. Chen, and Y. Zhou, "Tight regret bounds for infinite-armed linear contextual bandits," in *International Conference on Artificial Intelligence and Statistics*, 2021, pp. 370–378.

[17] S. Lamprier, T. Gisselbrecht, and P. Gallinari, "Profile-based bandit with unknown profiles," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2060–2099, 2018.

[18] S.-Y. Yun, J. H. Nam, S. Mo, and J. Shin, "Contextual multi-armed bandits under feature uncertainty," *arXiv preprint arXiv:1703.01347*, 2017.

[19] A. Badanidiyuru, J. Langford, and A. Slivkins, "Resourceful contextual bandits," in *Conference on Learning Theory*, 2014, pp. 1109–1134.

[20] D. Russo and B. Van Roy, "Learning to optimize via posterior sampling," *Mathematics of Operations Research*, vol. 39, no. 4, pp. 1221–1243, 2014.

[21] S. Daulton, S. Singh, V. Avadhanula, D. Dimmery, and E. Bakshy, "Thompson sampling for contextual bandit problems with auxiliary safety constraints," *arXiv:1911.00638*, 2019.

[22] S. Agrawal and N. R. Devanur, "Bandits with concave rewards and convex knapsacks," in *Proceedings of the fifteenth ACM conference on Economics and computation*, 2014, pp. 989–1006.

[23] B. A. McFarland, N. AlKhalifah, M. Bohn, J. Bubert, E. S. Buckler, I. Ciampitti, J. Edwards, D. Ertl, J. L. Gage, C. M. Falcon, *et al.*, "Maize genomes to fields (g2f): 2014–2017 field seasons: genotype, phenotype, climatic, soil, and inbred ear image datasets," *BMC research notes*, vol. 13, no. 1, pp. 1–6, 2020.

[24] A. Holzkämper, P. Calanca, and J. Fuhrer, "Identifying climatic limitations to grain maize yield potentials using a suitability evaluation approach," *Agricultural and Forest Meteorology*, vol. 168, pp. 149–159, 2013.

[25] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, 2009.