

# A One-Size-Fits-All Solution to Conservative Bandit Problems

Yihan Du,<sup>1</sup> Siwei Wang,<sup>1</sup> Longbo Huang<sup>1</sup>

<sup>1</sup> Tsinghua University

duyh18@mails.tsinghua.edu.cn, {wangsw2020, longbohuang}@tsinghua.edu.cn

## Abstract

In this paper, we study a family of conservative bandit problems (CBPs) with *sample-path* reward constraints, i.e., the learner’s reward performance must be at least as well as a given baseline at any time. We propose a general one-size-fits-all solution to CBPs and present its applications to three encompassed problems, i.e., conservative multi-armed bandits (CMAB), conservative linear bandits (CLB) and conservative contextual combinatorial bandits (CCCB). Different from previous works which consider high probability constraints on the expected reward, our algorithms guarantee sample-path constraints on the actual received reward, and achieve better theoretical guarantees ( $T$ -independent additive regrets instead of  $T$ -dependent) and empirical performance. Furthermore, we extend the results and consider a novel conservative *mean-variance* bandit problem (MVCBP), which measures the learning performance in both the expected reward and variability. We design a novel algorithm with  $O(1/T)$  normalized additive regrets ( $T$ -independent in the cumulative form) and validate this result through empirical evaluation.

## 1 Introduction

The multi-armed bandit (MAB) problem (Thompson 1933; Auer, Cesa-Bianchi, and Fischer 2002) is a classic online learning model that characterizes the exploration-exploitation trade-off in sequential decision making. While existing bandit algorithms achieve satisfactory regret bounds over the whole learning processes, they can perform wildly and lose much in the initial exploratory phase. This limitation has hindered their applications in real-world scenarios such as health sciences, marketing and finance, where it is important to guarantee safe and smooth algorithm behavior in initialization. Hence, studying bandit problems with safe (conservative) exploration contributes to solving this issue.

In this paper, we study the conservative bandit problems (CBPs) with *sample-path* reward constraints. Specifically, a learner is given a set of regular arms and a default arm. At each timestep, the learner chooses a regular arm or the default arm to play and receives a reward according to the played arm. The learning’s objective is to minimize the expected cumulative regret (equivalently, maximize the expected cumulative reward), while ensuring that the received

cumulative reward must stay above a fixed percentage of what one can obtain by always playing the default arm.

CBPs have extensive real-world applications including recommendation systems, company operation and finance. For instance, in finance, investors are offered various financial products including the fixed-income security such as bank deposit (default arm), and the fluctuating equity securities such as stocks (regular arms). While the fixed-income security is a safe and reasonable option, investors want to find better choices to earn higher returns. Meanwhile, compared to the returns they can obtain by simply depositing the money, investors do not want to lose too much when exploring other investment choices. CBPs provide an effective model for such exploration-exploitation trade-off with the safe exploration guarantees.

We propose a general one-size-fits-all solution GenCB for CBPs, and present its applications to three important CBP problems, i.e., conservative multi-armed bandits (CMAB), conservative linear bandits (CLB) and conservative contextual combinatorial bandits (CCCB). We provide theoretical analysis and empirical evaluations for these algorithms, and show that our algorithms outperform existing ones both theoretically and empirically. Table 1 presents the comparison of regret bounds between our algorithms and existing ones. In the table, each regret term contains two components, the first component incurred by regular arms and the second term due to playing the default arm. One can see that our algorithms possess better regret guarantees. Moreover, unlike existing algorithms that only provide high probability bounds with  $T$ -dependent conservative regrets, we not only obtain expected bounds but also have  $T$ -independent conservative regrets.

Our work distinguishes itself from previous conservative bandit works, e.g., (Wu et al. 2016; Kazerouni et al. 2017; Garcelon et al. 2020; Zhang, Li, and Liu 2019) in two aspects: (i) Previous works consider high probability guarantees on the expected reward. Such models cannot directly handle many risk-adverse tasks, e.g., a start-up does not wish to tolerate any failure probability to reach the basic earning under the debt, or an asset management company must perform better than the promised return. While one can choose a very small  $\delta$  in previous algorithms to provide high-probability guarantees, the  $\ln(1/\delta)$ -dependent regrets will boost accordingly. Instead, we focus on a certainty

Problem	Algorithm	Regret bound	Type
CMAB	GenCB-CMAB (ours)	$O(H \ln T + \frac{H}{\alpha} [\ln(\frac{H}{\alpha})]^2)$	E
CMAB	ConUCB (Wu et al. 2016)	$O(H \ln(\frac{T}{\delta}) + \sum_{i=1}^K \frac{1}{\alpha \max\{\Delta_i, \Delta_0 - \Delta_i\}} \ln(T/\delta))$	H
CMAB	Lower Bound (Wu et al. 2016)	$O(\max\{\frac{K}{\alpha}, \sqrt{KT}\})$	E
CLB	GenCB-CLB (ours)	$O(d \ln(T) \sqrt{T} + \frac{d^2}{\alpha} [\ln(\frac{d}{\alpha})]^2)$	E
CLB	CLUCB (Kazerouni et al. 2017)	$O(d \ln(\frac{T}{\delta}) \sqrt{T} + \frac{d^2}{\alpha} [\ln(\frac{d}{\alpha\delta})]^2)$	H
CLB	CLUCB2 (Garcelon et al. 2020)	$O(d \ln(\frac{T}{\delta}) \sqrt{T} + \frac{d^2}{\alpha^2} [\ln(\frac{d}{\alpha\delta})]^2)$	H
CCCB	GenCB-CCCB (ours)	$O(d \ln(KT) \sqrt{T} + \frac{(K+d)^2}{\alpha} [\ln(\frac{K+d}{\alpha})]^2)$	E
CCCB	C3UCB (Zhang, Li, and Liu 2019)	$O(d \ln(KT) \sqrt{T} + \frac{d}{\alpha} \sqrt{\frac{d}{K} \ln(\frac{K}{\delta} T)})$	H

Table 1. Comparison of regret bounds for CBPs. “Type” refers to the type of regret bounds. “E” and “H” denote the expected and high probability bounds, respectively. Here  $H = \sum_{i=1}^K \Delta_i^{-1}$ .  $d$  is the dimension in CLB and CCCB. For high probability bounds, the convention in the bandit literature is to choose  $\delta = 1/T$ . Note that our formulation focuses on a sample-path reward constraint, while the other results consider the constraints on the expected reward.

(sample-path) guarantee on the actual empirical reward. Doing so ensures safe exploration (our regret bounds do not contain  $\delta$ ) and better suits such tasks. (ii) Our problem formulation, solution and analysis offer a general framework for studying a family of CBPs, including CMAB (Wu et al. 2016), CLB (Kazerouni et al. 2017; Garcelon et al. 2020) and CCCB (Zhang, Li, and Liu 2019). Moreover, our algorithms achieve better theoretical and empirical performance than previous schemes.

We also extend our results to the mean-variance setting (Markowitz et al. 1952; Sani, Lazaric, and Munos 2012), called *conservative mean-variance bandit problem* (MV-CBP), which focuses on the balance between the expected reward and variability with safe exploration. Different from the typical CBPs which only consider the expected reward into learning performance, MV-CBP takes into account both the mean and variance of the arms, and is more suitable for practical tasks that are sensitive to reward fluctuations, e.g., clinical trials and finance. For example, many risk-adverse investors prefer stable assets (e.g., bonds) with satisfactory returns than volatile assets (e.g., derivatives) with high returns, and they do not want to suffer wild fluctuations when exploring different financial products.

Note that the mean-variance regret in MV-CBP (formally defined in Eq. (4) in next section) consists not only the gap of *mean-variance* (a combination of both measures) between the played arms and the optimal arm, but also an additional variance for playing arms with different means, called *exploration risk*, which requires alternative techniques beyond those in typical CBPs. To tackle this issue, we carefully adapt our solution and analysis for the CBPs and make non-trivial extensions. Our results offer new insight into algorithm design for mean-variance bandit problems.

Our contributions are summarized as follows.

- We study a family of CBPs with sample-path reward constraints, which encompasses previously studied

CMAB (Wu et al. 2016), CLB (Kazerouni et al. 2017; Garcelon et al. 2020) and CCCB (Kazerouni et al. 2017). We propose a general one-size-fits-all solution GenCB for CBPs, which can translate a standard bandit algorithm into a conservative bandit algorithm and achieve better ( $T$ -independent conservative regret rather than  $T$ -dependent) theoretical regret bounds than previous works in the three specific problems.

- We extend the conservative bandit formulation to a novel conservative mean-variance bandit setting, which characterizes the trade-off between the expected reward and variability. We propose an algorithm, MV-CUCB, and prove that it achieves an  $O(1/T)$  normalized additive regret for the extended problem.
- We conduct extensive experiments for the considered problems. The results match our theoretical bounds and demonstrate that our algorithms achieve the performance superiority compared to existing algorithms.

## 1.1 Related Work

**Conservative Bandit Literature.** Recently, there are several works (Wu et al. 2016; Kazerouni et al. 2017; Zhang, Li, and Liu 2019; Garcelon et al. 2020) studying bandit problems with conservative exploration constraints. Under the constraints on the expected rewards, (Wu et al. 2016) propose an algorithm ConUCB for CMAB. (Kazerouni et al. 2017) design an algorithm CLUCB for CLB and (Garcelon et al. 2020) further propose an improved algorithm CLUCB2. (Zhang, Li, and Liu 2019) present an algorithm C3UCB for CCCB. Under the stage-wise constraints, (Khezeli and Bitar 2020) restrict the expected reward at any timestep to stay above a given baseline. (Amani, Alizadeh, and Thrampoulidis 2019) confine the played arm at any timestep to stay in a given safe set. Under the interleaving constraint, (Katariya et al. 2019) require the chosen action at any timestep to perform better than the default ac-

tion when interleaving in the combinatorial semi-bandit setting. (Bubeck, Perchet, and Rigollet 2013) study the standard  $K$ -armed bandit problem with knowledge of the highest expected reward and the smallest gap, (Locatelli, Gutzeit, and Carpentier 2016) consider the thresholding pure exploration problem, and the settings and methods in both works are different from ours.

**Mean-variance Bandit Literature.** (Sani, Lazaric, and Munos 2012) open the mean-variance bandit literature which considers both the expected reward and variability into performance measures, and a series of follow-ups (Maillard 2013; Vakili, Boukouvalas, and Zhao 2019; Cardoso and Xu 2019) have emerged recently. To our best knowledge, this paper is the first to study the mean-variance bandit problem with conservative exploration.

## 2 Problem Formulation

In this section, we first review previous standard (non-conservative) bandit problems (SBPs) and then give the formulation of the Conservative Bandit Problems (CBPs).

**Standard Bandit Problems (SBPs).** In a standard bandit problem, a learner is given a set of arms  $\mathcal{X}$ , where each arm  $x \in \mathcal{X}$  has an unknown reward distribution in  $[0, 1]$  with mean of  $\mu_x$ . Each arm  $x$  at timestep  $t$  has a random reward  $r_{t,x} = \mu_x + \eta_{t,x}$ , where  $\eta_{t,x}$  is an independent random noise with respect to  $t$ . At each timestep  $t$ , the learner plays an arm  $x_t$  and only observes the reward  $r_{t,x_t}$  of the chosen arm. Let  $x_* = \operatorname{argmax}_{x \in \mathcal{X}} \mu_x$  denote the optimal arm. The learning performance over a time horizon  $T$  is measured by *expected cumulative regret*

$$\mathbb{E}[\mathcal{R}_T] = \mu_{x_*} T - \mathbb{E} \left[ \sum_{t=1}^T \mu_{x_t} \right] = \sum_{x \neq x_*} \mathbb{E}[N_x(T)] \Delta_x, \quad (1)$$

where  $\Delta_x = \mu_{x_*} - \mu_x$  and  $N_x(T)$  is the number of times arm  $x$  was played over time  $T$ . The regret characterizes the loss due to not always playing the optimal arm. The goal of standard bandit algorithms is to minimize Eq. (1).

**Conservative Bandit Problems (CBPs).** The CBPs provide an alternative default arm  $x_0$  to play. In this case, since playing  $x_0$  is a default (baseline) policy that the learner is familiar with, for ease of analysis we assume that  $x_0$  has a known constant reward  $0 < \mu_0 < \mu_{x_*}$  as previous works (Wu et al. 2016; Kazerouni et al. 2017; Zhang, Li, and Liu 2019) do.<sup>1</sup>

Then, during the learning process, the learner is required to ensure that the cumulative reward under the chosen policy is lower bounded by a fraction of the reward from always pulling the default arm. Specifically, given a parameter  $\alpha \in (0, 1)$ , for any timestep  $t$ , the learner's cumulative empirical reward should be least  $1 - \alpha$  fraction of the reward of always

<sup>1</sup>This assumption can be relaxed to that  $x_0$  has a random reward within a known interval  $[r_0^\ell, r_0^h]$  ( $r_0^\ell > 0$ ) by slightly changing the right-hand-side of the if statements in our algorithms, and our analysis procedure still works. While previous works can remove this assumption by estimating  $\mu_0$ , this is due to that their constraints are imposed on the expected reward.

playing  $x_0$ , i.e.,

$$\sum_{s=1}^t r_{s,x_s} \geq (1 - \alpha) \mu_0 t, \quad \forall t \in \{1, \dots, T\}. \quad (2)$$

Here  $\alpha$  controls the strictness of the constraint, i.e., how conservative we want the leaner to behave, and can be viewed as the weight we place on safety in exploration. The goal of conservative bandit algorithms is to minimize the expected cumulative regret (Eq. (1)) while satisfying the reward constraint (Eq. (2)).

We note that constraint (2) is a *sample-path* reward constraint, which is different from the high-probability constraints on the expected reward in prior works (Wu et al. 2016; Kazerouni et al. 2017; Garcelon et al. 2020; Zhang, Li, and Liu 2019). This setting is particularly useful when the practical tasks cannot tolerate higher losses than the baseline with certainty, e.g., health care and investment. On the other hand, it also imposes new challenges in algorithm design and regret analysis.

Our formulation is a general framework which encompasses various bandit problems from the prospective of conservative exploration. For example, in CMAB which studies a conservative version of the classic  $K$ -armed bandit problem (Thompson 1933; Auer, Cesa-Bianchi, and Fischer 2002; Agrawal and Goyal 2012),  $\mathcal{X} = [K]$  and  $\mu_x$  is an arbitrary value.<sup>2</sup> In CLB which considers the linear bandit problem (Dani, Hayes, and Kakade 2008; Abbasi-yadkori, Pál, and Szepesvári 2011) with conservative exploration,  $\mathcal{X}$  is a compact subset of  $\mathbb{R}^d$  and each arm  $x \in \mathbb{R}^d$  has an expected reward  $\mu_x = x^\top \theta^*$ , where  $\theta^* \in \mathbb{R}^d$  is an unknown parameter. In CCCB which investigates the contextual combinatorial bandit problem (Qin, Chen, and Zhu 2014) with the safe exploration requirement, there is a set of base arms  $[K]$  and  $\mathcal{X}$  is a collection of subsets of base arms, which represents certain combinatorial structure (e.g., matchings and paths). For each  $x \in \mathcal{X}$ ,  $\mu_x$  is associated with the expected rewards of its containing base arms. We will analyze the CBPs under specific bandit settings in the next section.

## 3 A General Solution to Conservative Bandits

In this section, we first present a general solution for CBPs, and its regret analysis. Then, we present its applications to three specific problems, i.e., CMAB, CLB and CCCB, and show that in all three cases, our algorithm achieves tighter bounds than existing algorithms.

Algorithm 1 illustrates the proposed solution to CBPs, called GenCB, which offers a general scheme for translating a standard non-conservative bandit algorithm  $\mathcal{A}_S$  into a conservative bandit algorithm. In the algorithm,  $m$  denotes the time horizon of  $\mathcal{A}_S$ , and the number of times we play the regular arms,  $r_S(t)$  denotes the cumulative reward from sampling regular arms, and  $N_0(t)$  denotes the number of times  $x_0$  is played up to time  $t$ .

The main idea of GenCB is to play regular arms as much as possible while ensuring the sample-path reward constraint

<sup>2</sup> $[K] \stackrel{\text{def}}{=} \{1, \dots, K\}$ .

---

**Algorithm 1:** General Solution to Conservative Bandits (GenCB)

---

**Input:** Standard bandit problem and algorithm  $\mathcal{A}_S$ , regular arms  $\mathcal{X}$ , default arm  $x_0$  with reward  $\mu_0$ , parameter  $\alpha$ .

```

1  $\forall t \geq 0, N_0(t) \leftarrow 0, r_S(t) \leftarrow 0, m \leftarrow 0;$ 
2 for  $t = 1, 2, \dots$  do
3   if  $r_S(t-1) + N_0(t-1)\mu_0 \geq (1-\alpha)\mu_0 t$  then
4      $m \leftarrow m + 1;$ 
5     Play an arm  $x_t$  according to  $\mathcal{A}_S$ , observe  $r_{t,x_t}$  and update the statistical information;
6      $N_0(t) \leftarrow N_0(t-1);$ 
7      $r_S(t) \leftarrow r_S(t-1) + r_{t,x_t};$ 
8   else
9     Play  $x_0$  and receive reward  $\mu_0$ ;
10     $N_0(t) \leftarrow N_0(t-1) + 1;$ 
11     $r_S(t) \leftarrow r_S(t-1);$ 

```

---

in the worst case, since playing the default arm cannot provide any information for identifying the optimal arm. At each time, GenCB checks if playing a regular arm can satisfy the sample-path reward constraint in the worst case (this pull feedbacks zero reward). If it can, we play a regular arm  $x_t$  according to  $\mathcal{A}_S$ , observe reward  $r_{t,x_t}$  and update the statistical information. Otherwise, we choose the default arm.

Different from previous conservative algorithms (Wu et al. 2016; Kazerouni et al. 2017; Zhang, Li, and Liu 2019), GenCB guarantees the constraint with certainty rather than with high probability, and GenCB uses the received cumulative reward rather than the lower confidence bound to check the constraint. Doing so makes our algorithm less conservative and boosts its empirical performance significantly (see Section 5 for empirical comparisons).

Next, we present the regret analysis for GenCB. Note that, the regret for CBPs can be decomposed into (i) the regret incurred by regular arms, and (ii) the regret due to playing the default arm, i.e., conservative regret. Since the analysis of the former is similar to that in SBPs, as in the conservative bandit literature (Wu et al. 2016; Kazerouni et al. 2017; Garcelon et al. 2020; Zhang, Li, and Liu 2019), we mainly focus the conservative regret. We remark that our analysis is different from those in prior works, and can be applied to several specific CBPs including CMAB, CLB and CCCB. We give the regret bound of GenCB as follows.

**Theorem 1.** *Given a standard bandit problem and a corresponding algorithm  $\mathcal{A}_S$  with regret  $\mathbb{E}[\mathcal{R}_T(\mathcal{A}_S)] \leq B(T)$ , GenCB (Algorithm 1) guarantees the sample-path reward constraint Eq. (2) and achieves a regret bound*

$$\mathbb{E}[\mathcal{R}_T(\text{GenCB})] \leq B(T) + C\Delta_0,$$

where  $C$  is a problem-specific constant independent of  $T$  and  $\Delta_0 = \mu_{x_*} - \mu_{x_0}$ .

*Proof.* First, it can be seen from the algorithm that the sample-path reward constraint Eq. (2) can be guaranteed. Next, we prove the regret bound of GenCB. We use  $\mathcal{S}_t$  to denote the set of timesteps up to time  $t$  during which we play regular arms and use  $m_t$  to denote its size. Let  $\tau$  denote the

last timestep we play  $x_0$ , i.e.,  $\tau$  is the last timestep such that  $r_S(\tau-1) + N_0(\tau-1)\mu_0 < (1-\alpha)\mu_0\tau$  holds. Rearranging the terms, and subtracting  $(1-\alpha)\mu_0N_0(\tau-1)$  from both sides (note that  $\tau = N_0(\tau-1) + m_{\tau-1} + 1$ ), we have

$$\begin{aligned} \alpha\mu_0N_0(\tau-1) &< (1-\alpha)\mu_0(m_{\tau-1}+1) - r_S(\tau-1) \\ &= (1-\alpha)\mu_0(m_{\tau-1}+1) \\ &\quad + \sum_{t \in \mathcal{S}_{\tau-1}} (\mu_{x_t} - r_{t,x_t}) - \sum_{t \in \mathcal{S}_{\tau-1}} \mu_{x_t}. \end{aligned} \quad (3)$$

$\sum_{t \in \mathcal{S}_{\tau-1}} (\mu_{x_t} - r_{t,x_t})$  is the deviation between the sum of  $m_{\tau-1}$  samples and their means. Using the Azuma-Hoeffding inequality, it can be upper bounded by  $F\sqrt{m_{\tau-1}\ln(1/\delta)}$  with probability at least  $1-\delta$ , for some constant  $F$  that varies in different settings. By setting  $\delta = 1/\tau$ , we can obtain an expected upper bound as  $F\sqrt{m_{\tau-1}\ln\tau} + 1$ . Taking expectation on both sides of Eq. (3), setting  $m = \mathbb{E}[m_{\tau-1}+1]$  and replacing  $\mathbb{E}[\sum_{t \in \mathcal{S}_{\tau-1}} \mu_{x_t}]$  with  $\mu_{x_*}\mathbb{E}[m_{\tau-1}] - \mathbb{E}[\mathcal{R}_{m_{\tau-1}}(\mathcal{A}_S)]$ , we have

$$\begin{aligned} \alpha\mu_0\mathbb{E}[N_0(\tau-1)] &< -(\Delta_0 + \alpha\mu_0)m + \mathbb{E}[B(m_{\tau-1})] \\ &\quad + \mathbb{E}[F\sqrt{m_{\tau-1}\ln(\tau)}] + 1 \\ &\stackrel{(a)}{<} -(\Delta_0 + \alpha\mu_0)m + B(m) + 2 \\ &\quad + F\sqrt{m\ln(\mathbb{E}[N_0(\tau-1)] + m)}, \end{aligned}$$

where (a) comes from Jensen's inequality. Note that since  $B(m)$  and  $F\sqrt{m\ln(\mathbb{E}[N_0(\tau-1)] + m)}$  are sublinear with respect to  $m$ , for any  $m \geq 2$ , the right-hand-side can be upper bounded by  $G[\ln(\sqrt{\mathbb{E}[N_0(\tau-1)]})]^2$  where  $G$  is a constant factor that only depends on problem parameters. Then, we obtain  $\mathbb{E}[N_0(\tau-1)] \leq \frac{G}{\alpha\mu_0}[\ln(\frac{G}{\alpha\mu_0})]^2$ . Thus,  $\mathbb{E}[N_0(T)] = \mathbb{E}[N_0(\tau)] = \mathbb{E}[N_0(\tau-1)] + 1 \leq C$ , where  $C \triangleq \frac{G}{\alpha\mu_0}[\ln(\frac{G}{\alpha\mu_0})]^2 + 1$  is independent of  $T$ .

Combining the regrets for  $\mathcal{A}_S$  and  $x_0$ , we obtain that  $\mathbb{E}[\mathcal{R}_T(\text{GenCB})] \leq B(T) + C\Delta_0$ .  $\square$

**Remark 1.** Theorem 1 shows that GenCB provides a general algorithmic and analytical framework for translating a standard bandit problem into a conservative bandit algorithm, and only generate an additional  $T$ -independent regret due to the reward constraint. To the best of our knowledge, this is the first general analysis procedure which works for a family of CBPs with sample-path reward constraints, and it provides an expected regret bound (rather than high probability bounds in (Wu et al. 2016; Kazerouni et al. 2017; Garcelon et al. 2020; Zhang, Li, and Liu 2019)) with  $T$ -independent conservative regret.

Below, we apply GenCB to three widely studied CBPs, i.e., CMAB, CLB and CCCB. Here we only present the main theorems, and defer the algorithm pseudo-codes and proofs to the supplementary material (Du, Wang, and Huang 2020).

### 3.1 Application to Conservative Multi-Armed Bandits (CMAB)

The conservative multi-armed bandit (CMAB) problem is a variation of the classic  $K$ -armed bandit model with conservative exploration (Wu et al. 2016), which has extensive

applications including clinical trials, online advertising and wireless network. In CMAB,  $\mathcal{X} = [K]$  and  $\mu_i$  ( $1 \leq i \leq K$ ) can be an arbitrary value. Without loss of generality, we assume  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$  and denote  $\mu_* \triangleq \mu_1$ .

We apply the GenCB algorithm with the UCB algorithm (Auer, Cesa-Bianchi, and Fischer 2002) to this setting, by replacing Line 5 in Algorithm 1 with  $x_t \leftarrow \text{argmax}_{i \in [K]} \left( \hat{\mu}_i + \sqrt{2 \ln m / N_i(t-1)} \right)$ , where  $\hat{\mu}_i$  is the reward empirical mean for arm  $i$ , and name this version of the algorithm GenCB-CMAB.

The main idea of GenCB-CMAB is to play the arm with the maximum upper confidence bound whenever the reward constraint is satisfied (otherwise we play the default arm). The regret bound for GenCB-CMAB is summarized below.

**Theorem 2.** *For the conservative multi-armed bandit problem, GenCB-CMAB guarantees the sample-path reward constraint Eq. (2) and achieves the regret bound*

$$O\left(H \ln T + \frac{H \Delta_0}{\alpha \mu_0 (\Delta_0 + \alpha \mu_0)} \left[ \ln \left( \frac{H}{\alpha \mu_0 (\Delta_0 + \alpha \mu_0)} \right) \right]^2\right),$$

where  $H = \sum_{i>1} \Delta_i^{-1}$ .

**Remark 2.** The first term owes to playing the regular arms, which is similar to the result in standard MAB (Auer, Cesa-Bianchi, and Fischer 2002), and the second term is caused by the default arm, i.e., the conservative regret, which is the main focus in conservative bandit study. Compared to the existing algorithm ConUCB (Wu et al. 2016), GenCB-CMAB only incurs a  $T$ -independent conservative regret rather than  $\ln T$  (see Table 1). Our result also matches the regret lower bound derived in (Wu et al. 2016) for CMAB with expected reward constraints, which also holds for our sample-path reward constraint setting.

### 3.2 Application to Conservative Linear Bandits (CLB)

The conservative linear bandit (CLB) (Kazerouni et al. 2017; Garcelon et al. 2020) problem considers the linear bandit problem (Dani, Hayes, and Kakade 2008; Abbasi-yadkori, Pál, and Szepesvári 2011) with safe exploration. In CLB where there is a linear structure among arms,  $\mathcal{X}$  is a compact subset of  $\mathbb{R}^d$  and  $\mu_x = x^\top \theta^*$ , where  $\theta^* \in \mathbb{R}^d$  is an unknown parameter. We make the common assumptions, i.e.,  $\|x\|_2 \leq L, \forall x \in \mathcal{X}$  and  $\|\theta^*\|_2 \leq S$ , as previous linear bandit papers (Dani, Hayes, and Kakade 2008; Abbasi-yadkori, Pál, and Szepesvári 2011; Kazerouni et al. 2017) do.

For CLB, we apply GenCB with the LinUCB algorithm (Abbasi-yadkori, Pál, and Szepesvári 2011) by replacing Line 5 in Algorithm 1 with  $(x_t, \hat{\theta}_t) \leftarrow \text{argmax}_{(x, \theta) \in \mathcal{X} \times \mathcal{C}_t} x^\top \theta$ . Here  $\mathcal{C}_t = \{\theta \in \mathbb{R}^d : \|\theta - \hat{\theta}_{t-1}\|_{V_{t-1}} \leq \sqrt{d \ln(2m^2(1+mL^2/\lambda))} + \sqrt{\lambda}S\}$  is a confidence ellipsoid that contains  $\theta^*$  with high probability, and we define  $\hat{\theta}_t = V_t^{-1} b_t$ ,  $V_t = \lambda I + \sum_{s=1}^t x_s x_s^\top$ ,  $b_t = \sum_{s=1}^t r_{s,x_s} x_s$  and  $\lambda \geq \max\{1, L^2\}$ .<sup>3</sup> We name this version of the algorithm GenCB-CLB, whose key idea is to play a

<sup>3</sup>  $\|x\|_V \stackrel{\text{def}}{=} \sqrt{x^\top V x}, \forall x \in \mathbb{R}^d, \forall V \in \mathbb{R}^{d \times d}$ .

regular arm according to the optimism in the face of uncertainty principle while ensuring the sample-path reward constraint. Below, we have the regret bound of GenCB-CLB.

**Theorem 3.** *For the conservative linear bandit problem, GenCB-CLB guarantees the sample-path reward constraint Eq. (2) and has the regret bound*

$$O\left(d \ln \left( \frac{LT}{\lambda} \right) \sqrt{T} + \frac{d^2 S^2 \lambda \Delta_0}{\alpha \mu_0 \tilde{\Delta}_0} \left[ \ln \left( \frac{dS\sqrt{\lambda}}{\alpha \mu_0 \tilde{\Delta}_0} \right) \right]^2\right),$$

where  $\tilde{\Delta}_0 = \Delta_0 + \alpha \mu_0$ .

**Remark 3.** Similarly, the first term is aligned with the result in standard linear bandits (Dani, Hayes, and Kakade 2008; Abbasi-yadkori, Pál, and Szepesvári 2011), and the second term is the conservative regret due to the default arm. While the existing algorithms CLUCB (Kazerouni et al. 2017) and CLUCB2 (Garcelon et al. 2020) have  $\ln(1/\delta)$ -dependent conservative regrets with high probability (do not contain  $T$  either), these results are of  $\ln T$  order when making the convention  $\delta = 1/T$ . In contrast, we provide an expected bound with a  $T$ -independent conservative regret.

### 3.3 Application to Conservative Contextual Combinatorial Bandits (CCCB)

The conservative contextual combinatorial bandit (CCCB) problem (Zhang, Li, and Liu 2019) investigates the contextual combinatorial bandit problem under the safe exploration requirement. In CCCB,  $\mathcal{X}$  is a collection of subsets of *base arms*  $x_1, \dots, x_K \in \mathbb{R}^d$  and generated from certain combinatorial structure (e.g., matchings and paths). The learner plays a *super arm* (subset of base arms)  $A_t \in \mathcal{X}$  or the default arm  $x_0$  at each timestep. The expected reward of base arm  $x_e$  is  $w_e^* = x_e^\top \theta^*$  and that of super arm  $A$  is  $f(A, \mathbf{w}^*)$ , where  $\theta^*$  is an unknown parameter and  $f$  satisfies two mild assumptions, i.e., monotonicity and Lipschitz continuous with parameter  $P$  (Qin, Chen, and Zhu 2014; Zhang, Li, and Liu 2019). Similar to CLB, we assume  $\|x\|_2 \leq L, \forall x \in \mathcal{X}$  and  $\|\theta^*\|_2 \leq S$ . At timestep  $t$ , the random reward of a base arm  $x_e$  and a super arm  $A$  are  $w_{t,e} = w_e^* + \eta_{t,e} \in [0, 1]$  and  $r_{t,A} = f(A, \mathbf{w}^*) + \eta_{t,A} \in [0, K]$ , respectively. After pulling super arm  $A_t$ , we receive the random reward  $r_{t,A_t}$  and observe a semi-bandit feedback, i.e.,  $w_{t,e}$  for each  $e \in A_t$ .

For CCCB, we apply GenCB with the C2UCB algorithm (Qin, Chen, and Zhu 2014), by replacing Line 5 in Algorithm 1 with  $A_t \leftarrow \text{argmax}_{A \in \mathcal{X}} f(A, \bar{\mathbf{w}}_t)$ . Here  $\bar{\mathbf{w}}_{t,e} = x_e^\top \hat{\theta}_{t-1} + (\sqrt{d \ln(2m^2(1+mKL^2/\lambda))} + \sqrt{\lambda}S) \|x_e\|_{V_{t-1}^{-1}}$  is the upper confidence bound of  $w_e^*$ , and we define  $\hat{\theta}_t = V_t^{-1} b_t$ ,  $V_t = \lambda I + \sum_{s=1}^t \sum_{e \in A_s} x_e x_e^\top$ ,  $b_t = \sum_{s=1}^t \sum_{e \in A_s} w_{s,e} x_e$  and  $\lambda \geq \max\{1, L^2\}$ . The key idea here is to play a super arm with the maximum upper confidence bound according to the historical observations on base arms. Theorem 4 below gives the regret bound of GenCB-CCCB.

**Theorem 4.** *For the contextual combinatorial bandit problem, GenCB-CCCB ensures the sample-path reward con-*

---

**Algorithm 2:** MV-CUCB

---

**Input:** Reugular arms  $[K]$ , default arm  $x_0$  with  $\text{MV}_0 = \rho\mu_0$ , parameters  $\alpha, \rho > \frac{2}{\alpha\mu_0}$ .

- 1  $\forall t \geq 0, \forall 0 \leq i \leq K, N_i(t) \leftarrow 0. m \leftarrow 0.$
- 2  $\widehat{\text{MV}}_0(\mathcal{A}) \leftarrow 0;$
- 3 **for**  $t = 1, 2, \dots$  **do**
- 4   **if**  $(t-1)\widehat{\text{MV}}_{t-1}(\mathcal{A}) - 2 \geq (1-\alpha)\text{MV}_0 t$  **then**
- 5      $m \leftarrow m + 1;$
- 6      $x_t \leftarrow \operatorname{argmax}_{i \in [K]} \left( \widehat{\text{MV}}_i + (5+\rho) \sqrt{\frac{\ln(12Km^3)}{2N_i(t-1)}} \right);$
- 7     Pull arm  $x_t$ , observe the random reward  $r_{t,x_t}$
- 8     and update  $\widehat{\text{MV}}_{x_t};$
- 9      $N_{x_t}(t) \leftarrow N_{x_t}(t-1) + 1$  and
- 10      $\forall 0 \leq i \leq K, i \neq x_t, N_i(t) \leftarrow N_i(t-1);$
- 11 **else**
- 12     Play  $x_0$  and receive reward  $\mu_0;$
- 13      $N_0(t) \leftarrow N_0(t-1) + 1$  and
- 14      $\forall 1 \leq i \leq K, N_i(t) \leftarrow N_i(t-1);$

---

straint Eq. (2) and achieves the regret bound

$$O\left(Pd \ln\left(\frac{KLT}{\lambda}\right) \sqrt{T} + \frac{D^2}{\alpha\mu_0\tilde{\Delta}_0} \left[ \ln\left(\frac{D}{\alpha\mu_0\tilde{\Delta}_0}\right) \right]^2\right),$$

where  $D = K + P\sqrt{\lambda}Sd$  and  $\tilde{\Delta}_0 = \Delta_0 + \alpha\mu_0$ .

**Remark 4.** The first term is consistent with the result in standard contextual combinatorial bandits (Qin, Chen, and Zhu 2014), and the second conservative regret term is due to playing the default arm. Compared to the state-of-the-art algorithm C3UCB (Zhang, Li, and Liu 2019), GenCB-CCCB provides a  $T$ -independent conservative regret, while C3UCB incurs a  $\ln T$  regret (see Table 1).

## 4 Conservative Mean-Variance Bandits

We now extend CBPs to the mean-variance (Sani, Lazaric, and Munos 2012; Maillard 2013; Cardoso and Xu 2019) setting (MV-CBP), which focuses on finding arms that achieve effective trade-off between the expected reward and variability. MV-CBP increments the typical conservative bandit model and better suits the tasks emphasizing on reward fluctuations. It also brings additional complications for algorithm design and regret analysis beyond GenCB.

### 4.1 Problem Formulation for MV-CBP

To introduce our MV-CBP formulation, we first review the standard mean-variance bandit setting (Sani, Lazaric, and Munos 2012). Each arm  $x \in [K]$  is associated with a measure mean-variance, which is formally defined as  $\text{MV}_x = \rho\mu_x - \sigma_x^2$ , where  $\sigma_x^2$  is the reward variance and  $\rho$  is a weight parameter. Let  $x_*^{\text{MV}} = \operatorname{argmax}_{x \in [K]} \text{MV}_x$  denote the mean-variance optimal arm. Given i.i.d. reward samples  $\{Z_{x,s}\}_{s=1}^t$  of arm  $x$ , we define the empirical mean-variance  $\widehat{\text{MV}}_{x,t} = \rho\hat{\mu}_{x,t} - \hat{\sigma}_{x,t}^2$ , where  $\hat{\mu}_{x,t} = \frac{1}{t} \sum_{s=1}^t Z_{x,s}$  and  $\hat{\sigma}_{x,t}^2 = \frac{1}{t} \sum_{s=1}^t (Z_{x,s} - \hat{\mu}_{x,t})^2$ .

For an algorithm  $\mathcal{A}$  and its sample path  $\{r_{t,x_t}\}_{t=1}^T$  over time horizon  $T$ , we define the empirical mean-variance  $\widehat{\text{MV}}_T(\mathcal{A}) = \rho\hat{\mu}_T(\mathcal{A}) - \hat{\sigma}_T^2(\mathcal{A})$ , where  $\hat{\mu}_T(\mathcal{A}) = \frac{1}{T} \sum_{t=1}^T r_{t,x_t}$  and  $\hat{\sigma}_T^2(\mathcal{A}) = \frac{1}{T} \sum_{t=1}^T (r_{t,x_t} - \hat{\mu}_T(\mathcal{A}))^2$ . Naturally, for algorithm  $\mathcal{A}$  over time  $T$ , we define the mean-variance regret  $\mathcal{R}_T^{\text{MV}}(\mathcal{A}) = \widehat{\text{MV}}_{x^*,T} - \widehat{\text{MV}}_T(\mathcal{A})$ , which is the difference of the mean-variance performance between  $\mathcal{A}$  and what we could have achieved by always playing  $x_*^{\text{MV}}$ .

Due to the difficulty of the  $\mathcal{R}_T^{\text{MV}}(\mathcal{A})$  metric, we follow the mean-variance bandit literature and use a more tractable mesure *mean-variance pseudo-regret* (Sani, Lazaric, and Munos 2012) defined as:

$$\widetilde{\mathcal{R}}_T^{\text{MV}}(\mathcal{A}) = \frac{1}{T} \sum_{x \neq x^*} N_{x,T} \Delta_x^{\text{MV}} + \frac{2}{T^2} \sum_{x \in \mathcal{X}} \sum_{y \neq x} N_{x,T} N_{y,T} \Gamma_{x,y}^2, \quad (4)$$

where  $N_{x,T}$  is a shorthand for  $N_x(T)$ ,  $\Delta_x^{\text{MV}} = \widehat{\text{MV}}_{x^*} - \widehat{\text{MV}}_x$  and  $\Gamma_{x,y} = \mu_x - \mu_y$ . It has been shown that any bound on  $\widetilde{\mathcal{R}}_T^{\text{MV}}(\mathcal{A})$  immediately translates into an bound on  $\mathcal{R}_T^{\text{MV}}(\mathcal{A})$  (Lemma 1 in (Sani, Lazaric, and Munos 2012)). Thus, most theoretical analysis (Sani, Lazaric, and Munos 2012; Maillard 2013; Cardoso and Xu 2019) on mean-variance bandits has been done via  $\widetilde{\mathcal{R}}_T^{\text{MV}}(\mathcal{A})$ . Note that, in MV-CBP the measures  $\widehat{\text{MV}}_T(\mathcal{A})$  and  $\mathcal{R}_T^{\text{MV}}(\mathcal{A})$  are both *normalized* quantities over  $T$ .

In addition to minimizing the regret, the learner is also required to guarantee the following mean-variance constraint:

$$\widehat{\text{MV}}_t(\mathcal{A}) \geq (1-\alpha)\text{MV}_0, \quad \forall t \in \{1, \dots, T\}. \quad (5)$$

Here  $\text{MV}_0$  denotes the mean-variance of our default arm  $x_0$  with known constant reward  $\mu_0$  and zero variance. The goal in MV-CBP is to minimize Eq. (4) while satisfying Eq. (5).

### 4.2 Algorithm for MV-CBP

We propose a novel algorithm named MV-CUCB for MV-CBP (illustrated in Algorithm 2). The main idea is to compute the upper confidence bound of mean-variance for each arm and select one according to the optimism principle whenever the constraint is not violated. Theorem 5 summarizes the performance results of MV-CUCB (see the supplementary material (Du, Wang, and Huang 2020) for its proof).

**Theorem 5.** For the conservative mean-variance multi-armed bandit problem with  $\alpha\text{MV}_0 > 2$ , MV-CUCB (Algorithm 2) ensures the mean-variance constraint Eq. (5) and achieves the following regret bound:<sup>4</sup>

$$\begin{aligned} \tilde{O}\left(\frac{\rho^2 \ln(KT)}{T} \left( H_1 + H_2 + \frac{\rho^2 \ln(KT)}{T} H_3 \right) \right. \\ \left. + \frac{\rho^3 \sqrt{K} (H_1^{\text{MV}} + 4H_2^{\text{MV}}) + (\rho^4 K H_3^{\text{MV}} + \rho K) \tilde{\Delta}_0^{\text{MV}}}{(\alpha\text{MV}_0 - 2)\tilde{\Delta}_0^{\text{MV}} T} \cdot \Delta_0^{\text{MV}} \right), \end{aligned}$$

where  $H_1^{\text{MV}} = \sum_{i>1} (\Delta_i^{\text{MV}})^{-1}$ ,  $H_2^{\text{MV}} = \sum_{i>1} (\Delta_i^{\text{MV}})^{-2}$ ,  $H_3^{\text{MV}} = \sum_{i>1} \sum_{j>1, j \neq i} (\Delta_i^{\text{MV}} \Delta_j^{\text{MV}})^{-2}$  and  $\tilde{\Delta}_0^{\text{MV}} = \Delta_0^{\text{MV}} + \alpha\text{MV}_0$ .

---

<sup>4</sup> $\tilde{O}$  omits the logarithmic terms that are independent of  $T$ .

**Remark 5.** Since a pull of  $x_0$  not only accumulates  $MV_0$  but also causes an exploration risk (bounded by 2 for reward distributions in  $[0, 1]$ ) due to the switch between different arms, we need the mild assumption  $\alpha MV_0 > 2$  to guarantee that a pull of  $x_0$  will not violate the constraint. Recall that the result is a normalized regret over  $T$ , the first term of  $O(\ln T/T)$  order owes to regular arms, which agrees with the previous non-conservative mean-variance result (Sani, Lazaric, and Munos 2012). The second term is the conservative regret for satisfying the constraint, which is of only  $O(1/T)$  order and independent of  $T$  in the cumulative form. To our best knowledge, Theorem 5 is the first result for conservative bandits with mean-variance objectives.

## 5 Experiments

We conduct experiments for our algorithms in four problems, i.e., CMAB, CLB, CCCB and MV-CBP, with a wide range of parameter settings. Due to space limit, only partial results are presented here (see the supplementary material (Du, Wang, and Huang 2020) for full results).

In all experiments, we assume the rewards to take i.i.d. Bernoulli values. For CMAB, we set  $K \in \{24, 72, 144\}$ ,  $\alpha \in \{0.05, 0.1, 0.15\}$ ,  $\mu_0 = 0.7$  and  $\mu_1, \dots, \mu_K$  as an arithmetic sequence from 0.8 to 0.2. For CLB and CCCB, we set  $d \in \{5, 7, 9\}$ ,  $\alpha \in \{0.01, 0.02, 0.03\}$ ,  $K = 2d$  and  $f(A, \mathbf{w}^*) = \sum_{e \in A} w_e^*$ . For MV-CBP, we use the same parameter settings as CMAB and additionally set  $\rho \in \{10, 30, 60\}$ . For each algorithm, we perform 50 independent runs and present the average (middle curve), maximum (upper curve) and minimum (bottom curve) cumulative regrets across runs. For each figure, we also zoom in the initial exploratory phase in the sub-figure to compare algorithm performance in this phase.

**Experiments for CBPs.** In the experiments for CMAB (Figure 1(a)), CLB (Figure 1(b)) and CCCB (Figure 1(c)), we compare GenCB-CMAB, GenCB-CLB and GenCB-CCCB to previous CBP algorithms CUCB (Wu et al. 2016), CLUCB (Kazerouni et al. 2017) and C3UCB (Zhang, Li, and Liu 2019), the standard bandit algorithms UCB (Auer, Cesa-Bianchi, and Fischer 2002), LinUCB (Abbasi-yadkori, Pál, and Szepesvári 2011) and C2UCB (Qin, Chen, and Zhu 2014), and the conservative baseline  $(1-\alpha)\mu_0$ , respectively.

We see that, in the exploration phase, existing non-conservative algorithms suffer higher losses than the baseline, while our algorithms and previous CBP algorithms achieve similar performance as (or better than) the baseline due to the conservative constraints. However, since previous CBP algorithms use lower confidence bounds (rather than the empirical rewards in ours) to check the constraints, they are forced to play the default arm more and act more conservatively compared to ours.

In the exploitation phase, when compared to non-conservative algorithms, our algorithms have additional regrets that keep constant as  $T$  increases, which matches our  $T$ -independent conservative regret bounds. Compared to previous CBP algorithms, our schemes achieve significantly better performance, since we play the default arm less and enjoy a lower conservative regret.

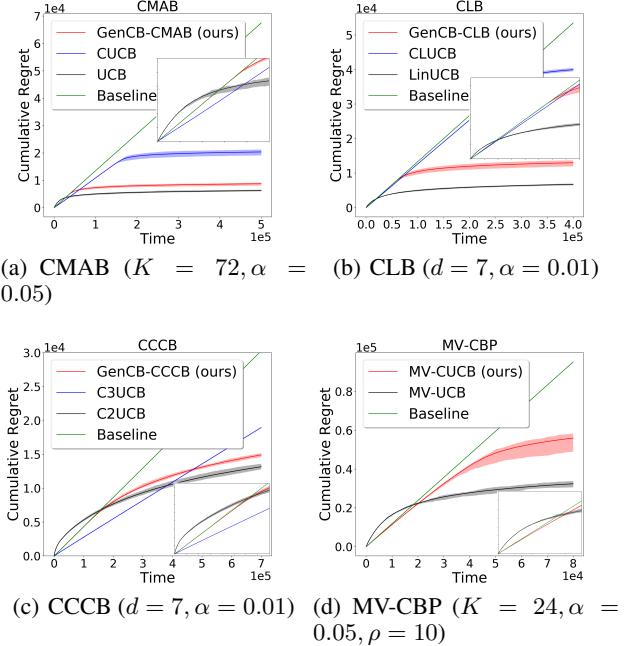


Figure 1. Experiments for the studied problems, i.e., CMAB, CLB, CCCB and MV-CBP.

**Experiments for MV-CBP.** In the experiments for MV-CBP (Figure 1(d)), we present the mean-variance regret in the cumulative form  $T \cdot \bar{\mathcal{R}}_T^{MV}$  for clarity of comparison. Since MV-CUCB is the first algorithm for MV-CBP, we compare it with the standard mean-variance bandit algorithm MV-UCB and the baseline  $(1-\alpha)MV_0T$ . We can see that, in the exploration phase, MV-UCB suffers from a higher regret than the baseline while MV-CUCB follows the baseline closely. One also sees that MV-CUCB achieves this with only an additional constant overall regret compared to MV-UCB, which matches our  $T$ -independent bound of conservative regret.

## 6 Conclusion and Future Works

In this paper, we propose a general solution to a family of conservative bandit problems (CBPs) with sample-path reward constraints, and present its applications to three encompassed problems, i.e., conservative multi-armed bandits (CMAB), conservative linear bandits (CLB) and conservative contextual combinatorial bandits (CCCB). We show that our algorithms outperform existing ones both theoretically (incurs  $T$ -independent conservative regrets rather than  $T$ -dependent) and empirically. Moreover, we study a novel extension of CBPs to the mean-variance setting (MV-CBP) and develop an algorithm with  $O(1/T)$  normalized conservative regret ( $T$ -independent in the cumulative form). We also validate this result through empirical evaluation.

There are several directions worth further investigation. One is to consider more general conservative mean-variance bandits other than the  $K$ -armed setting, e.g., a contextual extension. Another direction is to consider other practical conservative constraints which capture the safe exploration requirement in real-world applications.

## Acknowledgements

This work is supported in part by the National Natural Science Foundation of China Grant 61672316, the Zhongguancun Haishu Institute for Frontier Information Technology and the Turing AI Institute of Nanjing.

## Ethical Impact

In this paper, we study a family of conservative bandit problems and present algorithms with theoretical guarantees and experimental results. While our work mainly focuses on the theoretical analysis, it may have potential social impacts on the applications including finance and clinical trials. For example, our algorithms may help risk-adverse investors choose financial products, with the objective of obtaining high cumulative returns while guaranteeing a certain baseline during exploration. We believe that this work does not involve any ethical issue.

## References

- Abbasi-yadkori, Y.; Pál, D.; and Szepesvári, C. 2011. Improved Algorithms for Linear Stochastic Bandits. In *Advances in Neural Information Processing Systems*, 2312–2320.
- Agrawal, S.; and Goyal, N. 2012. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, 39–1.
- Amani, S.; Alizadeh, M.; and Thrampoulidis, C. 2019. Linear Stochastic Bandits Under Safety Constraints. In *Advances in Neural Information Processing Systems*, 9256–9266.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine Learning* 47(2-3): 235–256.
- Bubeck, S.; Perchet, V.; and Rigollet, P. 2013. Bounded regret in stochastic multi-armed bandits. In *Conference on Learning Theory*, 122–134.
- Cardoso, A. R.; and Xu, H. 2019. Risk-averse stochastic convex bandit. In *International Conference on Artificial Intelligence and Statistics*, 39–47.
- Dani, V.; Hayes, T.; and Kakade, S. M. 2008. Stochastic Linear Optimization under Bandit Feedback. In *Conference on Learning Theory*.
- Du, Y.; Wang, S.; and Huang, L. 2020. A One-Size-Fits-All Solution to Conservative Bandit Problems. volume abs/2012.07341. URL <https://arxiv.org/abs/2012.07341>.
- Garcelon, E.; Ghavamzadeh, M.; Lazaric, A.; and Pirotta, M. 2020. Improved Algorithms for Conservative Exploration in Bandits. In *AAAI Conference on Artificial Intelligence*.
- Katariya, S.; Kveton, B.; Wen, Z.; and Potluru, V. 2019. Conservative Exploration using Interleaving. In *International Conference on Artificial Intelligence and Statistics*.
- Kazerouni, A.; Ghavamzadeh, M.; Yadkori, Y. A.; and Van Roy, B. 2017. Conservative contextual linear bandits. In *Advances in Neural Information Processing Systems*, 3910–3919.
- Khezeli, K.; and Bitar, E. 2020. Safe Linear Stochastic Bandits. In *AAAI Conference on Artificial Intelligence*.
- Locatelli, A.; Gutzeit, M.; and Carpentier, A. 2016. An optimal algorithm for the Thresholding Bandit Problem. In *International Conference on Machine Learning*, 1690–1698.
- Maillard, O.-A. 2013. Robust risk-averse stochastic multi-armed bandits. In *International Conference on Algorithmic Learning Theory*, 218–233. Springer.
- Markowitz, H. M.; et al. 1952. Portfolio Selection. *Journal of Finance* 7(1): 77–91.
- Qin, L.; Chen, S.; and Zhu, X. 2014. Contextual Combinatorial Bandit and its Application on Diversified Online Recommendation. In *International Conference on Data Mining*, 461–469.
- Sani, A.; Lazaric, A.; and Munos, R. 2012. Risk-aversion in multi-armed bandits. In *Advances in Neural Information Processing Systems*, 3275–3283.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25(3/4): 285–294.
- Vakili, S.; Boukouvalas, A.; and Zhao, Q. 2019. Decision variance in risk-averse online learning. In *Conference on Decision and Control*, 2738–2744. IEEE.
- Wu, Y.; Shariff, R.; Lattimore, T.; and Szepesvári, C. 2016. Conservative bandits. In *International Conference on Machine Learning*, 1254–1262.
- Zhang, X.; Li, S.; and Liu, W. 2019. Contextual Combinatorial Conservative Bandits. *arXiv preprint:1911.11337*.