# Budget-Constrained Bandits over General Cost and Reward Distributions

**Semih Cayci**
ECE, The Ohio State University

**Atilla Eryilmaz**
ECE, The Ohio State University

**R. Srikant**
CSL and ECE, UIUC

## Abstract

We consider a budget-constrained bandit problem where each arm pull incurs a random cost, and yields a random reward in return. The objective is to maximize the total expected reward under a budget constraint on the total cost. The model is general in the sense that it allows correlated and potentially heavy-tailed cost-reward pairs that can take on negative values as required by many applications. We show that if moments of order $(2+\gamma)$ for some $\gamma > 0$ exist for all cost-reward pairs, $O(\log B)$ regret is achievable for a budget $B > 0$. In order to achieve tight regret bounds, we propose algorithms that exploit the correlation between the cost and reward of each arm by extracting the common information via linear minimum mean-square error estimation. We prove a regret lower bound for this problem, and show that the proposed algorithms achieve tight problem-dependent regret bounds, which are optimal up to a universal constant factor in the case of jointly Gaussian cost and reward pairs.

## 1 Introduction

Multi-armed bandit problem (MAB) has been the prominent model for the exploration-and-exploitation dilemma since its introduction (Robbins, 1952; Lai and Robbins, 1985; Berry and Fristedt, 1985). Due to the universality of the dilemma, bandit algorithms have found a broad area of applications from medical trials and dynamic pricing to ad allocation. As a common feature of all MAB instances, each action depletes a cost from a limited budget, and a random reward

is obtained in return. In such a setting, the aim of the decision maker is to balance the exploration and exploitation at every step so as to maximize the cumulative reward until depleting the budget. In the classical MAB setting, each action is assumed to consume a known deterministic amount of resource, i.e., one time-slot. However, in many problems of interest, different tasks consume different and random amount of resources, which can be unbounded and potentially correlated with the reward. The applications of this extended setting include routing in communications and task scheduling in computing systems, where the controller sequentially makes a selection among multiple arms (alternative paths or task types) so as to maximize the total reward (i.e., throughput) within a given time budget. In these applications, the cost (i.e., completion time) and reward of each arm pull can be potentially correlated and heavy-tailed (Harchol-Balter, 2000; Jelenković and Tan, 2013).

In this paper, we investigate the unique dynamics of this extended budget-constrained bandit setting with general cost and reward distributions. Unlike the classical stochastic MAB problem, each action incurs a random cost and yields a random reward in our model. Under a budget constraint $B$, the objective of the controller is to maximize the expected cumulative reward until the total cost exceeds the budget. As we will see, the correlation and variability of the cost-reward pairs can have a substantial impact on the performance in this bandit setting, which we incorporate in the design of learning algorithms for near-optimal performance. Many of our results are obtained for a very general setting where the cost and reward can be correlated and heavy-tailed, but sharper results are presented for some interesting special cases.

### 1.1 Main Contributions

The main objective in this paper is to design efficient algorithms that achieve provably tight regret bounds in an extended setting of correlated and potentially heavy-tailed cost and reward. Our main contributions are as follows:

1. **Exploiting the correlation:** One of the key contributions in this work is to use a linear minimum mean square (LMMSE) estimator to extract and exploit the correlation between the cost and reward of an arm (see Section 4.2). Furthermore, we incorporate the effect of variability in cost-reward pairs through variance. Consequently, we achieve provably tight problem-dependent regret bounds in an extended setting of unbounded cost and reward.

2. **Extension to unbounded cost and reward:** We develop novel design and analysis methods for the setting of unbounded and potentially heavy-tailed cost and reward pairs, and show that $O\big(\log(B)\big)$ regret is achievable if moments of order $2 + \gamma$ exist for some $\gamma > 0$ for all cost and reward pairs (see Section 4.3).

3. **Regret lower bounds:** We establish a regret lower bound for the budget-constrained bandit problem (see Section 5). By using this result, we obtain explicit regret lower bounds for jointly Gaussian cost-reward distributions. Consequently, we prove that the algorithms we propose in this paper achieve tight regret bounds, which are optimal up to a constant factor in the case of jointly Gaussian cost and reward.

### 1.2 Related Work

The classical stochastic multi-armed bandit problem, which is a specific case of the model we study in this paper, has been extensively studied in the literature. For detailed discussion on the basic model, we refer to (Bubeck et al., 2012; Berry and Fristedt, 1985).

The budget-constrained MAB problem and its variants were investigated in a variety of papers. In (Tran-Thanh et al., 2012) and (Combes et al., 2015), budget-constrained multi-armed bandit problem is investigated where each arm pull incurs an arm-dependent and deterministic cost. In (Guha and Munagala, 2009), the budgeted-bandit problem with deterministic costs is investigated from a Bayesian perspective, and constant-factor approximation algorithms are proposed. In (György et al., 2007), the continuous-time extension of the MAB problem with side information is investigated, which is an early example for the budget-constrained bandit problem. In (Badanidiyuru et al., 2013; Agrawal and Devanur, 2014), the bandit problem under multiple budget constraints is examined, and problem-independent regret bounds of order $\tilde{O}(\sqrt{B})$ are obtained. Bandits with knapsacks have been extended to other bandit settings (Agrawal and Devanur, 2016; Badanidiyuru et al., 2014; Sankararaman and Slivkins, 2017; Ding et al., 2013). In (Xia et al.,

2015, 2016), the budget-constrained MAB problem is explored in a similar setting to ours. In these works, the cost and reward of each arm are supported in $[0, 1]$, and the correlation between them is not exploited. In (Cayci et al., 2019), the authors consider a variation of the budget-constrained bandit problem where the controller has the option to interrupt an ongoing cycle for a faster alternative. The interruption mechanism brings significantly different dynamics to the problem that is investigated in this paper.

Bandits with heavy-tailed reward distributions are considered in (Liu and Zhao, 2011; Bubeck et al., 2013). These papers are still in the scope of the classical MAB setting: the budget is consumed deterministically at rate 1 by each action, so the dynamics of the random resource consumption with heterogeneous statistics are not included in the model.

## 2 System Setup

In this paper, we consider a bandit problem with $K$ arms. The set of arms is denoted by $\mathbb{K} = \{1, 2, \ldots, K\}$. Each arm $k \in \mathbb{K}$ is described by a two-dimensional random process $\{(X_{n,k}, R_{n,k}) : n \geq 1\}$ that is independent from other arms. If arm $k$ is chosen at $n$-th epoch, it incurs a cost of $X_{n,k}$ and yields a reward of $R_{n,k}$, where both are learned via a bandit feedback only after the decision is made. The controller has a cost budget $B > 0$, and tries to maximize the expected cumulative reward it receives by sampling the arms wisely under this budget constraint.

The pair $(X_{n,k}, R_{n,k})$ is assumed to be independent and identically distributed over $n$, but the cost $X_{n,k}$ and reward $R_{n,k}$ can be positively correlated. We allow $X_{n,k}$ to take on negative values, but the drift is assumed to be positive, i.e., there exists $\mu_* > 0$ such that $\mathbb{E}[X_{n,k}] \geq \mu_* > 0$ for all $k$.

Let $\pi$ be an algorithm that yields a sequence of arm pulls $\{I_n^\pi \in \mathbb{K} : n \geq 1\}$. Under $\pi$, the history until epoch $n$ is the following filtration:

$$\mathcal{F}_n^\pi = \sigma(\{(X_{j,k}, R_{j,k}) : I_j^\pi = k, 1 \leq j \leq n\}), \quad (1)$$

where $\sigma(X)$ denotes the sigma-field of a random variable $X$. We call an algorithm $\pi$ admissible if $\pi$ is non-anticipating, i.e., $\{I_n^\pi = k\} \in \mathcal{F}_{n-1}^\pi$ for all $k, n$. The set of all admissible policies is denoted as $\Pi$.

The total cost incurred in $n$ epochs under an admissible policy $\pi \in \Pi$ is a controlled random walk which is defined as $S_n^\pi = \sum_{i=1}^n X_{i, I_i^\pi}$. The arm pulling process under an algorithm $\pi$ continues until the budget $B$ is depleted. We assume that the reward corresponding to the final epoch during which the budget is depleted is gathered by the controller. Thus, the total number

of pulls under $\pi$ is defined as follows:

$$N_\pi(B) = \inf\left\{n : S_n^\pi > B\right\}. \qquad (2)$$

Note that the total number of pulls $N_\pi(B)$ is a stopping time adapted to the filtration $\{(\mathcal{F}_t^\pi) : t \geq 0\}$. With these definitions, the cumulative reward under a policy $\pi$ can be written as follows:

$$\mathrm{REW}_\pi(B) = \sum_{i=1}^{N_\pi(B)} R_{i,I_i^\pi}. \qquad (3)$$

The objective in this paper is to design algorithms that achieve maximum $\mathbb{E}[\mathrm{REW}_\pi(B)]$, or equivalently minimum regret, which is defined as follows:

$$Reg_\pi(B) = \mathbb{E}[\mathrm{REW}_{\pi^{\mathrm{opt}}}(B)] - \mathbb{E}[\mathrm{REW}_\pi(B)], \qquad (4)$$

where $\pi^{\mathrm{opt}}(B)$ denotes the optimal policy:

$$\pi^{\mathrm{opt}}(B) \in \arg\max_{\pi' \in \Pi} \mathbb{E}[\mathrm{REW}_{\pi'}(B)],$$

for any $B > 0$.

In the following section, we investigate the optimal policy that maximizes the expected cumulative reward when all arm distributions are known, and provide low-complexity approximations that have desirable performance characteristics.

## 3 Approximations of the Oracle

The optimization problem described in Section 2 is a variant of the unbounded knapsack problem, and it is known that similar stochastic control problems are PSPACE-hard (Badanidiyuru et al., 2013; Papadimitriou and Tsitsiklis, 1999). In order to find a tractable benchmark, we will consider approximation algorithms with provably good performance in this section.

The main quantity of interest will be the reward rate, which is defined as follows:

$$r_k = \frac{\mathbb{E}[R_{1,k}]}{\mathbb{E}[X_{1,k}]}, \ \ k \in \mathbb{K}. \qquad (5)$$

Intuitively, if arm $k$ is chosen persistently until the budget $B > 0$ is depleted, the cumulative reward becomes $r_k B + o(B)$ as $B \to \infty$. The additive $o(B)$ term is $O(1)$ if $\mathbb{E}[(X_{1,k}^+)^2] < \infty$ by Lorden's inequality (Asmussen, 2008). Hence, pulling the arm with the highest reward rate is a logical choice.

In the following, we prove that the optimality gap is $O(1)$ under mild moment conditions, which covers the case of heavy-tailed cost-reward pairs.

**Definition 1** (Optimal Static Algorithm). *Let $k^*$ be the arm with the highest reward rate:*

$$k^* \in \arg\max_{k \in \mathbb{K}} r_k.$$

*The optimal static policy, denoted by $\pi^*$, pulls $k^*$ until the budget is depleted: $I_n^{\pi^*} = k^*$ for all $n \leq N_{\pi^*}(B)$.*

The main result of this section is the following proposition, which implies that $\pi^*$ is a plausible approximation algorithm for $\pi^{\mathrm{opt}}(B)$ for all $B > 0$ under mild moment conditions.

**Assumption 1.** *There exists $\gamma > 0$ such that $\mathbb{E}[(X_{1,k}^+)^{2+\gamma}] < \infty$ for all $k \in \mathbb{K}$.*

**Proposition 1** (Optimality Gap for $\pi^*$). *Under Assumption 1, there exists a constant*

$$G^\star = G^\star\left(\min_k \mathbb{E}[X_{1,k}], \max_k Var(X_{1,k})\right) < \infty,$$

*independent of $B$ such that the following holds:*

$$\max_{\pi \in \Pi} \ \mathbb{E}[\mathrm{REW}\pi(B)] - \mathbb{E}[\mathrm{REW}\pi^*(B)] \leq G^\star, \qquad (6)$$

*for any $B > 0$. Consequently, $\pi^*$ is asymptotically optimal as $B \to \infty$.*

*Proof.* The proof of Proposition 1 is based on tools from stochastic control, and is given in Appendix A. $\square$

Proposition 1 implies that the optimality gap of the optimal static policy is a constant with respect to the budget $B$, which depends on the first- and second-order moments of the cost. This extends the result presented in (Xia et al., 2016) for bounded and strictly positive costs to unbounded costs with positive drift that can take on negative values. Also, for small $B$ values, there can be dynamic policies that outperform this simple static policy (Dean et al., 2004). However, the optimality gap is still $O(1)$ for these dynamic policies, therefore we consider $\pi^*$ for its simplicity and efficiency.

Now that we have an accurate approximation for the oracle, we propose the first and basic algorithms that assume the knowledge of second-order moments.

## 4 Algorithms for Known Second-Order Moments

In this section, we will assume that the second-order moments of all cost-reward pairs are known by the decision maker. First, in Section 4.2, we will consider the case $(X_{n,k}, R_{n,k})$ are jointly Gaussian, and propose a learning algorithm that achieves tight regret bound on
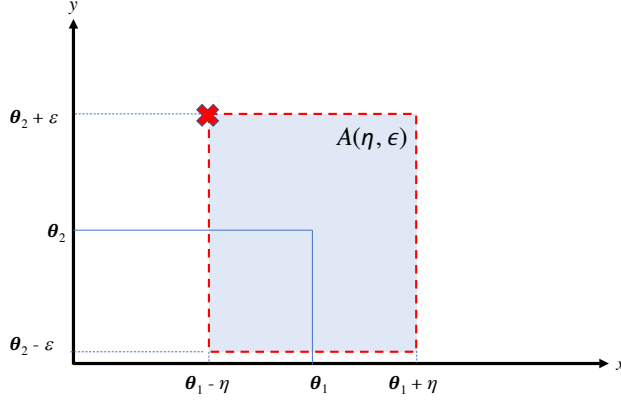
Figure 1: If $(\hat{\theta}_1, \hat{\theta}_2)$ is in the high-probability set $A(\eta, \epsilon)$, then the maximum deviation of $\hat{r} = \frac{\hat{\theta}_2}{\hat{\theta}_1}$ from $r$ is $\frac{\lambda(\epsilon + r\eta)}{\theta_1}$, and it is achieved at the marked corner.

the order of $O(\log(B))$ by using the correlation information. Then, in Section 4.3, we will study the general case where the cost and reward can be unbounded and potentially heavy-tailed, and propose algorithms that achieve the same regret bounds (up to a constant) as the sub-Gaussian case.

The following proposition provides a basis for the algorithm design and analysis throughout the paper.

### 4.1 Preliminaries: Rate Estimation

Let $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$ be a pair of unknown constants for which $r = \frac{\theta_2}{\theta_1}$ is to be estimated. The following proposition yields a useful device to obtain concentration results for r from concentration results for $\theta_1$ and $\theta_2$ for this estimation procedure.

**Proposition 2** (Rate Estimation). *Let $\hat{\theta}_1$ and $\hat{\theta}_2$ be estimators for $\theta_1 > 0, \theta_2 \geq 0$, respectively. If*

$$\eta \in \left(0, \frac{\theta_1(\lambda - 1)}{\lambda}\right), \tag{7}$$

*for some $\lambda > 1$, then we have the following result:*

$$\mathbb{P}\left(|r - \frac{\hat{\theta}_2}{\hat{\theta}_1}| > \frac{\lambda(\epsilon + r\eta)}{\theta_1}\right) \leq \mathbb{P}(|\hat{\theta}_1 - \theta_1| > \eta)$$
$$+ \mathbb{P}(|\hat{\theta}_2 - \theta_2| > \epsilon).$$

Therefore, if $\hat{\theta}_1$ and $\hat{\theta}_2$ both achieve exponential convergence rate, then $\frac{\hat{\theta}_2}{\hat{\theta}_1}$ converges to $r$ exponentially fast. The intuition behind the proposition is illustrated in Figure 4.1.

**Remark 1** (Stability of the rate estimator). The condition $\eta < \theta_1$, i.e., sufficient concentration of the estimator around the true parameter $\theta_1$, is crucial for

Proposition 2. Note that if the variability of the mean estimator is high and thus $A(\eta, \epsilon)$ intersects with the $y$-axis, then the above bound is useless as $\hat{r}$ can have arbitrarily large deviations from $r$.

In the following, we propose algorithms under the assumption that the second-order moments for each arm $k$ is known by the controller.

### 4.2 Sub-Gaussian Case: Algorithm UCB-B1

The main idea behind UCB-B1 is to use an upper confidence bound for the reward rate $r_k$. Let $T_k(n)$ be the number of pulls for arm $k$ in the first $n$ stages and $\hat{r}_{k,n} = \frac{\max\{0, \hat{\mathbb{E}}_n[R_k]\}}{\max\{b, \hat{\mathbb{E}}_n[X_k]\}}$ where

$$\hat{\mathbb{E}}_n[X_k] = \frac{1}{T_k(n)} \sum_{i=1}^n \mathbb{I}\{I_i = k\} X_{i,k},$$

$$\hat{\mathbb{E}}_n[R_k] = \frac{1}{T_k(n)} \sum_{i=1}^n \mathbb{I}\{I_i = k\} R_{i,k},$$

and $b \leq \mathbb{E}[X_{1,k}]/2$ for all $k$. Instead of estimating $\mathbb{E}[X_{1,k}]$ and $\mathbb{E}[R_{1,k}]$ separately from the samples of $(X_{n,k}, R_{n,k})$, the correlation between $X_{n,k}$ and $R_{n,k}$ can be exploited to tighten the upper confidence bound for $r_k$. This is achieved by estimating $R_{n,k}$ by a linear estimator $\omega X_{n,k}$ so as to minimize $Var(R_{n,k} - \omega X_{n,k})$. Let

$$V(X_{1,k}, R_{1,k}) = \min_{\omega \in \mathbb{R}} Var(R_{1,k} - \omega X_{1,k}). \tag{8}$$

If $Var(X_{n,k}) > 0$, we have:

$$\omega_k = \arg\min_{\omega \in \mathbb{R}} Var(R_{1,k} - \omega X_{1,k}),$$
$$= \frac{Cov(X_{1,k}, R_{1,k})}{Var(X_{1,k})}, \tag{9}$$

by the orthogonality principle (Poor, 2013), and the optimal value of the objective is given by:

$$V(X_{1,k}, R_{1,k}) = Var(R_{1,k}) - \omega_k^2 Var(X_{1,k}).$$

If $Var(X_{n,k}) = 0$, we have $V(X_{1,k}, R_{1,k}) = Var(R_{1,k})$. This implies that $\omega_k$ and $V$ can be computed from the second-order moments of $(X_{n,k}, R_{n,k})$, which are assumed to be given in this section. For simplicity, we assume $\omega_k \leq r_k$ for all $k$ throughout the paper.

For non-negative $(M_X, M_R, L)$ that will be specified later, let

$$\epsilon_{k,n}^{\text{B}} = \frac{2\alpha M_R \log(n)}{3 T_k(n)} + \sqrt{L\alpha \frac{V(X_{1,k}, R_{1,k}) \log(n)}{T_k(n)}},$$

$$\eta_{k,n}^{\text{B}} = \frac{2\alpha M_X \log(n)}{3 T_k(n)} + \sqrt{L\alpha \frac{Var(X_{1,k}) \log(n)}{T_k(n)}}.$$

Then, if $S_n^\pi < B$, i.e., there is a remaining budget, then the `UCB-B1` Algorithm pulls an arm at stage $n+1$ according to:

$$I_{n+1} \in \arg\max_k \left\{ \widehat{r}_{k,n} + \widehat{c}_{k,n}^{\texttt{B1}} \right\},$$

where

$$\widehat{c}_{k,n}^{\texttt{B}} = 1.4 \frac{\epsilon_{k,n}^{\texttt{B}} + (\widehat{r}_{k,n} - \omega_k)\eta_{k,n}^{\texttt{B}}}{(\widehat{\mathbb{E}}_n[X_k])^+}$$

if the stability condition (7) holds for $\eta = \eta_{k,n}^{\texttt{B}}$ and $\lambda = 1.28$, and $\widehat{c}_{k,n}^{\texttt{B}} = \infty$ otherwise.

The regret performance of `UCB-B1` is presented in the following theorem.

**Theorem 1** (Regret Upper Bound for `UCB-B1`). *Let* $\Delta_k = r^* - r_k,$

$$\sigma_k^2 = V(X_{1,k}, R_{1,k}) + (r^* - \omega_k)^2 Var(X_{1,k}), \quad (10)$$

*for all* $k \in \mathbb{K}$ *and recall that* $\mu_* = \min_k \mathbb{E}[X_{1,k}].$

1. ***Bounded Cost and Reward:*** *If* $|X_{1,k}| \leq M_X$, $|R_{1,k}| \leq M_R$ *a.s.,* $\alpha > 2$ *and* $L = 2$, *then the regret under* `UCB-B1` *is upper bounded as:*

$$Reg_{\pi^{\texttt{B1}}}(B) \leq \alpha \sum_{k:\Delta_k > 0} \log\left(\frac{2B}{\mu_*}\right) C_k^{\texttt{B1}} + O(1), \quad (11)$$

*where* $M_k = M_R + r_k M_X$ *and*

$$C_k^{\texttt{B1}} = \frac{42\sigma_k^2}{\Delta_k \mathbb{E}[X_{1,k}]} + 42M_k + 21M_X\Delta_k,$$

*for all* $k$.

2. ***Jointly Gaussian Cost and Reward:*** *Let* $(X_{n,k}, R_{n,k})$ *be jointly Gaussian with known second-order moments. Then,* `UCB-B1`, $\alpha > 2$, $M_X = M_R = 0$ *and* $L = \frac{1}{2}$ *yields the following regret bound:*

$$Reg_{\pi^{\texttt{B1}}}(B) \leq \alpha \sum_{k:\Delta_k > 0} \log\left(\frac{2B}{\mu_*}\right) \frac{11\sigma_k^2}{\Delta_k \mathbb{E}[X_{1,k}]} + O(1),$$
$$(12)$$

*where* $\sigma_k$ *is defined in* (10).

*Proof.* The detailed proof, which will provide basis for the analysis of other algorithms proposed in this work, can be found in Appendix C. Note that the total reward is a controlled and stopped random walk with potentially unbounded support. Thus, the regret analysis requires new methods from the theory of martingales and stopped random walks. As such, we follow a proof strategy based on establishing a high-probability upper bound for $N_\pi(B)$, which can be found in Appendix B. □

## 4.3 Heavy-Tailed Case: Algorithm `UCB-M1`

In this subsection, we design a general algorithm that achieves the regret in the sub-Gaussian case (up to a constant) under the mild moment condition that $\mathbb{E}[(X_{1,k}^+)^{2+\gamma}] < \infty$ for all $k$.

The empirical mean estimator played a central role in the design of the `UCB-B1` Algorithm for sub-Gaussian distributions, which is proved to achieve $O(\log(B))$ regret. However, if we consider heavy-tailed distributions, the empirical mean estimator fails to achieve exponential convergence rate due to the frequent outliers (Bubeck et al., 2013). The median-based estimators, introduced in (Nemirovsky and Yudin, 1983) provide an elegant method to boost the convergence speed in mean estimation. The idea of boosting the confidence of weak independent estimators by taking the median was extended to general point estimation problems (beyond the mean estimation) in (Minsker et al., 2015). In the following, we will use a variation of this method in the design of median-based rate estimators.

Consider arm $k \in \mathbb{K}$ at stage $n$. For

$$m = \lfloor 3.5\alpha \log(n) \rfloor + 1,$$

we partition the observed samples $\{(X_{i,k}, R_{i,k}) : I_i = k, \ 1 \leq i \leq n\}$ into index sets $G_1, G_2, \ldots, G_m$ of size $\lfloor T_k(n)/m \rfloor$ each. Then, for each $j \in \{1, 2, \ldots, m\}$, let $\tilde{r}_{k,G_j} = \frac{\max\{\widehat{\mathbb{E}}_{G_j}[R_k], 0\}}{\max\{\widehat{\mathbb{E}}_{G_j}[X_k], b\}}$ where $b \leq \mathbb{E}[X_{1,k}]/2$, and

$$\widehat{\mathbb{E}}_{G_j}[X_k] = \sum_{i \in G_j} \frac{X_{i,k}}{|G_j|}, \qquad \widehat{\mathbb{E}}_{G_j}[R_k] = \sum_{i \in G_j} \frac{R_{i,k}}{|G_j|}.$$

The median-based rate estimator for arm $k$ at stage $n$ is thus

$$\overline{r}_{k,n} = \underset{1 \leq j \leq m}{\text{median}} \ \tilde{r}_{k,G_j}.$$

The deviations in the cost and reward are as follows:

$$\epsilon_{k,n}^{\texttt{M}} = 11\sqrt{\alpha \frac{V(X_{1,k}, R_{1,k}) \log(n)}{T_k(n)}},$$

$$\eta_{k,n}^{\texttt{M}} = 11\sqrt{\alpha \frac{Var(X_{1,k}) \log(n)}{T_k(n)}}.$$

Therefore, the decision at stage $(n+1)$ under `UCB-M1` is as follows:

$$I_{n+1} \in \arg\max_k \left\{ \overline{r}_{k,n} + \widehat{c}_{k,n}^{\texttt{M}} \right\} \qquad (13)$$

where

$$\widehat{c}_{k,n}^{\texttt{M}} = \frac{2\sqrt{2}(\epsilon_{k,n}^{\texttt{M}} + (\overline{r}_{k,n} - \omega_k)\eta_{k,n}^{\texttt{M}})}{\left( \underset{1 \leq j \leq m}{\text{median}} \ \widehat{\mathbb{E}}_{G_j}[X_k] \right)^+},$$

if the condition (7) is satisfied for $\eta = \underset{1 \leq j \leq m}{\text{median}} \, \widehat{\mathbb{E}}_{G_j}[X_k]$ and $\lambda = 1.28$.

For `UCB-M1`, we have the following regret upper bound.

**Theorem 2** (Regret Upper Bound for `UCB-M1`). *If the following moment conditions hold:*

- $\mathbb{E}[(X_{1,k}^+)^{2+\gamma}] < \infty$, *for all* $k$,

- $Var(R_{1,k}) < \infty$, *for all* $k$,

*then the regret under* `UCB-M1` *satisfies the following upper bound:*

$$Reg_{\pi^{\text{M1}}}(B) \leq \alpha \sum_{k:\Delta_k > 0} \log\left(\frac{2B}{\mu_*}\right) \frac{C\sigma_k^2}{\Delta_k \mathbb{E}[X_{1,k}]} + O(1),$$
(14)

*where* $\sigma_k$ *is as defined in* (10) *and* $C > 0$ *is a constant.*

*Proof.* The proof uses tools from the theory of martingales and stopped random walks, and can be found in Appendix B and Appendix C. $\square$

**Remark 2.** We have the following observations from Theorem 1 and 2:

- If $Var(X_{1,k}) \downarrow 0$ and $\mathbb{E}[X_{1,k}] = 1$, the regret upper bounds match with the existing regret bounds for the stochastic bandit problem.

- Note that for positively correlated $X_{n,k}$ and $R_{n,k}$, one can ignore the correlation and use an upper confidence bound based on the separate estimation of $X_{n,k}$ and $R_{n,k}$. From Theorem 1, it can be observed that this scheme leads to a loss of $O\left(\sum_k Cov(X_{1,k}, R_{1,k})\right)$. Moreover, as it will be seen in the next section, this is nearly the best way of exploiting the correlation in the case of jointly Gaussian cost and reward pairs.

- The `UCB-M1` Algorithm achieves the same regret upper bound as the `UCB-B1` Algorithm up to a constant with much less moment assumptions: while `UCB-B1` requires sub-Gaussianity, `UCB-M1` requires only existence of moments of order $(2+\gamma)$ for some $\gamma > 0$ for the costs, and second-order moments for the rewards. However, the constant that multiplies the $O(\log B)$ term is much higher in `UCB-M1` than `UCB-B1`, which can be viewed as the cost of generality.

- If the cost is deterministic, i.e., $Var(X_{1,k}) = 0$, then the regret is monotonically decreasing in $\Delta_k$ as $O\left(\frac{\log B}{\Delta_k}\right)$ for each arm $k$. However, for random costs, since $r^* = r_k + \Delta_k$, the regret bounds have an additive term scaling linearly in $\Delta_k$ as

$O\left(\log\left(\frac{2B}{\mu_*}\right) \sum_k \frac{Var(X_{1,k})}{\mathbb{E}[X_{1,k}]}\Delta_k\right)$, which might seem strange at first since the separability of a suboptimal arm $k$ increases with its corresponding $\Delta_k$. This is a unique phenomenon observed in the case of stochastic costs: recall from Remark 1 that the rate estimator is unstable when the confidence interval for the estimation of $\mathbb{E}[X_{1,k}]$ is large, and thus it incurs $\mathbb{E}[X_{1,k}]\Delta_k$ regret per pull since rate estimation is unreliable. As it will be seen in Corollary 1, the same term appears with the same coefficient in the regret lower bound for jointly Gaussian cost-reward pairs, which implies that it is inevitable at least in that case.

## 5 Regret Lower Bound for Admissible Policies

In this section, we will propose regret lower bounds for the budget-constrained bandit problem based on (Lai and Robbins, 1985). In the specific case of jointly Gaussian cost-reward pairs, we can determine a lower bound explicitly, which provides useful insight about the impact of variability and correlation on the regret.

In order to establish a regret lower bound, assume that the joint distribution of $\{(X_{n,k}, R_{n,k}) : n \geq 1\}$ is parametrized by $\theta_k \in \Theta_k$ for some parameter space $\Theta_k$, i.e., $(X_{n,k}, R_{n,k}) \sim P_{\theta_k}$. For any $k \in \mathbb{K}$ and $\theta \in \Theta_k$, let $r_k(\theta) = \frac{\mathbb{E}_\theta[R_{1,k}]}{\mathbb{E}_\theta[X_{1,k}]}$ be the reward rate (i.e., reward per unit cost). Furthermore, for a given bandit instance $\vec{\theta} = (\theta_1, \theta_2, \ldots, \theta_K)$, let $r^* = \max_k r_k(\theta_k)$ be the optimal reward rate, and $\Delta_k = r^* - r_k(\theta_k)$. For admissible policies, we have the following regret lower bound, which is an extension of Lai-Robbins style regret lower bounds for the stochastic bandit problem (Lai and Robbins, 1985; Burnetas and Katehakis, 1996).

**Theorem 3** (Regret Lower Bound). *Suppose that* $\mathbb{E}[(X_{1,k})^{2+\gamma}] < \infty$ *for some* $\gamma > 0$ *and* $Var(R_{1,k}) < \infty$ *hold for all* $k$. *Assume that the following conditions are satisfied by* $P_{k,\theta}$ *for any* $k$:

1. *If* $r_k(\theta_1) > r_k(\theta_2)$, *then* $D(P_{k,\theta_2}||P_{k,\theta_1}) < \infty$,

2. *(Denseness)* $r_k(\Theta_k) = \{r_k(\theta) : \theta \in \Theta_k\}$ *is dense,*

3. *(Continuity)* $\theta \mapsto D(P_{k,\theta_k}||P_{k,\theta})$ *is a continuous mapping.*

*For a given bandit instance* $\vec{\theta} = (\theta_1, \theta_2, \ldots, \theta_K)$, *if* $\pi \in \Pi$ *is a policy such that* $\mathbb{E}[T_k^\pi(n)] = o(n^\alpha)$ *for any* $\alpha > 0$ *and* $k$ *such that* $r_k(\theta_k) < r^*$, *then we have the following lower bound:*

$$\liminf_{B \to \infty} \frac{Reg_\pi(B)}{\log(B)} \geq \frac{1}{2} \sum_{k:\Delta_k > 0} \frac{\mathbb{E}[X_{1,k}]\Delta_k}{D_k^\star},$$
(15)

where $D_k^\star$ is the solution to the following optimization problem:

$$D_k^\star = \min_{\theta \in \Theta_k} D(P_{k,\theta_k} || P_{k,\theta}) \text{ subject to } r_k(\theta) \geq r^*.$$

*Proof.* The proof can be found in Appendix E. □

The regret lower bound has an explicit form if the cost and reward distributions of each arm is jointly Gaussian with a known covariance matrix.

**Corollary 1** (Jointly Gaussian Cost and Reward). *Let* $(X_{n,k}, R_{n,k})$ *be jointly Gaussian:*

$$(X_{n,k}, R_{n,k}) \sim \mathcal{N}(\mu_k, \Sigma_k),$$

*for all* $k \in \mathbb{K}$ *where* $\mu_k = (\mathbb{E}[X_{n,k}], \mathbb{E}[R_{n,k}])$ *and*

$$\Sigma_k = \begin{pmatrix} Var(X_{n,k}) & Cov(X_{n,k}, R_{n,k}) \\ Cov(X_{n,k}, R_{n,k}) & Var(R_{n,k}) \end{pmatrix}.$$

*If* $\Sigma_k$ *is known and* $\mu_k$ *is unknown by the controller for all* $k \in \mathbb{K}$, *we have the following regret lower bound for the Gaussian case:*

$$\liminf_{B \to \infty} \frac{Reg_\pi(B)}{\log(B)} \geq \sum_{k:\Delta_k>0} \frac{\sigma_k^2}{\mathbb{E}[X_{1,k}]\Delta_k}, \qquad (16)$$

*where* $\sigma_k^2$ *is defined in* (10).

*Proof.* For known $\Sigma_k$, we have $D_k^\star = \frac{(\mathbb{E}[X_{1,k}]\Delta_k)^2}{2\sigma_k^2}$ for $\theta_k = \mu_k$ and $\Theta_k = \mathbb{R}_+^2$. Using this in Theorem 3 yields the result. □

**Remark 3** (Optimality of UCB-B1 and UCB-M1). *Comparing* (12) *and* (14) *with* (16), *we can deduce that* UCB-B1 *and* UCB-M1 *achieve optimal regret up to a universal constant for the case of jointly Gaussian cost and reward pairs with known covariance matrix.*

# 6 Algorithms for Unknown Second-Order Moments

In Section 4, we proposed algorithms under the assumption that the second-order moments are known for each arm $k$. However, in practice, these second-order moments are unknown, and therefore to be estimated from the samples collected via bandit feedback. In this section, we will propose algorithms that use these second-order moment estimates to achieve tight regret bounds.

The general strategy in the development of the algorithms in this section is to use empirical estimates for the second-order moments that appear in UCB-B1 as a surrogate.

## 6.1 Bounded and Uncorrelated Cost and Reward: UCB-B2

For clarity, we first consider the case $X_{n,k}$ and $R_{n,k}$ are uncorrelated for all $k$ and $X_{n,k} \in [0, M_X]$ and $R_{n,k} \in [0, M_R]$ almost surely for known $M_X, M_R > 0$. In this case, we will propose an algorithm based on a variant of the empirical Bernstein inequality, which was introduced in (Audibert et al., 2009).

For any $k$, let the variance estimate $\widehat{V}_{k,n}(X_k)$ be defined as follows:

$$\widehat{V}_{k,n}(X_k) = \frac{1}{T_k(n)} \sum_{i=1}^n \mathbb{I}\{I_i = k\}(X_{i,k} - \widehat{\mathbb{E}}_n[X_{1,k}])^2,$$

where $\widehat{\mathbb{E}}_n[X_k]$ is the empirical mean of the observations up to epoch $n$.

The bias terms in UCB-B2 are defined as follows:

$$\epsilon_{k,n}^{\texttt{B2}} = \sqrt{\frac{2\widehat{V}_{k,n}(R_k)\log(n^\alpha)}{T_k(n)}} + \frac{3M_R\log(n^\alpha)}{T_k(n)},$$

$$\eta_{k,n}^{\texttt{B2}} = \sqrt{\frac{2\widehat{V}_{k,n}(X_k)\log(n^\alpha)}{T_k(n)}} + \frac{3M_X\log(n^\alpha)}{T_k(n)}.$$

Let $\widehat{r}_{k,n}$ be the empirical reward rate estimator in Section 4.2, and

$$\widehat{c}_{k,n}^{\texttt{B2}} = 1.4 \frac{\epsilon_{k,n}^{\texttt{B2}} + \widehat{r}_{k,n}\eta_{k,n}^{\texttt{B2}}}{(\widehat{\mathbb{E}}_n[X_k])^+}, \qquad (17)$$

if the condition (7) is satisfied for $\lambda = 1.28$ ($\widehat{c}_{k,n}^{\texttt{B2}} = \infty$ otherwise). Then, at stage $n+1$, the following decision is made under UCB-B2:

$$I_{n+1} \in \arg\max_k \left\{\widehat{r}_{k,n} + \widehat{c}_{k,n}^{\texttt{B2}}\right\}.$$

The lack of knowledge for the second-order statistics loosen the upper confidence bound for the rate estimator, which in turn increases the regret. In the following, we provide the regret upper bounds for UCB-B2 to gain insight about the impact of using variance estimates on the performance of the algorithm.

**Theorem 4** (Regret Upper Bound for UCB-B2). *Let* $\sigma_k$ *and* $M_k$ *be as defined in Theorem 1. Then, we have the following upper bound for the regret under* UCB-B2:

$$Reg_{\pi^{\texttt{B2}}}(B) \leq \alpha \sum_{k:\Delta_k>0} \log\left(\frac{2B}{\mu_*}\right)(C_k^{\texttt{B1}} + \delta C_k) + O(1), \qquad (18)$$

*where*

$$\delta C_k = 21\left(\frac{M_X^4 \Delta_k \mu_k}{Var^2(X_{1,k})} + \frac{Var(X_{1,k})\Delta_k}{\mu_k}\right). \qquad (19)$$

*for* $\mu_k = \mathbb{E}[X_{1,k}]$.

The proof of Theorem 4 can be found in Appendix F.

**Remark 4** (Impact of Unknown Variances)**.** The additional terms are caused by the stability of the rate estimator: since we use a variance estimate in the upper confidence bound of $X_{n,k}$, the rate estimator suffers from a longer period of instability, which increases the regret coefficient proportional to $\Delta_k$.

### 6.2 Learning the Correlation: UCB-B2C

Finally we consider the case $(X_{n,k}, R_{n,k})$ are bounded and correlated, but the second-order moments are unknown. In the absence of correlation, our goal was to estimate $Var(R_{1,k})$ and $Var(X_{1,k})$ from the samples of $(X_{n,k}, R_{n,k})$. When there is a correlation, we have an optimization problem: we need to establish confidence bounds for the LMMSE estimator $\omega_k$ defined in (9) as well as the minimum variance $Var(R_{1,k} - \omega_k X_{1,k})$ by using the samples of $(X_{n,k}, R_{n,k})$ observed via bandit feedback. We take a loss minimization approach in the statistical learning setting to estimate these quantities.

For any $k \in \mathbb{K}$, let the empirical LMMSE estimator be defined as follows:

$$\widehat{\omega}_{k,n} = \arg\min_{\omega' \in \mathbb{R}} \ \widehat{L}_{k,n}(\omega)$$

where the empirical loss function is the following:

$$\widehat{L}_{k,n}(\omega) = \sum_{i=1}^{n} \frac{\mathbb{I}\{I_i = k\}}{T_k(n)} \Big( R_i - \widehat{\mathbb{E}}_n[R] - \omega\big(X_i - \widehat{\mathbb{E}}_n[X]\big)\Big)^2.$$

It can be shown that $\widehat{\omega}_{k,n} \to \omega_k$ if $T_k(n) \to \infty$ as $n \to \infty$, and moreover the convergence rate is exponential and tight concentration bounds for $\widehat{\omega}_{k,n}$ and $\widehat{L}_{k,n}(\widehat{\omega}_{k,n})$ can be established. Let $M_Z = M_R + \overline{\omega} M_X$ where $\overline{\omega} > \max_k \ \omega_k$ is a given parameter, and let

$$\nu_{k,n}(\omega_k) = \frac{1.4 M_X M_Z}{Var(X_{1,k})} \sqrt{\frac{\log n^\alpha}{T_k(n)}}, \qquad (20)$$

$$\nu_{k,n}(L_k) = M_Z^2 \sqrt{\frac{2 \log n^\alpha}{T_k(n)}}. \qquad (21)$$

Then, it can be shown that $-\widehat{\omega}_{k,n} + \nu_{k,n}(\omega_k)$ and $\widehat{L}_{k,n}(\widehat{\omega}_{k,n}) + \nu_{k,n}(\omega_k)$ are high-probability upper bounds for $-\omega_k$ and $\min_\omega \ Var(R_{1,k} - \omega X_{1,k})$, respectively, for large enough $T_k(n)$.

The bias terms in UCB-B2C are defined as follows:

$$\epsilon_{k,n}^{\text{B2C}} = \sqrt{\frac{2\widehat{L}_{k,n}(\widehat{\omega}_{k,n}) \log(n^\alpha)}{T_k(n)}} + \frac{3 M_Z \log(n^\alpha)}{T_k(n)},$$

$$\eta_{k,n}^{\text{B2C}} = \sqrt{\frac{2\widehat{V}_{k,n}(X_k) \log(n^\alpha)}{T_k(n)}} + \frac{3 M_X \log(n^\alpha)}{T_k(n)}.$$

Then, at stage $n + 1$, the following decision is made under UCB-B2C:

$$I_{n+1} \in \arg\max_k \ \left\{ \widehat{r}_{k,n} + \widehat{c}_{k,n}^{\text{B2C}} \right\},$$

where

$$\widehat{c}_{k,n}^{\text{B2C}} = 1.4 \frac{\epsilon_{k,n}^{\text{B2C}} + (\widehat{r}_{k,n} - \widehat{\omega}_{k,n})^+ \eta_{k,n}^{\text{B2C}}}{\big(\widehat{\mathbb{E}}_n[X_k]\big)^+},$$

if the stability condition (7) is satisfied for $\lambda = 1.28$, and $\widehat{c}_{k,n}^{\text{B2C}} = \infty$ otherwise.

In the following, we investigate the impact of using second-order moment estimates on the regret of UCB-B2C. The proof can be found in Appendix G.

**Theorem 5** (Regret Upper Bound for UCB-B2C)**.** *Let $C_k^{\text{B1}}$ be defined as in Theorem 1. Then, we have the following upper bound for the regret under UCB-B2:*

$$Reg_{\pi^{\text{B2C}}}(B) \le \alpha \sum_{k:\Delta_k>0} \log\left(\frac{2B}{\mu_*}\right)(C_k^{\text{B1}} + \delta C_k') + O(1),$$

*where*

$$\delta C_k' = \delta C_k + 42\Big( \frac{M_Z M_X}{\sqrt{Var(X_{1,k})}} + \frac{M_X^4 \Delta_k \mu_k}{Var^2(X_{1,k})} \Big). \quad (22)$$

*for $\mu_k = \mathbb{E}[X_{1,k}]$ and $\delta C_k$ defined in (19).*

Note that the regret of UCB-B2C converges to the regret of UCB-B2, and they both approach to the performance of the UCB-B1 Algorithm as $\Delta_k \downarrow 0$.

## 7 Conclusions

In this paper, we considered a very general setting for the budgeted bandit problem where each action incurs a potentially correlated and heavy-tailed cost-reward pair. We proved that positive expected cost and existence of moments of order $2 + \gamma$ for some $\gamma > 0$ suffice for $O(\log B)$ regret for a given budget $B > 0$. For known second-order moments, we proposed two algorithms named UCB-B1 and UCB-M1 that exploit the correlation between cost and reward by using an LMMSE estimator. By proposing a regret lower bound, we proved that UCB-B1 and UCB-M1 achieve order optimality, and moreover they achieve optimal regret up to a universal constant for the specific case of jointly Gaussian cost and reward pairs, which underlines the significance of second-order moments and correlation in the regret performance. For the case of bounded cost and reward with unknown second-order moments, we proposed learning algorithms UCB-B2 and UCB-B2C that estimate variances as well as LMMSE estimator to approach the performance of UCB-B1. We investigated the effect of using these estimates as surrogates in the absence of second-order moments, and showed that they approach the performance of UCB-B1 in certain cases.

## Acknowledgements

## References

S. Agrawal and N. Devanur. Linear contextual bandits with knapsacks. In *Advances in Neural Information Processing Systems*, pages 3450–3458, 2016.

S. Agrawal and N. R. Devanur. Bandits with concave rewards and convex knapsacks. In *Proceedings of the fifteenth ACM conference on Economics and computation*, pages 989–1006. ACM, 2014.

S. Asmussen. *Applied probability and queues*, volume 51. Springer Science & Business Media, 2008.

J.-Y. Audibert, R. Munos, and C. Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.

A. Badanidiyuru, R. Kleinberg, and A. Slivkins. Bandits with knapsacks. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 207–216. IEEE, 2013.

A. Badanidiyuru, J. Langford, and A. Slivkins. Resourceful contextual bandits. In *Conference on Learning Theory*, pages 1109–1134, 2014.

D. A. Berry and B. Fristedt. Bandit problems: sequential allocation of experiments (monographs on statistics and applied probability). *London: Chapman and Hall*, 5:71–87, 1985.

S. Bubeck, N. Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.

S. Bubeck, N. Cesa-Bianchi, and G. Lugosi. Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717, 2013.

A. N. Burnetas and M. N. Katehakis. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.

S. Cayci, A. Eryilmaz, and R. Srikant. Learning to control renewal processes with bandit feedback. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(2):43, 2019.

R. Combes, C. Jiang, and R. Srikant. Bandits with budgets: Regret lower bounds and optimal algorithms. *ACM SIGMETRICS Performance Evaluation Review*, 43(1):245–257, 2015.

B. C. Dean, M. X. Goemans, and J. Vondrdk. Approximating the stochastic knapsack problem: The benefit of adaptivity. In *45th Annual IEEE Symposium on Foundations of Computer Science*, pages 208–217. IEEE, 2004.

W. Ding, T. Qin, X.-D. Zhang, and T.-Y. Liu. Multi-armed bandit with budget constraint and variable costs. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.

S. Guha and K. Munagala. Multi-armed bandits with metric switching costs. In *International Colloquium on Automata, Languages, and Programming*, pages 496–507. Springer, 2009.

A. György, L. Kocsis, I. Szabó, and C. Szepesvári. Continuous time associative bandit problems. In *IJCAI*, pages 830–835, 2007.

M. Harchol-Balter. Task assignment with unknown duration. In *Proceedings 20th IEEE International Conference on Distributed Computing Systems*, pages 214–224. IEEE, 2000.

P. R. Jelenković and J. Tan. Characterizing heavy-tailed distributions induced by retransmissions. *Advances in Applied Probability*, 45(1):106–138, 2013.

T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.

K. Liu and Q. Zhao. Multi-armed bandit problems with heavy-tailed reward distributions. In *2011 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 485–492. IEEE, 2011.

S. Minsker et al. Geometric median and robust estimation in banach spaces. *Bernoulli*, 21(4):2308–2335, 2015.

A. S. Nemirovsky and D. B. Yudin. Problem complexity and method efficiency in optimization. 1983.

C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of optimal queuing network control. *Mathematics of Operations Research*, 24(2):293–305, 1999.

H. V. Poor. *An introduction to signal detection and estimation*. Springer Science & Business Media, 2013.

H. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58(5):527–535, 1952.

K. A. Sankararaman and A. Slivkins. Combinatorial semi-bandits with knapsacks. *arXiv preprint arXiv:1705.08110*, 2017.

L. Tran-Thanh, A. Chapman, A. Rogers, and N. R. Jennings. Knapsack based optimal policies for

budget–limited multi–armed bandits. In *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012.

Y. Xia, H. Li, T. Qin, N. Yu, and T.-Y. Liu. Thompson sampling for budgeted multi-armed bandits. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.

Y. Xia, W. Ding, X.-D. Zhang, N. Yu, and T. Qin. Budgeted bandit problems with continuous random costs. In *Asian conference on machine learning*, pages 317–332, 2016.