



HAL
open science

Contextual Bandits with Knapsacks for a Conversion Model

Zhen Li, Gilles Stoltz

► **To cite this version:**

Zhen Li, Gilles Stoltz. Contextual Bandits with Knapsacks for a Conversion Model. Thirty-sixth Conference on Neural Information Processing Systems, 2022, New Orleans, United States. hal-03683289v2

HAL Id: hal-03683289

<https://hal.science/hal-03683289v2>

Submitted on 29 Sep 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Contextual Bandits with Knapsacks for a Conversion Model

Zhen Li

BNP Paribas, 16 boulevard des Italiens, 75009 Paris, France
zhen.li@bnpparibas.com

Gilles Stoltz

Université Paris-Saclay, CNRS, Laboratoire de mathématiques d'Orsay, 91405, Orsay, France
gilles.stoltz@universite-paris-saclay.fr
HEC Paris, 1 rue de la Libération, 78350 Jouy-en-Josas, France
stoltz@hec.fr

Abstract

We consider contextual bandits with knapsacks, with an underlying structure between rewards generated and cost vectors suffered. We do so motivated by sales with commercial discounts. At each round, given the stochastic i.i.d. context \mathbf{x}_t and the arm picked a_t (corresponding, e.g., to a discount level), a customer conversion may be obtained, in which case a reward $r(a, \mathbf{x}_t)$ is gained and vector costs $\mathbf{c}(a_t, \mathbf{x}_t)$ are suffered (corresponding, e.g., to losses of earnings). Otherwise, in the absence of a conversion, the reward and costs are null. The reward and costs achieved are thus coupled through the binary variable measuring conversion or the absence thereof. This underlying structure between rewards and costs is different from the linear structures considered by Agrawal and Devanur [2016] (but we show that the techniques introduced in the present article may also be applied to the case of these linear structures). The adaptive policies exhibited solve at each round a linear program based on upper-confidence estimates of the probabilities of conversion given a and \mathbf{x} . This kind of policy is most natural and achieves a regret bound of the typical order $(\text{OPT}/B)\sqrt{T}$, where B is the total budget allowed, OPT is the optimal expected reward achievable by a static policy, and T is the number of rounds.

1 Introduction and Literature Review

We consider the framework of stochastic multi-armed bandits, which has been extensively studied since the early works by Thompson [1933] and Robbins [1952]. Two recent (and complementary) surveys summarizing the latest research in the field were written by Lattimore and Szepesvári [2020] and Slivkins [2019]. On the one hand, we are particularly interested in the setting of *contextual* stochastic multi-armed bandits, preferably with some structural assumptions on the dependency between rewards and contexts: linear models (again, a rich literature, see, among many others, Chu et al. [2011] and Abbasi-Yadkori et al. [2011], whose work marked a turning point), and, for $[0, 1]$ -valued rewards, logistic models (Filippi et al. [2010] and Fauray et al. [2020]). On the other hand, we are also particularly interested in stochastic multi-armed bandits *with knapsacks*, i.e., with cumulative vector-cost constraints to be abided by on top of maximizing the accumulated rewards. The setting was introduced by Badanidiyuru et al. [2013, 2018] and a comprehensive summary of the results achieved since then may be found in Slivkins [2019, Chapter 10]. The intersection of these two frameworks of interest is called *contextual bandits with knapsacks* [CBwK] and is the focus of the present article.

Literature review on CBwK. The first approach to CBwK, by Badanidiyuru et al. [2014] and Agrawal et al. [2016], assumes a joint stochastic generation of triplets of contexts-rewards-costs, with no specific underlying structure, and makes the problem tractable by using as a benchmark a finite set of static policies. As noted by Agrawal and Devanur [2016], picking this finite set may be uneasy, which is why they introduce instead a structural assumption of linear modeling: the (unknown) expected rewards and cost vectors depend linearly on the contexts.

We consider a different modeling assumption, motivated by sales with commercial discounts (see Appendix A): general (known) reward and cost functions are considered but they are coupled via a 0/1-valued factor, called a (customer) conversion, obtained as the realization of a Bernoulli variable with parameter $P(a, \mathbf{x})$ depending on the context \mathbf{x} observed (customer’s characteristics) and the action a taken (discount level offered). The probabilities $P(a, \mathbf{x})$ are themselves modeled by a logistic regression, whose parameters may be learned through an adaptation of the techniques by Filippi et al. [2010] and Fauray et al. [2020]. We do so in the first phase of the adaptive policy introduced in this article. More details on the comparison of the new setting considered to known settings of CBwK may be found in Section 2.2.

Primal-dual approach. The second phase of the adaptive policy exhibited uses the primal-dual approach to a convex optimization problem—actually, a simple optimization problem given by a linear program. This approach was already used in various ways for bandits with knapsacks, including CBwK, to define policies based on the dual problem: this is explicit in the LagrangeBwK policy of Immorlica et al. [2019] and is implicit in the reward-minus-weighted-cost approach of Agrawal and Devanur [2016] and Agrawal et al. [2016], as we underline in the proof sketch of Section 4 as well as in the discussion of Section 6. However, we only use the primal-dual approach in the analysis and state our adaptive policy directly in terms of the primal problem, where we substituted upper-confidence estimates of the probabilities $P(a, \mathbf{x})$. We therefore end up with a most natural adaptive policy, which mimics the optimal static policy used as a benchmark. This direct primal statement of the policy actually also works for the setting of linear CBwK studied by Agrawal and Devanur [2016], as we show in Section 6. Policies based on such direct primal statements were already considered for bandits with knapsacks (see Li et al. [2021] and references therein) but do not seem easily extendable to CBwK.

Outline and main contributions. The first contribution of this article is a new structured setting of CBwK, based on a coupling between general rewards and cost vectors through conversions modeled based on a logistic regression; we present and discuss it in Section 2.1 (and explain its origins in Appendix A of the supplementary material). The adaptive policy introduced is described in Section 3. Its first phase consists of learning the parameter of logistic regression and is adapted from Fauray et al. [2020]. Its second phase—and this is the second contribution of this article—directly solves a primal problem with optimistic conversion probabilities. The analysis, which we believe is concise, elegant, and natural, is provided in Sections 4 (when the context distribution ν is known) and 5 (when ν is unknown). As mentioned above, Section 6 draws the consequences of our second contribution for linear CBwK.

Notation. Throughout the article, vectors are denoted with bold symbols. In particular, $\mathbf{0}$ and $\mathbf{1}$ denote the vectors with all components equal to 0 and 1, respectively. With no additional subscript, $\|\mathbf{v}\|$ denotes the Euclidean norm of a vector \mathbf{v} , while a subscript given by a non-negative symmetric matrix M refers to $\|\mathbf{v}\|_M = \sqrt{\mathbf{v}^T M \mathbf{v}}$.

2 Learning Protocol and Motivation

We describe the learning protocol and objectives considered (Section 2.1) and explain why it is not covered by earlier works (Section 2.2). We also detail (Appendix A in the supplementary material) how this learning protocol was defined based on an industrial motivation in the banking sector: market share expansion for loans by granting discounts, under commercial budget constraints.

2.1 Learning Protocol and Modeling Assumptions

We consider a finite action set \mathcal{A} , including a special action a_{null} called no-op, and a finite context set $\mathcal{X} \subseteq \mathbb{R}^n$. (We discuss and mitigate finiteness of \mathcal{X} in Section 2.2.) A scalar reward function $r : \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$ and a vector-valued cost function $\mathbf{c} : \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]^d$ evaluate the performance

of actions given the contexts. There are several sources of costs to control: each corresponds to a component of \mathbf{c} . We assume that these functions are known, and (with no loss of generality) that their ranges are $[0, 1]$ and $[0, 1]^d$. The no-op action induces null reward and costs: $r(a_{\text{null}}, \mathbf{x}) = 0$ and $\mathbf{c}(a_{\text{null}}, \mathbf{x}) = \mathbf{0}$ for all $\mathbf{x} \in \mathcal{X}$.

Contexts—which correspond, for instance, to customers’ characteristics, see Appendix A—are drawn sequentially according to some distribution ν , which may be known or unknown (we will deal with both cases). At each round $t \geq 1$, upon observing the context $\mathbf{x}_t \in \mathcal{X}$ drawn, the learner picks an action $a_t \in \mathcal{A}$, which corresponds, for instance, to an offer made to the customer t . If the latter accepts the offer, an event which we denote $y_t = 1$, then the learner obtains a reward $r(a_t, \mathbf{x}_t)$ and suffers some costs $\mathbf{c}(a_t, \mathbf{x}_t)$. When the customer declines the offer, we set $y_t = 0$, and null reward and costs are obtained. Thus, in both cases, the reward and costs may be written as $r(a_t, \mathbf{x}_t) y_t$ and $\mathbf{c}(a_t, \mathbf{x}_t) y_t$. We call y_t the conversion and now explain how it is modeled.

Modeling conversions. We model each conversion y_t as an independent Bernoulli random drawn, with parameter $P(a_t, \mathbf{x}_t)$ depending on the context \mathbf{x}_t and action $a_t \neq a_{\text{null}}$. We further assume that these probabilities may be written as a logistic regression model, i.e., there exists a known transfer function $\varphi : \mathcal{A} \setminus \{a_{\text{null}}\} \times \mathcal{X} \rightarrow \mathbb{R}^m$ and some unknown parameter $\boldsymbol{\theta}_* \in \mathbb{R}^m$ such that

$$\forall \mathbf{x} \in \mathcal{X}, \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \quad P(a, \mathbf{x}) = \eta(\boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta}_*), \quad \text{where } \eta(x) = 1/(1+e^{-x}). \quad (1)$$

We assume that $\boldsymbol{\varphi}$ is normalized in a way that its Euclidean norm satisfies $\|\boldsymbol{\varphi}\| \leq 1$ and that a bounded convex set Θ containing $\boldsymbol{\theta}_*$ is known. Such a modeling is natural and opens the toolbox of logistic bandits; see Faury et al. [2020] and references cited therein. We however note (and discuss this fact in Appendix C) that the logistic regression model above is slightly different from the one by Faury et al. [2020].

The concept of a conversion y for a round when the no-op action a_{null} is played is void, and thus, we leave the probabilities $P(a_{\text{null}}, \mathbf{x})$ undefined, though by an abuse of notation, these quantities might appear but always multiplied by a 0, given, e.g., by indicator functions like $\mathbb{1}_{\{a \neq a_{\text{null}}\}}$, null rewards $r(a_{\text{null}}, \mathbf{x})$, or null costs $\mathbf{c}(a_{\text{null}}, \mathbf{x})$.

Policies: static vs. adaptive. The learner is given a number of rounds T and a maximal budget B (the same for all cost components, with no loss of generality: up to some normalization). A static policy is a function $\pi : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$, where $\mathcal{P}(\mathcal{A})$ is the set of probability distributions over \mathcal{A} . As is traditional in the literature of CBwK (we recall below why this is the case), we take as benchmark the static policy π^* with largest expected cumulative rewards under the condition that its cumulative costs abide by the budget constraints in expectation. More formally, π^* achieves the maximum defining

$$\begin{aligned} \text{OPT}(\nu, P, B) = & \max_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ \text{under } & T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \leq B \mathbf{1}, \end{aligned} \quad (2)$$

where $\mathbb{E}_{\mathbf{X} \sim \nu}$ denotes an expectation solely over random contexts \mathbf{X} following distribution ν , where $\pi_a(\mathbf{X})$ denotes the probability mass put by $\pi(\mathbf{X})$ on $a \in \mathcal{A}$, and where \leq is understood component-wise. Of course, the sums in the two expectations above are taken indifferently over \mathcal{A} or $\mathcal{A} \setminus \{a_{\text{null}}\}$.

The learner uses an adaptive policy, i.e., a sequence of measurable functions $\mathbf{p}_t : \mathcal{H}^{t-1} \times \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$ indexed by $t \geq 1$, where $\mathcal{H} = \mathcal{X} \times \mathcal{A} \times \{0, 1\}$. Indeed, the history available to the learner at the beginning of the round $t \geq 2$ is summarized by $h_{t-1} = (\mathbf{x}_s, a_s, y_s)_{s \leq t-1}$, and we define h_0 as the empty vector. Such a policy draws the action a_t for round $t \geq 1$ independently at random according to $\mathbf{p}_t(h_{t-1}, \mathbf{x}_t)$. We impose hard budget constraints on adaptive policies: they must satisfy

$$\sum_{t \leq T} \mathbf{c}(a_t, \mathbf{x}_t) y_t \leq B \mathbf{1} \quad \text{a.s.}$$

Such adaptive policies are called feasible in the literature. To abide by these hard constraints, we may restrict our attention to adaptive policies that pick Dirac masses on a_{null} whenever one component of the cumulative costs is larger than $B - 1$. At the same time, an adaptive policy should maximize the cumulative rewards obtained or, equivalently, minimize its regret:

$$R_T = \text{OPT}(\nu, P, B) - \sum_{t \leq T} r(a_t, \mathbf{x}_t) y_t.$$

BOX A: CONTEXTUAL BANDITS WITH KNAPSACKS [CBwK] FOR A CONVERSION MODEL

Known parameters: finite action set \mathcal{A} including a no-op action a_{null} ; finite context set $\mathcal{X} \subseteq \mathbb{R}^n$; scalar reward function $r : \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$; vector-valued cost function $\mathbf{c} : \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]^d$; number T of rounds; total budget constraint $B > 0$.

Possibly unknown parameters: context distribution ν on \mathcal{X} ; probability of conversion given action and context $P : \mathcal{A} \setminus \{a_{\text{null}}\} \times \mathcal{X} \rightarrow [0, 1]$, modeled as $P(a, \mathbf{x}) = \eta(\boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta}_*)$ for some known transfer function $\boldsymbol{\varphi} : \mathcal{A} \setminus \{a_{\text{null}}\} \times \mathcal{X} \rightarrow \mathbb{R}^m$, with $\|\boldsymbol{\varphi}\| \leq 1$, and some unknown parameter $\boldsymbol{\theta}_* \in \mathbb{R}^m$, lying in a known bounded convex set Θ .

For rounds $t = 1, 2, 3, \dots, T$:

1. Context $\mathbf{x}_t \sim \nu$ is drawn independently of the past;
2. Learner observes \mathbf{x}_t and picks an action $a_t \in \mathcal{A}$;
3. Conversion $y_t \in \{0, 1\}$ is drawn according to $\text{Ber}(P(a_t, \mathbf{x}_t))$;
4. Learner observes y_t , gets reward $r(a_t, \mathbf{x}_t) y_t$, and suffers costs $\mathbf{c}(a_t, \mathbf{x}_t) y_t$.

Goals: Maximize $\sum_{t \leq T} r(a_t, \mathbf{x}_t) y_t$ while controlling $\sum_{t \leq T} \mathbf{c}(a_t, \mathbf{x}_t) y_t \leq B \mathbf{1}$

It may be proved (along the same lines as Agrawal and Devanur [2016, Appendix B] do for a different model) that the optimal static policy π^* obtains, on average and in expectation, a cumulative reward at least as good as the best feasible adaptive policy.

Summary. A summary of the learning protocol and of the goals is provided in Box A. We note here that rewards gained and vector costs suffered at round t in the case $y_t = 1$ of a conversion could be stochastic with expectations $r(a_t, \mathbf{x}_t)$ and $\mathbf{c}(a_t, \mathbf{x}_t)$: our analysis and the regret bounds would be unchanged, as long as the expectation functions r and \mathbf{c} are known.

2.2 Discussion and Comparison to Existing Learning Protocols

The setting described above may be reduced to the general setting of CBwK, as introduced by Badanidiyuru et al. [2014] and Agrawal et al. [2016]. Indeed, introduce independent Bernoulli variables $y_{t,a}$ with parameters $P(a, \mathbf{x}_t)$, for all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$, and set $y_{t,a_{\text{null}}} = 0$. The vectors

$$\left(\mathbf{x}_t, (r_t(a))_{a \in \mathcal{A}}, (\mathbf{c}_t(a))_{a \in \mathcal{A}} \right), \quad \text{where } r_t(a) = r(a, \mathbf{x}_t) y_{t,a} \quad \text{and} \quad \mathbf{c}_t(a) = \mathbf{c}(a, \mathbf{x}_t) y_{t,a}$$

are i.i.d., and upon picking action $a_t \in \mathcal{A}$, the obtained and observed rewards and cost vectors equal $r_t(a_t)$ and $\mathbf{c}_t(a_t)$. When \mathcal{X} is discrete, we may consider the set Π of base policies that map \mathcal{X} to $\{\delta_a : a \in \mathcal{A}\}$, the set of Dirac masses at some $a \in \mathcal{A}$. The convex hull of Π is the set of all static policies $\mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$, against which we would like our policy to compete; but the adaptive policies by Badanidiyuru et al. [2014] and Agrawal et al. [2016] only compete with respect to the best single element in Π , not the best convex combination of elements of Π .

The setting of linear CBwK (Agrawal and Devanur [2016]) provides a structural link between contexts and expected rewards and cost vectors, but in a linear way that is incomparable to the setting of CBwK for a conversion model introduced above. More details are given in Section 6. We also mention that linear and logistic structural links between contexts (prices) and rewards or costs were studied in a non-contextual setting (i.e., not in CBwK) by Miao et al. [2021]. Their strategy bears some resemblance to the one by Agrawal and Devanur [2016], in particular, both consider an online convex optimization strategy as a subroutine.

All mentioned references consider a no-op action a_{null} . (It could be replaced by the existence of a standard action $a_{\text{no-cost}}$ always achieving null costs and possibly some positive rewards.)

On the contrary, none of the mentioned references assumes that the context \mathcal{X} set is finite. This is a technical necessity for a part of the adaptive policy introduced; see the discussion of computational complexity at the end of Section 3. But somehow, considering a finite set Π of policies, as in Badanidiyuru et al. [2014] and Agrawal et al. [2016], is a counterpart to assuming finiteness of \mathcal{X} . Also, Appendix F actually mitigates this restriction that \mathcal{X} is finite: learning the logistic parameter $\boldsymbol{\theta}_*$ may be achieved with continuous contexts (see Phase 1 in Section 3); only the subsequent

optimization part (Phase 2 in Section 3) requires finiteness of \mathcal{X} . We may thus well discretize only \mathcal{X} for this Phase 2, which is exactly what Appendix F performs. This mitigation comes with possible theoretical guarantees as Sections 4 and 5 reveal that the errors $\varepsilon_t(a, \mathbf{x})$ for learning θ_* and P , obtained as outcomes of the first step of the analyses, are carried over in the subsequent steps, where the optimization part is evaluated.

3 Description of the Adaptive Policy Considered

At each stage $t \geq 1$, the policy first updates an estimator $\hat{\theta}_{t-1}$ of θ_* based on the history h_{t-1} available so far, based on an adaptation of the Logistic-UCB1 algorithm by Faury et al. [2020], and deduces estimators $\hat{P}_{t-1}(a, \mathbf{x})$ and upper confidence bounds $U_{t-1}(a, \mathbf{x})$ of the probabilities $P(a, \mathbf{x})$. The policy then solves the corresponding estimated version of the optimization problem (2). We now describe the corresponding two steps. In the description below, quantities that depend on information available at round $t-1$ (respectively, t) are indexed by $t-1$ (respectively, t).

Phase 0: In case the cost constraints are about to be violated. To make sure cost constraints are never violated, whenever at least one of the components of the current cumulative costs is larger than $B-1$ and could possibly be larger than B at the end of round t , we play a_{null} (and we actually do so for the rest of the rounds). This corresponds to defining $\mathbf{p}_t(h_{t-1}, \mathbf{x}) = \delta_{a_{\text{null}}}$ for all $\mathbf{x} \in \mathcal{X}$, where $\delta_{a_{\text{null}}}$ denotes the Dirac mass on a_{null} . Otherwise, we proceed as described below in Phase 1 and Phase 2.

Phase 1: Learning θ_* via an adapted Logistic-UCB1. This first phase depends on a regularization parameter $\lambda > 0$ and on upper-confidence bonuses $\varepsilon_t(a, \mathbf{x}) > 0$, both to be specified by the analysis.

At rounds $t \geq 2$, we first maximize a regularized log-likelihood of the history h_{t-1} :

$$\tilde{\theta}_{t-1} \in \operatorname{argmax}_{\theta \in \mathbb{R}^m} \sum_{s=1}^{t-1} \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \left(y_s \ln \eta(\varphi(a_s, \mathbf{x}_s)^\top \theta) + (1-y_s) \ln(1-\eta(\varphi(a_s, \mathbf{x}_s)^\top \theta)) \right) - \frac{\lambda}{2} \|\theta\|^2. \quad (3)$$

In the expression above, we read that we only gather information about θ_* at those rounds s when $a_s \neq a_{\text{null}}$. When $\tilde{\theta}_{t-1}$ does not belong to Θ , an ad hoc projection step corrects for this, if needed:

$$\hat{\theta}_{t-1} \in \operatorname{argmin}_{\theta \in \Theta} \left\| \Psi_{t-1}(\theta) - \Psi_{t-1}(\tilde{\theta}_{t-1}) \right\|_{W_{t-1}(\theta)^{-1}}, \quad (4)$$

$$\begin{aligned} \text{where} \quad \Psi_{t-1}(\theta) &= \sum_{s=1}^{t-1} \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \eta(\varphi(a_s, \mathbf{x}_s)^\top \theta) \varphi(a_s, \mathbf{x}_s) + \lambda \theta \\ \text{and} \quad W_{t-1}(\theta) &= \lambda \mathbf{I}_m + \sum_{s=1}^{t-1} \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \dot{\eta}(\varphi(a_s, \mathbf{x}_s)^\top \theta) \varphi(a_s, \mathbf{x}_s) \varphi(a_s, \mathbf{x}_s)^\top. \end{aligned} \quad (5)$$

We recall that the function $\dot{\eta}$ denotes the derivative of η , i.e., $\dot{\eta}(x) = e^{-x}/(1+e^{-x})^2$. We have $\dot{\eta} = \eta(1-\eta)$.

By plug-in, we finally define estimators and upper-confidence bounds of the probabilities $P(a, \mathbf{x})$ for $a \neq a_{\text{null}}$ and all $\mathbf{x} \in \mathcal{X}$:

$$\hat{P}_{t-1}(a, \mathbf{x}) = \eta(\varphi(a, \mathbf{x})^\top \hat{\theta}_{t-1}) \quad \text{and} \quad U_{t-1}(a, \mathbf{x}) = \min\{\hat{P}_{t-1}(a, \mathbf{x}) + \varepsilon_{t-1}(a, \mathbf{x}), 1\}.$$

For a_{null} , no estimators or upper-confidence bounds need to be defined, as the quantities $P(a_{\text{null}}, \mathbf{x})$ are actually undefined.

Phase 2: Sampling, via solving an optimization problem with expected constraints. This phase relies on a conservative-budget parameter denoted by B_T , which is only slightly smaller than B and whose exact value is to be specified by the analysis.

We start with the case of a known context distribution ν . At round $t=1$, we play an arbitrary action in $\mathcal{A} \setminus \{a_{\text{null}}\}$. At rounds $t \geq 2$, if at least one component of the cumulative vector costs suffered so far is larger than $B-1$, we pick $a_t = a_{\text{null}}$. Otherwise, we pick for $\mathbf{p}_t(h_{t-1}, \cdot)$ the solution of the optimization problem $\operatorname{OPT}(\nu, U_{t-1}, B_T)$ defined¹ in (2), and draw a_t according to $\mathbf{p}_t(h_{t-1}, \mathbf{x}_t)$.

¹In the definition (2) of $\operatorname{OPT}(\nu, U_{t-1}, B_T)$, expectations are only over $\mathbf{X} \sim \nu$ and not over the random variable U_{t-1} ; more comments and explanations on this fact may be found in Appendix B.3.

BOX B: LOGISTIC-UCB1 FOR DIRECT SOLUTIONS TO OPT PROBLEMS

Parameters: regularization parameter $\lambda > 0$; conservative-budget parameter B_T ; upper-confidence bonuses $\varepsilon_s(a, \mathbf{x}) > 0$, for $s \geq 1$ and $(a, \mathbf{x}) \in (\mathcal{A} \setminus \{a_{\text{null}}\}) \times \mathcal{X}$.

Round $t = 1$: play an arbitrary action $a_1 \in \mathcal{A} \setminus \{a_{\text{null}}\}$

At rounds $t \geq 2$:

Phase 0 If $\sum_{s \leq t-1} c(a_s, \mathbf{x}_s) y_s \leq (B-1)\mathbf{1}$ is violated, then $\mathbf{p}_t(h_{t-1}, \mathbf{x}) = \delta_{a_{\text{null}}}$ for all \mathbf{x}

Phase 1 Otherwise, compute a maximum-likelihood estimator $\tilde{\boldsymbol{\theta}}_{t-1}$ of $\boldsymbol{\theta}_*$ according to (3), compute its projection $\hat{\boldsymbol{\theta}}_{t-1}$ onto Θ according to (4), and define, for $a \neq a_{\text{null}}$:

$$\hat{P}_{t-1}(a, \mathbf{x}) = \eta(\boldsymbol{\varphi}(a, \mathbf{x})^\top \hat{\boldsymbol{\theta}}_{t-1}) \quad \text{and} \quad U_{t-1}(a, \mathbf{x}) = \min\left\{\hat{P}_{t-1}(a, \mathbf{x}) + \varepsilon_{t-1}(a, \mathbf{x}), 1\right\}$$

Phase 2 Compute the solution $\mathbf{p}_t(h_{t-1}, \cdot)$ of

$$\begin{aligned} \text{OPT}(\tilde{\nu}, U_{t-1}, B_T) &= \max_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} T \mathbb{E}_{\mathbf{X} \sim \tilde{\nu}} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ \text{under} \quad & T \mathbb{E}_{\mathbf{X} \sim \tilde{\nu}} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \leq B_T \mathbf{1}, \end{aligned}$$

where $\tilde{\nu}$ denotes either ν (when it is known) or its empirical estimate $\hat{\nu}_t$ in (6)

Draw an arm $a_t \sim \mathbf{p}_t(h_{t-1}, \mathbf{x}_t)$.

When the context distribution is unknown, we rather pick for $\mathbf{p}_t(h_{t-1}, \cdot)$ the solution of the optimization problem $\text{OPT}(\hat{\nu}_t, U_{t-1}, B_T)$, where

$$\hat{\nu}_t = \frac{1}{t} \sum_{s=1}^t \delta_{\mathbf{x}_s}, \quad (6)$$

with $\delta_{\mathbf{x}}$ denoting the Dirac mass at $\mathbf{x} \in \mathcal{X}$. Since \mathbf{x}_t is revealed at the beginning of round t , before we pick an action, we may indeed use $\hat{\nu}_t$ at round t .

Summary and discussion of the computational complexity. We summarize the considered adaptive policy in Box B and now discuss its computational complexity.

As $\ln \varphi$ and $\ln(1 - \varphi)$ are strictly concave and smooth, the maximum-likelihood step (3) of Phase 1 consists of maximizing a strictly concave and smooth function over \mathbb{R}^m , which may be performed efficiently. The projection step (4) of Phase 1 is however an issue, both with the version of Logistic-UCB1 discussed here and with the earlier approach by Filippi et al. [2010, Section 3]. The latter and Fauray et al. [2020, Section 4.1] both underline that the projection step (4) is a complex optimization problem that however does not often need to be solved in practice, as they usually observe $\tilde{\boldsymbol{\theta}}_{t-1} \in \Theta$. Our numerical experiments concur with this statement (but admittedly, they rely on choosing a rather large value of Θ).

On the contrary, Phase 2 of the adaptive policy consists of solving a linear program with $|\mathcal{X}| \times |\mathcal{A}|$ constraints, where where $|\mathcal{X}|$ and $|\mathcal{A}|$ denote the cardinality of \mathcal{X} and \mathcal{A} , respectively—see the detailed rewriting (13) in the supplementary material. Therefore, the computational complexity of Phase 2 is polynomial (of weak order) in $|\mathcal{X}| \times |\mathcal{A}|$. To achieve this acceptable complexity we had however to restrict our attention to finite sets of contexts \mathcal{X} , which requires in practice segmenting countable or continuous context sets into finitely many clusters, for instance. We do so in our experiments.

Simulation study. A simulation study on partially simulated but realistic data may be found in Appendix F. The underlying dataset is the standard “default of credit card clients” dataset of UCI [2016], initially provided by Yeh and Lien [2009]. (It may be used under a Creative Commons Attribution 4.0 International [CC BY 4.0] license.)

4 Analysis for a Known Context Distribution ν

Since Θ is bounded, the following quantity, standardly introduced in the context of logistic bandits (see Faury et al. [2020] and references therein), is finite, though possibly large:

$$\kappa = \sup \left\{ \frac{1}{\dot{\eta}(\varphi(a, \mathbf{x})^\top \boldsymbol{\theta})} : \mathbf{x} \in \mathcal{X}, a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \boldsymbol{\theta} \in \Theta \right\} < +\infty.$$

We denote by $\|\Theta\| = \max\{\|\boldsymbol{\theta}\| : \boldsymbol{\theta} \in \Theta\}$ the maximal Euclidean norm of an element in Θ .

By construction, given that individual cost vectors lie in $[0, 1]^d$ and due to its ‘‘Phase 0’’, the adaptive policy considered always satisfies the budget constraints. The bound on rewards reads as follows.

Theorem 1. *In the setting of Box A of Section 2.1, we consider the adaptive policy of Box B of Section 3 assuming that the distribution of the contexts is known, i.e., with $\tilde{\nu} = \nu$. We set a confidence level $1 - \delta \in (0, 1)$ and use parameters $\lambda = m \ln(1 + T/m)$,*

$$B_T = B - 2 - \sqrt{2T \ln(4d/\delta)},$$

and $\varepsilon_t(a, \mathbf{x})$ stated in (9) of the supplementary material. Then, provided that $T \geq 2m$ and $B > 4 + 2\sqrt{2T \ln(4d/\delta)}$, we have, with probability at least $1 - 2\delta$,

$$\text{OPT}(\nu, P, B) - \sum_{t \leq T} r(a_t, \mathbf{x}_t) y_t \leq \left(4 + 2\sqrt{2T \ln \frac{4d}{\delta}} \right) \frac{\text{OPT}(\nu, P, B)}{B} + E_T + \sqrt{2T \ln \frac{4}{\delta}} + 1,$$

where the closed-form expression of $E_T = \mathcal{O}(m\sqrt{T} \ln T)$ is in (35) of the supplementary material.

We will rather discuss the bound of the more general Theorem 2 (to be stated and proved in Section 5) than the one of Theorem 1. We provide a proof sketch in Section 4.1 and discuss the main technical novelty in Section 4.2.

4.1 Proof Sketch for Theorem 1

The detailed proof of Theorem 1 may be found in Appendix B. We provide here an overview thereof, highlighting the four main ingredients. The third and fourth steps benefited from some inspiration drawn from the proof techniques of Agrawal and Devanur [2016]. The first step is an adaptation of Lemmas 1 and 2 by Faury et al. [2020].

First, the mentioned adaptation provides values of the parameters $\varepsilon_t(a, \mathbf{x})$ such that, with probability at least $1 - \delta$,

$$\forall t \geq 1, \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \left| \widehat{P}_t(a, \mathbf{x}) - P(a, \mathbf{x}) \right| \leq \varepsilon_t(a, \mathbf{x}),$$

hence

$$U_t(a, \mathbf{x}) - 2\varepsilon_t(a, \mathbf{x}) \leq P(a, \mathbf{x}) \leq U_t(a, \mathbf{x}),$$

while $\sum_{t \leq T} \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}$ is of order \sqrt{T} up to poly-logarithmic terms.

Second, the Phase 2 formulation of the strategy, in a primal form, is equivalently restated in a dual form. For each round $t \geq 2$, strong duality holds and entails the existence of a vector $\boldsymbol{\beta}_t^{\text{budg},*} \in \mathbb{R}^d$ such that $\mathbf{p}_t(h_{t-1}, \cdot)$ may be identified as the argmax over $\pi : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$ of

$$\mathbb{E}_{\mathbf{X} \sim \nu} \left[T \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) + \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \beta_{\mathbf{x},a}^{\text{p-pos},*} \pi_a(\mathbf{x}) \right].$$

By exploiting the KKT conditions, we are able to get rid of the double sum above and finally get a \mathcal{X} -pointwise characterization of $\mathbf{p}_t(h_{t-1}, \cdot)$: for all $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} \mathbf{p}_t(h_{t-1}, \mathbf{x}) &\in \operatorname{argmax}_{\mathbf{q} \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) q_a \\ &= \operatorname{argmax}_{\mathbf{q} \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ U_{t-1}(a, \mathbf{x}) q_a. \end{aligned}$$

Non-negative parts $(\cdot)_+$ may be introduced thanks to the existence of the no-op action a_{null} . The distributions $\mathbf{p}_t(h_{t-1}, \mathbf{x})$ may therefore be interpreted as maximizing some upper-confidence bound on penalized gains (rewards minus some scalarized costs); the dual variables $\beta_t^{\text{budg},*}$ play a role similar to the Z parameter of Agrawal and Devanur [2016, Section 3.3] in terms of weighing gains versus costs. In passing, we also prove

$$\text{OPT}(\nu, U_{t-1}, B_T) \geq B_T (\beta_t^{\text{budg},*})^\top \mathbf{1}$$

based on the KKT conditions. The latter inequality is comparable in spirit to the bound of Agrawal and Devanur [2016, Corollary 3], relating Z to $\text{OPT}(\nu, P, B)/B$.

Third, for $t \geq 2$, whenever the policy $\mathbf{p}_t(h_{t-1}, \cdot)$ is obtained by solving the optimization problem $\text{OPT}(\nu, U_{t-1}, B_T)$ of Phase 2 and by independence of \mathbf{x}_t and h_{t-1} , we have

$$\begin{aligned} \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} &= \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \mathbf{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ &= \mathbb{E} [r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) \mid h_{t-1}]. \end{aligned}$$

Therefore, repeated applications of the Hoeffding-Azuma inequality and the inequalities of the first step entail that, up to quantities of the order of \sqrt{T} ,

$$\begin{aligned} \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} &\approx \sum_{t=2}^T r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) \\ &\lesssim \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbf{1}_{\{a_t \neq a_{\text{null}}\}} + \sum_{t=2}^T r(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t) \lesssim \sum_{t=2}^T r(a_t, \mathbf{x}_t) y_t. \end{aligned}$$

We thus only need to control $\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T}$, which may be assumed ≥ 0 .

The value $B_T = B - 2 - \sqrt{2T \ln(4d/\delta)}$ and similar Hoeffding-Azuma-based arguments show that with high probability, the budget limit $B - 1$ is indeed never reached and that we always compute $\mathbf{p}_t(h_{t-1}, \cdot)$ in the way indicated by Phase 2.

Fourth, we collect all bounds together. We start with

$$\sum_{t=2}^T \frac{B_T}{T} (\beta_t^{\text{budg},*})^\top \mathbf{1} \leq \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \leq \text{OPT}(\nu, P, B).$$

We then exploit the dual characterization of $\mathbf{p}_t(h_{t-1}, \cdot)$ and the control $P \leq U_{t-1}$ to get that with high probability, for all $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} \sum_{a \in \mathcal{A}} (r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x})) U_{t-1}(a, \mathbf{x}) \mathbf{p}_{t,a}(h_{t-1}, \mathbf{x}) \\ \geq \sum_{a \in \mathcal{A}} (r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x})) P(a, \mathbf{x}) \pi_a^*(\mathbf{x}). \end{aligned}$$

After integration over $\mathbf{X} \sim \nu$ and substituting of the definitions of π^* and $\mathbf{p}_{t,a}(h_{t-1}, \cdot)$, as well as the equality stemming from the KKT conditions, we have

$$\begin{aligned} &= \text{OPT}(\nu, U_{t-1}, B_T) / T \\ &\underbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \mathbf{p}_{t,a}(h_{t-1}, \mathbf{X}) \right]}_{\text{OPT}(\nu, U_{t-1}, B_T) / T} \\ &\quad - \underbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \mathbf{p}_{t,a}(h_{t-1}, \mathbf{X}) \right]}_{(B_T/T) (\beta_t^{\text{budg},*})^\top \mathbf{1}} \\ &\geq \underbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}_{=\text{OPT}(\nu, P, B) / T} - \underbrace{(\beta_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}_{\leq (B/T) \mathbf{1}}. \end{aligned}$$

Rearranging and summing over $2 \leq t \leq T$, we obtain

$$\sum_{t=2}^T \frac{\text{OPT}(\nu, P, B) - \text{OPT}(\nu, U_{t-1}, B_T)}{T} \leq \sum_{t=2}^T \frac{B - B_T}{T} (\beta_t^{\text{budg},*})^\top \mathbf{1} \leq \left(\frac{B}{B_T} - 1 \right) \text{OPT}(\nu, P, B),$$

where we substituted the first inequality stated in this fourth step. This concludes the proof.

4.2 Discussion on the Main Technical Novelties

As should be clear from the comments at the beginning of Section 4.1, the technical novelty lies in the second step of the proof of Theorem 1. On the one hand, we are able to directly analyze a strategy stated in a primal form, which is a more natural formulation. On the other hand, doing so, we are also able to avoid the issues that come with dual formulations, relying, e.g., on some critical parameter Z , as in Agrawal and Devanur [2016, Theorem 3], to trade off rewards and costs. This parameter Z should be of order OPT/B and has to be learned, e.g., through \sqrt{T} initial exploration rounds. (More details are to be found in Section E.3.) In our analysis, this parameter Z is superseded by dual optimal variables $\beta_t^{\text{budg},*} \geq \mathbf{0}$, that are only used in the analysis and not to state the policy, unlike in Agrawal and Devanur [2016]. Put differently, the clever use in this context of KKT conditions is the main technical novelty. On a side note, we are also able to take care in an explicit and detailed fashion of the no-op action a_{null} , whose specific treatment is often unaddressed in the literature.

5 Analysis for an Unknown Context Distribution ν

When the context distribution ν is unknown, we simply estimate it through its empirical frequencies (6). The regret bound is almost unchanged: an additional mild factor of, e.g., $2|\mathcal{X}|\sqrt{2T \ln(2T|\mathcal{X}|/\delta)}$ appears in the \sqrt{T} term multiplying $\text{OPT}(\nu, P, B)/B$. This term comes from some uniform deviation argument stated in (7) and 8.

Theorem 2. *In the setting of Box A of Section 2.1, we consider the adaptive policy of Box B of Section 3 with $\tilde{\nu} = \hat{\nu}_t$ at rounds $t \geq 2$. We set a confidence level $1 - \delta \in (0, 1)$ and use parameters $\lambda = m \ln(1 + T/m)$, a working budget of*

$$B - b_T, \quad \text{where} \quad b_T = 2 + \sqrt{2T \ln(4d/\delta)} + |\mathcal{X}|\sqrt{2T \ln(2T|\mathcal{X}|/\delta)},$$

and $\varepsilon_t(a, \mathbf{x})$ stated in (9) of the supplementary material. Then, provided that $T \geq 2m$ and $B > 2b_T$, we have, with probability at least $1 - 3\delta$,

$$\text{OPT}(\nu, P, B) - \sum_{t \leq T} r(a_t, \mathbf{x}_t) y_t \leq 2b_T \left(1 + \frac{\text{OPT}(\nu, P, B)}{B} \right) + E_T,$$

where the expression of $E_T = \mathcal{O}(m\sqrt{T} \ln T)$ may be found in (35) of the supplementary material.

The order of magnitude of the regret bound is $(m + |\mathcal{X}|\text{OPT}(\nu, P, B)/B)\sqrt{T} \ln T$, which is reminiscent of all known regret upper bounds for CBwK (e.g., the ones by Badanidiyuru et al. [2014] and Agrawal et al. [2016], for general CBwK, and Agrawal and Devanur [2016] for linear CBwK, see Section 6). The factor $|\mathcal{X}|$ may be improved, see below, but this is a detail. A discussion on exhibiting corresponding lower bounds is to be found at the end of Section 6.

A detailed proof of Theorem 2 is provided in Appendix D of the supplementary material. It follows closely the proof of Theorem 1, with modifications mostly consisting of relating quantities of the form

$$\mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} [f(\mathbf{X})] \text{ vs. } \mathbb{E}_{\mathbf{X} \sim \nu} [f(\mathbf{X})], \quad \text{where, e.g., } f(\mathbf{X}) = \sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \mathbf{p}_{t,a}(h_{t-1}, \mathbf{X}).$$

To do so, we use that for all functions $f : \mathcal{X} \rightarrow [0, 1]$,

$$\forall t \leq T, \quad \left| \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} [f(\mathbf{X})] - \mathbb{E}_{\mathbf{X} \sim \nu} [f(\mathbf{X})] \right| \leq \sum_{\mathbf{x} \in \mathcal{X}} |\hat{\nu}_t(\mathbf{x}) - \nu(\mathbf{x})| \stackrel{\text{def}}{=} \|\hat{\nu}_t - \nu\|_1, \quad (7)$$

where $\|\hat{\nu}_t - \nu\|_1$ is the total variation distance between $\hat{\nu}_t$ and ν . In Appendix D, we upper bound the latter, for the sake of simplicity, in a crude way by applying $T|\mathcal{X}|$ times the Hoeffding-Azuma inequality (once for each $1 \leq t \leq T$ and $\mathbf{x} \in \mathcal{X}$) and obtain that with probability at least $1 - \delta$,

$$\forall t \leq T, \quad \|\hat{\nu}_t - \nu\|_1 \leq |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}. \quad (8)$$

The $|\mathcal{X}|\sqrt{2T\ln(2T|\mathcal{X}|/\delta)}$ term in the regret bound of Theorem 2 appears as the sum over $t \leq T$ of the deviation bounds (8). The bounds (8) may actually be improved into bounds of the order of $\sqrt{|\mathcal{X}|/t}$, via some Cauchy-Schwarz bound and a deviation argument in Banach spaces by Pinelis [1994], or by more direct techniques described by Devroye [1983, Lemma 3] and Berend and Kontorovich [2012]. In any case, the regret bound of Theorem 2 automatically benefits from such improvements, by replacing the $2|\mathcal{X}|\sqrt{2T\ln(2T|\mathcal{X}|/\delta)}$ term in the bound therein by the (sum over $t \leq T$ of the) better uniform deviation bounds.

6 Extension to Linear Contextual Bandits with Knapsacks

This section is a brief summary of Appendix E. We explain therein how the adaptive policy of Box B may be adapted to the setting of linear CBwK, introduced by Agrawal and Devanur [2016], where the bounded rewards r_t and vector costs \mathbf{c}_t are independently generated at each round according to bounded distributions with respective expectations $\bar{r}(a_t, \mathbf{x}_t)$ and $\bar{\mathbf{c}}(a_t, \mathbf{x}_t)$, depending linearly on (a transfer function φ of) the contexts: for all $a \neq a_{\text{null}}$ and $\mathbf{x} \in \mathcal{X}$, for all components i of $\bar{\mathbf{c}}$,

$$\bar{r}(a, \mathbf{x}) = \varphi(a, \mathbf{x})^T \boldsymbol{\mu}_* \quad \text{and} \quad \bar{c}_i(a, \mathbf{x}) = \varphi(a, \mathbf{x})^T \boldsymbol{\theta}_{*,i}.$$

We consider the same benchmark $\text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B)$ as Agrawal and Devanur [2016] and are able to exhibit a similar $(\text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B)/B)m\sqrt{T}\ln T$ regret bound, with however a slight relaxation on the order of magnitude required for B . We do so with a strategy that we deem more direct and natural, inspired from the one of Box B, where in Phase 1 a LinUCB-type (Abbasi-Yadkori et al. [2011]) estimation of the parameters is performed, and where in Phase 2, a direct solution to an OPT problem with estimated parameters is performed. The parameters are upper-confidence functions U_{t-1} on \bar{r} and lower-confidence vector functions \mathbf{L}_{t-1} on $\bar{\mathbf{c}}$.

The main advantage of our approach in the case of linear contextual bandits is exactly as described in Section 4.2: avoiding the critical parameter Z of Agrawal and Devanur [2016, Theorem 3], which is used to trade off rewards and costs. The main limitation of our approach is the assumption of a finite context set \mathcal{X} , which is required to make the Phase-2 linear program tractable.

7 Future Work

We conclude this article with a list of issues to be further investigated.

First, as discussed in Section 2.2, the restrictions of finiteness should be alleviated: finiteness of the context set in the setting of this article, or finiteness of the set of benchmark policies in other settings (see Badanidiyuru et al., 2014 and Agrawal et al., 2016).

Second, we only dealt with \leq budget constraints (and do does the literature so far). Direct approaches to constraints of the form \geq remain to be further investigated.

A third series of questions to be clarifies concerns regret lower bounds, and more generally, the tightness of the results—in particular, the required conditions on budget sizes. Earlier references for contextual bandits with knapsacks did also not provide lower bounds statements that were simultaneously optimal (i.e., matching the obtained upper bounds) and general (i.e., valid for all problems with a given number of rounds T , a given budget B , and a given value OPT for the optimal expected reward achievable by a static policy). Badanidiyuru et al. [2014, comments after Theorem 1] merely indicates that the obtained regret upper bound is optimal in some regimes, e.g., when the budget B grows linearly with the number of rounds T . Agrawal and Devanur [2016, comments after Theorem 1] only compares the obtained regret upper bound to the case of no budget constraints. In particular, as far as the orders of magnitude in T are concerned, it is unclear whether the $(\text{OPT}/B)\sqrt{T}$ rates achieved (up to poly-logarithmic factors) in the present article and in the two mentioned references are optimal. These rates do not match the optimal rates in the case of no contexts, which were stated and proved by Badanidiyuru et al. [2013].

Acknowledgments and Disclosure of Funding

Zhen Li and Gilles Stoltz have no direct funding to acknowledge other than the salaries paid by their employers, BNP Paribas, CNRS, and HEC Paris. They have no competing interests to declare.

References

- Default of credit card clients. UCI Machine Learning Repository, 2016. URL <https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients>.
- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems (NeurIPS'11)*, volume 24, 2011.
- S. Agrawal and N. Devanur. Linear contextual bandits with knapsacks. In *Advances in Neural Information Processing Systems (NeurIPS'16)*, volume 29, 2016.
- S. Agrawal, N.R. Devanur, and L. Li. An efficient algorithm for contextual bandits with knapsacks, and an extension to concave objectives. In *Proceedings of the 29th Annual Conference on Learning Theory (COLT'16)*, volume PMLR:49, pages 4–18, 2016.
- A. Badanidiyuru, R. Kleinberg, and A. Slivkins. Bandits with knapsacks. In *IEEE 54th Annual Symposium on Foundations of Computer Science (FOCS'13)*, pages 207–216, 2013. Latest extended version available at arXiv:1305.2545, dated September 2017.
- A. Badanidiyuru, J. Langford, and A. Slivkins. Resourceful contextual bandits. In *Proceedings of the 27th Conference on Learning Theory (COLT'14)*, volume PMLR:35, pages 1109–1134, 2014.
- A. Badanidiyuru, R. Kleinberg, and A. Slivkins. Bandits with global convex constraints and objective. *Journal of the ACM*, 65(3):1–55, 2018.
- D. Berend and A. Kontorovich. On the convergence of the empirical distribution, 2012. Preprint, arXiv:1205.6711.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- T. Chen and C. Guestrin. XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, 2016.
- W. Chu, L. Li, L. Reyzin, and R. Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS'11)*, volume 15, pages 208–214, 2011.
- L. Devroye. The equivalence of weak, strong and complete convergence in l_1 for kernel density estimates. *Annals of Statistics*, 11(3):896–904, 1983.
- L. Faury, M. Abeille, C. Calauzenes, and O. Fercoq. Improved optimistic algorithms for logistic bandits. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*, volume PMLR:119, pages 3052–3060, 2020.
- S. Filippi, O. Cappé, A. Garivier, and C. Szepesvári. Parametric bandits: The generalized linear case. In *Advances in Neural Information Processing Systems (NeurIPS'10)*, volume 23, 2010.
- N. Immorlica, K.A. Sankararaman, R. Schapire, and A. Slivkins. Adversarial bandits with knapsacks. In *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS'19)*, pages 202–219, 2019.
- T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- X. Li, C. Sun, and Y. Ye. The symmetry between arms and knapsacks: A primal-dual approach for bandits with knapsacks. In *Proceedings of the 38th International Conference on Machine Learning (ICML'21)*, pages 6483–6492, 2021.
- S. Miao, Y. Wang, and J. Zhang. A general framework for resource constrained revenue management with demand learning and large action space. *Available at SSRN 3841273*, 2021.
- I. Pinelis. Optimum bounds for the distributions of martingales in Banach spaces. *Annals of Probability*, 22(4):1679–1706, 1994.
- H. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58(5):527–535, 1952.

- A. Slivkins. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12 (1-2):1–286, 2019.
- W.R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- Z. Xu and V.-A. Truong. Reoptimization algorithms for contextual bandits with knapsack constraints, 2019. URL <http://www.columbia.edu/~vt2196/OnlineLearningAllocation8.pdf>.
- I.C. Yeh and C.H. Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications*, 36(2):2473–2480, 2009.
- M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning (ICML'03)*, volume PMLR:119, pages 3052–3060, 2003.

Supplementary Material for “Contextual Bandits with Knapsacks for a Conversion Model”

The appendices of this article contain the following elements.

- Appendix A provides the industrial motivation behind the setting described in Section 2.1, namely, within the banking industry, a market share expansion for loans.
- Appendix B contains the proof of Theorem 1, i.e., the regret bound in case of a known distribution ν , except for a lemma on learning the parameter of logistic bandits, provided in Appendix C.
- Appendix C states and proves the indicated lemma; both the statement and the proof are mere adaptations of Faury et al. [2020, Lemmas 1 and 2].
- Appendix D contains the proof of Theorem 2, i.e., explains how to adapt the proof of Appendix B to the case of an unknown distribution ν .
- Appendix E details the claims of Section 6, i.e., the extension of the techniques introduced to the setting of linear CBwK.
- Appendix F reports a simulation study based on realistic data.

A Industrial Motivation: Market Share Expansion for Loans

We describe the industrial problem we faced and which led to the setting described in Section 2.1.

Incentives and discounts are common practices in many industries to achieve some business objectives; however, there is usually a limit on the number or/and total volume of discounts that can be granted, so companies need to select carefully who should receive them. We consider, for instance, the banking industry, with a business objective of market share expansion: achieving the highest possible volume of loan subscription (total subscribed amount). Note that in practice, all loan applications need to go first through a risk-assessment process, and offers are only made if the bank considers that the loan will not put the customer's solvability at risk. We assume that all the clients concerned here have already gone through this process and are eligible for getting a loan from a given bank. We formulate the problem in a sequential manner as follows.

At each round $t \geq 1$, a client asks for a credit product. Her/his characteristics are denoted by $\mathbf{x}_t \in \mathbb{R}^n$, and encompass the socio-demographic profile, the loan request (amount $x_{\text{am},t}$, duration $x_{\text{dur},t}$), etc. It is reasonable to model these characteristics as independent draws from a common (possibly unknown) distribution ν . Based on \mathbf{x}_t , the bank will suggest some standard interest rate $\text{ir}(\mathbf{x}_t)$ based on its pricing rules; the detail of the rules is not relevant and we assume that the underlying function ir is given. If client t accepts the offered rate $\text{ir}(\mathbf{x}_t)$ and subscribes to the loan, an event which we denote by $y_t = 1$ and call a conversion, the bank gets a sales performance (gain on volume) $x_{\text{am},t}$. Otherwise, the client declines the offer, which we denote by $y_t = 0$ and the bank gets a null reward.

Actually, to improve the chances of a conversion, the bank may also offer a discount $a_t \in (0, 1]$, or lack of discount $a_t = 0$, on the interest rate. If it offers a discount $a_t > 0$ and $y_t = 1$, the bank will suffer a loss of earnings, equal to $a_t \text{ir}(\mathbf{x}_t) \text{out}(\mathbf{x}_t)$, where $\text{out}(\mathbf{x}_t)$ denotes the total outstanding amounts. This loss of earnings is considered a promotion cost. These promotion costs are summed up to previous such costs and should usually not exceed a fixed-in-advance budget $B_2 > 0$. Also, there is usually a fixed-in-advance limit $B_1 > 0$ on the total number of clients who can subscribe with a discount.

Given that the customers' characteristics are i.i.d., it is indeed reasonable to assume that y_t follows some Bernoulli distribution with unknown probability $P(a_t, \mathbf{x}_t)$. Of course, the higher the discount, the higher the probability of a conversion.

We summarize the setting with the notation of Section 2.1. We assume that discounts are picked in a finite grid $\mathcal{D} = \{j/D : j \in \{0, \dots, D\}\}$, so that the action set equals $\mathcal{A} = \mathcal{D} \cup \{a_{\text{null}}\}$. At each round $t \geq 1$, given the customer's characteristics \mathbf{x}_t and the discount $a_t \in [0, 1]$ picked by the bank, the latter receives the following reward and suffers the following costs:

$$\begin{aligned} & r(a_t, \mathbf{x}_t) y_t, & \text{where} & \quad r : (a, \mathbf{x}) \mapsto x_{\text{am}}, \\ \text{and} & \quad c(a_t, \mathbf{x}_t) y_t, & \text{where} & \quad c : (a, \mathbf{x}) \mapsto (\mathbb{1}_{\{a \neq 0\}}, a \text{ir}(\mathbf{x}_t) \text{out}(\mathbf{x}_t)). \end{aligned}$$

The first component of the cumulative cost vector measures the total number subscriptions with discounts, and the second component reports the total promotion costs. The bank wants to enforce

$$\sum_{t=1}^T c(a_t, \mathbf{x}_t) y_t = \sum_{t=1}^T (\mathbb{1}_{\{a_t \neq 0\}}, a_t \text{ir}(\mathbf{x}_t) \text{out}(\mathbf{x}_t)) y_t \leq (B_1, B_2)$$

while maximizing the sum of the achieved rewards.

Normalizations in $[0, 1]$ both for rewards and cost components may be achieved by considering the maximal amount M_{am} and outstanding M_{out} that the bank would allow, and by considering

$$r : (a, \mathbf{x}) \mapsto x_{\text{am}}/M_{\text{am}} \quad \text{and} \quad c : (a, \mathbf{x}) \mapsto (\mathbb{1}_{\{a \neq 0\}}, a \text{ir}(\mathbf{x}_t) \text{out}(\mathbf{x}_t)/M_{\text{out}})$$

with the alternative budget $(B_1, B_2/M_{\text{out}})$. A single budget parameter $B = \min\{B_1, B_2/M_{\text{out}}\}$ may be considered by a final normalization: by dividing the first cost component by $B_1/B > 1$ if $B_1 > B$, or the second cost component by $(B_2/M_{\text{out}})/B > 1$ if $B_2/M_{\text{out}} > B$, respectively.

B Detailed Proof of the Regret Bound in Case of a Known Distribution ν : Proof of Theorem 1

The proof is divided into four steps.

B.1 First Step: Defining Confidence Intervals on the Probabilities $P(a, \mathbf{x})$

The keystone of this step is the following lemma, adapted from Faury et al. [2020]: it provides guarantees for the adapted version of Logistic-UCB1 defined in Phase 1 of the adaptive policy studied in this article. The lemma actually holds for any sampling strategy of the arms, not just the one used in Phase 2 of the adaptive policy.

The reasons of the adaptations, lying in different settings being considered, as well as a detailed proof, are provided in Appendix C. We recall that we denote

$$\kappa = \sup \left\{ \frac{1}{\dot{\eta}(\boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta})} : \mathbf{x} \in \mathcal{X}, a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \boldsymbol{\theta} \in \Theta \right\},$$

and that since $\eta = \eta(1 - \eta) \in [0, 1/4]$, we always have $\kappa \geq 4$. We also recall the notation for the maximal Euclidean norm of an element in Θ :

$$\|\Theta\| = \max_{\boldsymbol{\theta} \in \Theta} \|\boldsymbol{\theta}\|.$$

Lemma 1 (combination of Lemmas 1 and 2 of Faury et al. [2020] with minor adjustments, detailed in Appendix C). *Assume that $\kappa < +\infty$. Fix any sampling strategy and consider the version of Logistic-UCB1 given by (3)–(5). For all $\delta \in (0, 1)$, there exists an event $\mathcal{E}_{\text{prob}, \delta}$ with probability at least $1 - \delta$ and such that over $\mathcal{E}_{\text{prob}, \delta}$:*

$$\forall t \geq 1, \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \left| \widehat{P}_t(a, \mathbf{x}) - P(a, \mathbf{x}) \right| \leq \gamma_{t, \lambda, \delta} \sqrt{\kappa(\|\Theta\| + 1/2)} \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{V_t^{-1}},$$

$$\text{where} \quad \gamma_{t, \lambda, \delta} = \sqrt{\lambda}(\|\Theta\| + 1/2) + \frac{2}{\sqrt{\lambda}} \ln \left(\frac{2^m}{\delta} \left(1 + \frac{t}{4m\lambda} \right)^{m/2} \right)$$

$$\text{and} \quad V_t = \sum_{s=1}^t \boldsymbol{\varphi}(a_s, \mathbf{x}_s) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} + \kappa \lambda \mathbf{I}_m.$$

Associated confidence intervals for the $P(a, \mathbf{x})$. Thanks to this lemma, we consider the upper-confidence bonuses

$$\varepsilon_t(a, \mathbf{x}) = \gamma_{t, \lambda, \delta} \sqrt{\kappa(\|\Theta\| + 1/2)} \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{V_t^{-1}}. \quad (9)$$

On the event $\mathcal{E}_{\text{prob}, \delta}$ of Lemma 1, we have, for all $t \geq 1$, all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$, and all $\mathbf{x} \in \mathcal{X}$: on the one hand,

$$U_t(a, \mathbf{x}) = \min \left\{ \widehat{P}_t(a, \mathbf{x}) + \varepsilon_t(a, \mathbf{x}), 1 \right\} \geq \min \{ P(a, \mathbf{x}), 1 \} = P(a, \mathbf{x}) \quad (10)$$

and on the other hand,

$$U_t(a, \mathbf{x}) = \min \left\{ \widehat{P}_t(a, \mathbf{x}) + \varepsilon_t(a, \mathbf{x}), 1 \right\} \leq \widehat{P}_t(a, \mathbf{x}) + \varepsilon_t(a, \mathbf{x}) \leq P(a, \mathbf{x}) + 2\varepsilon_t(a, \mathbf{x}). \quad (11)$$

Control of the sum of upper-confidence bonuses. According to the proof sketch of Section 4, it only remains to control the sum of the upper-confidence bonuses at observed contexts \mathbf{x}_t and played actions a_t . Note that at rounds $t \geq 2$, we use the bonuses $\varepsilon_{t-1}(a, \mathbf{x})$. We prove that

$$2 \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \leq \underbrace{\gamma_{T, \lambda, \delta} \sqrt{\kappa(4\|\Theta\| + 2)} \sqrt{2mT \max \left\{ 1, \frac{1}{\kappa\lambda} \right\} \ln \left(1 + \frac{T}{\kappa\lambda m} \right)}}_{\stackrel{\text{def}}{=} E_T}. \quad (12)$$

To that end, we first note that $\gamma_{t-1,\lambda,\delta} \leq \gamma_{T,\lambda,\delta}$:

$$2 \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \leq \gamma_{T,\lambda,\delta} \sqrt{\kappa(4\|\Theta\| + 2)} \sum_{t=2}^T \|\varphi(a_t, \mathbf{x}_t)\|_{V_{t-1}^{-1}} \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}.$$

Given that $\|\varphi\| \leq 1$ by assumption, a direct application of Lemma 2 below (a classical result of basic algebra for linear bandits) ensures that

$$\sum_{t=1}^T \|\varphi(a_t, \mathbf{x}_t)\|_{V_{t-1}^{-1}}^2 \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \leq 2m \max \left\{ 1, \frac{1}{\kappa\lambda} \right\} \ln \left(1 + \frac{\sum_{t=1}^T \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}}{\kappa\lambda m} \right).$$

A Cauchy-Schwarz inequality thus entails

$$\begin{aligned} \sum_{t=1}^T \|\varphi(a_t, \mathbf{x}_t)\|_{V_{t-1}^{-1}} \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} &\leq \sqrt{\sum_{t=1}^T \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}} \sqrt{2m \max \left\{ 1, \frac{1}{\kappa\lambda} \right\} \ln \left(1 + \frac{\sum_{t=1}^T \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}}{\kappa\lambda m} \right)} \\ &\leq \sqrt{T} \sqrt{2m \max \left\{ 1, \frac{1}{\kappa\lambda} \right\} \ln \left(1 + \frac{T}{\kappa\lambda m} \right)}, \end{aligned}$$

which concludes this step.

Lemma 2 (Elliptic potential and determinant-trace inequality, cf. Lemmas 10 and 11 of Abbasi-Yadkori et al. [2011], Lemmas 15 and 16 of Faury et al. [2020]). *For all $\lambda > 0$ and all sequences $\mathbf{u}_1, \mathbf{u}_2, \dots$ of vectors in \mathbb{R}^m with $\|\mathbf{u}_s\| \leq 1$, defining $U_0 = \lambda \mathbf{I}_m$ and for $t \geq 1$,*

$$U_t = \lambda \mathbf{I}_m + \sum_{s=1}^t \mathbf{u}_s (\mathbf{u}_s)^\top,$$

we have, for all $\tau \geq 1$:

$$\sum_{t=1}^{\tau} \|\mathbf{u}_t\|_{U_{t-1}^{-1}}^2 \leq 2m \max \left\{ 1, \frac{1}{\lambda} \right\} \ln \left(1 + \frac{\tau}{\lambda m} \right).$$

B.2 Second Step: Dual Formulation of the Sampling Phase (Phase 2) and Consequences

In this step, we consider a round $t \geq 2$ for which the cost constraints of Phase 0 of the adaptive policy are not violated and the optimization problem $\text{OPT}(\nu, U_{t-1}, B_{t,T})$ is to be solved; its solution is the policy $\mathbf{p}_t(h_{t-1}, \cdot)$ used to sample a_t according to $\mathbf{p}_t(h_{t-1}, \mathbf{x}_t)$.

We first rewrite in its dual form the optimization problem $\text{OPT}(\nu, U_{t-1}, B_{t,T})$ and show that strong duality holds. As a consequence, there exists a vector $\beta_t^{\text{budg},*} \in \mathbb{R}^d$ such that $\mathbf{p}_t(h_{t-1}, \cdot)$ may be identified as

$$\operatorname{argmax}_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \mathbb{E}_{\mathbf{X} \sim \nu} \left[T \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) + \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \beta_{\mathbf{x},a}^{\text{pp-},*} \pi_a(\mathbf{x}) \right].$$

By exploiting the KKT conditions, we are able to get rid of the double sum above and finally get a \mathcal{X} -pointwise characterization of $\mathbf{p}_t(h_{t-1}, \cdot)$: for all $\mathbf{x} \in \mathcal{X}$,

$$\mathbf{p}_t(h_{t-1}, \mathbf{x}) \in \operatorname{argmax}_{q \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) q_a,$$

where, with no impact, we may replace the $r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x})$ by their non-negative parts. The distributions $\mathbf{p}_t(h_{t-1}, \mathbf{x})$ may therefore be interpreted as maximizing some upper-confidence bound on penalized gains (rewards minus some scalarized costs).

We also prove $\text{OPT}(\nu, U_{t-1}, B_T) \geq B_T (\beta_t^{\text{budg},*})^\top \mathbf{1}$ based on the KKT conditions.

Primal form of the optimization problem $\text{OPT}(\nu, U_{t-1}, B_T)$. Since \mathcal{X} is a finite set (this is actually the key place where we need this assumption), the optimization problem $\text{OPT}(\nu, U_{t-1}, B_T)$ may be stated as the *opposite* of

$$\begin{aligned} & \min_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} -T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ \text{under} \quad & T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \leq B_T \mathbf{1}, \\ & \forall a \in \mathcal{A}, \forall \mathbf{x} \in \mathcal{X}, \quad \pi_a(\mathbf{x}) \geq 0, \\ & \forall \mathbf{x} \in \mathcal{X}, \quad \sum_{a \in \mathcal{A}} \pi_a(\mathbf{x}) = 1. \end{aligned}$$

Thanks to the no-op action $a_{\text{null}} \in \mathcal{A}$, which is used to model abstention and results in null rewards and costs, i.e., $r(a_{\text{null}}, \mathbf{x}) = 0$ and $c(a_{\text{null}}, \mathbf{x}) = \mathbf{0}$ for all $\mathbf{x} \in \mathcal{X}$, we may relax the third constraint into

$$\forall \mathbf{x} \in \mathcal{X}, \quad \sum_{a \in \mathcal{A}} \pi_a(\mathbf{x}) \leq 1.$$

Indeed, any vector $(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}$ satisfying the constraint with ≤ 1 can be transformed into a vector $(\pi'_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}$ for which the expected reward and the first and second constraints remain identical while the third constraint is satisfied with $= 1$: by adding the necessary probability mass to the components $\pi_{a_{\text{null}}}(\mathbf{x})$.

In the sequel, we consider this primal problem with the ≤ 1 constraint:

$$\begin{aligned} -\text{OPT}(\nu, U_{t-1}, B_T) &= \min_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} -T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ \text{under} \quad & T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \leq B_T \mathbf{1}, \quad (13) \\ & \forall a \in \mathcal{A}, \forall \mathbf{x} \in \mathcal{X}, \quad \pi_a(\mathbf{x}) \geq 0, \\ & \forall \mathbf{x} \in \mathcal{X}, \quad \sum_{a \in \mathcal{A}} \pi_a(\mathbf{x}) \leq 1. \end{aligned}$$

It forms a convex optimization problem, as its objective and its constraints are all affine (Boyd and Vandenberghe [2004, Section 4.2.1]).

Lagrangian (dual) form of the optimization problem. Denote by β^{budg} , $\beta^{\text{p-pos}}$, $\beta^{\text{p-sum}}$ the vector dual variables associated with the constraints on budget [budg], non-negative probability [p-pos], and sum of probabilities [p-sum], respectively. The vectors $\beta^{\text{p-pos}}$ and $\beta^{\text{p-sum}}$ have components $\beta_{\mathbf{x},a}^{\text{p-sum}}$ and $\beta_{\mathbf{x},a}^{\text{p-pos}}$ indexed by $\mathbf{x} \in \mathcal{X}$ and $a \in \mathcal{A}$.

We define the Lagrangian associated with our primal problem:

$$\begin{aligned} & \mathcal{L}_t \left((\pi_a(\mathbf{x}))_{a,\mathbf{x}}, \beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}} \right) \\ &= -T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ &+ (\beta^{\text{budg}})^T \left(T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] - B_T \mathbf{1} \right) \quad (14) \\ &+ \sum_{\mathbf{x} \in \mathcal{X}} \beta_{\mathbf{x}}^{\text{p-sum}} \left(\sum_{a \in \mathcal{A}} \pi_a(\mathbf{x}) - 1 \right) - \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \beta_{\mathbf{x},a}^{\text{p-pos}} \pi_a(\mathbf{x}). \end{aligned}$$

The dual problem consists of maximizing

$$\inf_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} \mathcal{L}_t \left((\pi_a(\mathbf{x}))_{a,\mathbf{x}}, \beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}} \right)$$

under the constraints that all components of the $\beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}}$ are non-negative, which we denote in a vector-wise manner by

$$\beta^{\text{budg}} \geq \mathbf{0}, \quad \beta^{\text{p-sum}} \geq \mathbf{0}, \quad \beta^{\text{p-pos}} \geq \mathbf{0}.$$

Strong duality and consequences. We explain why the so-called Slater's condition (Boyd and Vandenberghe, 2004, Section 5.2.3) holds; it entails that the value $-\text{OPT}(\nu, U_{t-1}, B_T)$ of the primal problem equals the value of the Lagrangian dual problem. The primal problem is convex, all its constraints are affine, with domain $\mathbb{R}^{\mathcal{A} \times \mathcal{X}}$: Slater's condition therefore reduces to feasibility. And feasibility of the constraints is clear by taking Dirac masses on a_{null} , i.e., $\pi_a(\mathbf{x}) = \delta_{a_{\text{null}}}$ for all $\mathbf{x} \in \mathcal{X}$. Since the values of the primal and dual problems are clearly larger than $-T > -\infty$, Slater's condition also implies that the dual optimal value is achieved at a dual feasible set of parameters, i.e., that the constrained supremum and the infimum defining the dual problem are a minimum and a maximum, respectively. We may therefore summarize the consequences of Slater's condition as follows:

$$\begin{aligned} -\text{OPT}(\nu, U_{t-1}, B_T) = & \max_{\beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}}} \min_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} \mathcal{L}_t \left((\pi_a(\mathbf{x}))_{a, \mathbf{x}}, \beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}} \right) \\ & \text{under } \beta^{\text{budg}} \geq \mathbf{0}, \quad \beta^{\text{p-sum}} \geq \mathbf{0}, \quad \beta^{\text{p-pos}} \geq \mathbf{0}. \end{aligned}$$

Because of strong duality and the existence of a dual feasible set of parameters, the max-min above equals its min-max counterpart (Boyd and Vandenberghe [2004, Sections 5.4.1 and 5.4.2]):

$$-\text{OPT}(\nu, U_{t-1}, B_T) = \min_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} \max_{\beta^{\text{budg}} \geq \mathbf{0}, \beta^{\text{p-sum}} \geq \mathbf{0}, \beta^{\text{p-pos}} \geq \mathbf{0}} \mathcal{L}_t \left((\pi_a(\mathbf{x}))_{a, \mathbf{x}}, \beta^{\text{budg}}, \beta^{\text{p-sum}}, \beta^{\text{p-pos}} \right).$$

We let $\beta_t^{\text{budg},*} \geq \mathbf{0}$, $\beta_t^{\text{p-sum},*} \geq \mathbf{0}$, and $\beta_t^{\text{p-pos},*} \geq \mathbf{0}$ be an optimal dual solution and recall that $\mathbf{p}_t(h_{t-1}, \cdot)$ denote an optimal primal solution, which, with no loss of generality, may be assumed satisfying the ≤ 1 constraints with equality. From Boyd and Vandenberghe [2004, Section 5.4.2], this pair of solutions forms a saddle-point for the Lagrangian; in particular,

$$\begin{aligned} -\text{OPT}(\nu, U_{t-1}, B_T) &= \mathcal{L}_t(\mathbf{p}_t(h_{t-1}, \cdot), \beta_t^{\text{budg},*}, \beta_t^{\text{p-sum},*}, \beta_t^{\text{p-pos},*}) \\ &= \min_{(\pi_a(\mathbf{x})) \in \mathbb{R}^{\mathcal{A} \times \mathcal{X}}} \mathcal{L}_t \left((\pi_a(\mathbf{x})), \beta_t^{\text{budg},*}, \beta_t^{\text{p-sum},*}, \beta_t^{\text{p-pos},*} \right) \\ &= \min_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \mathcal{L}_t(\pi, \beta_t^{\text{budg},*}, \beta_t^{\text{p-sum},*}, \beta_t^{\text{p-pos},*}). \end{aligned} \tag{15}$$

The distribution $\mathbf{p}_t(h_{t-1}, \cdot)$ played thus appears as the argument of the minimum above.

Substituting the definition (14) of \mathcal{L}_t into the characterization (15), rearranging the first two terms of \mathcal{L}_t , noting that the third term of \mathcal{L}_t is null, and discarding the constant term $B_T (\beta_t^{\text{budg},*})^\top \mathbf{1}$, we get:

$$\begin{aligned} & \mathbf{p}_t(h_{t-1}, \cdot) \in \\ & \underset{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})}{\text{argmin}} \mathbb{E}_{\mathbf{X} \sim \nu} \left[-T \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) + (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) - \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \beta_{\mathbf{x}, a}^{\text{p-pos},*} \pi_a(\mathbf{x}) \right] \\ = & \underset{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})}{\text{argmax}} \mathbb{E}_{\mathbf{X} \sim \nu} \left[T \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) + \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \beta_{\mathbf{x}, a}^{\text{p-pos},*} \pi_a(\mathbf{x}) \right]. \end{aligned} \tag{16}$$

We further simplify this alternative definition by showing that the sums of the $\beta_{\mathbf{x}, a}^{\text{p-pos},*} \pi_a(\mathbf{x})$ may be omitted. We do so by exploiting the KKT conditions.

KKT conditions: statement. The Karush–Kuhn–Tucker (KKT) conditions (Boyd and Vandenberghe [2004, Section 5.5.3]) for the primal optimal $\mathbf{p}_t(h_{t-1}, \cdot): \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$ and the dual optimal $\beta_t^{\text{budg},*} \geq \mathbf{0}$, $\beta_t^{\text{p-sum},*} \geq \mathbf{0}$, and $\beta_t^{\text{p-pos},*} \geq \mathbf{0}$ imply the following conditions: first, complementary slackness, which reads

$$\forall a \in \mathcal{A}, \forall \mathbf{x} \in \mathcal{X}, \quad \beta_{t, \mathbf{x}, a}^{\text{p-pos},*} p_{t, a}(h_{t-1}, \mathbf{x}) = 0 \tag{17}$$

and

$$(\beta_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t, a}(h_{t-1}, \mathbf{X}) \right] = \frac{B_T}{T} (\beta_t^{\text{budg},*})^\top \mathbf{1}; \tag{18}$$

second, a stationary condition, based on the fact that the gradient of the Lagrangian function (14) with respect to the $\pi_a(\mathbf{x})$ vanishes. Denoting by $\nu(\mathbf{x})$ the probability mass put by ν on \mathbf{x} , we have:

$$\begin{aligned} \forall a \in \mathcal{A}, \quad \forall \mathbf{x} \in \mathcal{X}, \quad & T r(a, \mathbf{x}) U_{t-1}(a, \mathbf{x}) \nu(\mathbf{x}) \\ &= (\beta_t^{\text{budg},*})^\top \left(T \mathbf{c}(a, \mathbf{x}) U_{t-1}(a, \mathbf{x}) \nu(\mathbf{x}) \right) + \beta_{t,\mathbf{x}}^{\text{p-sum},*} - \beta_{t,\mathbf{x},a}^{\text{p-pos},*}. \end{aligned} \quad (19)$$

KKT conditions: first consequence—final characterization of $\mathbf{p}_t(h_{t-1}, \cdot)$. For now, we only exploit (17): this equality and the fact that $\beta_{t,\mathbf{x},a}^{\text{p-pos},*} \pi_a(\mathbf{x})$ is always non-negative show that, as announced, the characterization (16) may be further simplified into

$$\begin{aligned} \mathbf{p}_t(h_{t-1}, \cdot) &\in \operatorname{argmax}_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \mathbb{E}_{\mathbf{X} \sim \nu} \left[T \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ &= \operatorname{argmax}_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right]. \end{aligned} \quad (20)$$

In the characterization (20), as we got rid of the cross terms, the maximization may be carried out in a \mathcal{X} -pointwise manner, i.e., by separately computing each probability distribution $\mathbf{p}_t(h_{t-1}, \mathbf{x}) \in \mathcal{P}(\mathcal{A})$. More formally, for each $\mathbf{x} \in \mathcal{X}$,

$$\mathbf{p}_t(h_{t-1}, \mathbf{x}) \in \operatorname{argmax}_{q \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) q_a. \quad (21)$$

In this characterization, $r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x})$ appears as a penalized gain in case a client with characteristics \mathbf{x} converts, and $\mathbf{p}_t(h_{t-1}, \mathbf{x})$ is obtained by combining an upper-confidence-bound estimation $U_{t-1}(a, \mathbf{x})$ of the conversion rate $P(a, \mathbf{x})$ with this penalized gain.

Denote by $(z)_+ = \max\{z, 0\}$ the non-negative part of $z \in \mathbb{R}$ and fix $\mathbf{x} \in \mathcal{X}$. Given that the no-op action a_{null} is such that $r(a_{\text{null}}, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a_{\text{null}}, \mathbf{x}) = 0$ and $U_{t-1}(a, \mathbf{x}) \geq 0$ for all $a \in \mathcal{A}$, in view of the objective, any distribution q should move the probability mass q_a on an action $a \in \mathcal{A}$ with $r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) < 0$ to a_{null} . As a consequence,

$$\begin{aligned} \max_{q \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) q_a(\mathbf{x}) \\ = \max_{q \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ U_{t-1}(a, \mathbf{x}) q_a(\mathbf{x}) \end{aligned} \quad (22)$$

and

$$\mathbf{p}_t(h_{t-1}, \mathbf{x}) \in \operatorname{argmax}_{q \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\beta_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ U_{t-1}(a, \mathbf{x}) q_a. \quad (23)$$

KKT conditions: a second consequence. We first exploit the stationary condition (19). Multiplying both sides of this equality by $p_{t,a}(h_{t-1}, \mathbf{x})$, summing over $\mathbf{x} \in \mathcal{X}$ and $a \in \mathcal{A}$, we obtain an equality between expectations with respect to $\mathbf{X} \sim \nu$:

$$\begin{aligned} & T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ &= (\beta_t^{\text{budg},*})^\top \left(T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] \right) \\ &+ \sum_{\mathbf{x} \in \mathcal{X}} \beta_{t,\mathbf{x}}^{\text{p-sum},*} \underbrace{\sum_{a \in \mathcal{A}} p_{t,a}(h_{t-1}, \mathbf{x})}_{=1} - \sum_{\mathbf{x} \in \mathcal{X}} \sum_{a \in \mathcal{A}} \underbrace{\beta_{t,\mathbf{x},a}^{\text{p-pos},*} p_{t,a}(h_{t-1}, \mathbf{x})}_{=0}, \end{aligned} \quad (24)$$

where the equality to 0 indicated in the right-hand side correspond to (17). We now substitute (18) into (24) and obtain

$$T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] = B_T (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{1} + \sum_{\mathbf{x} \in \mathcal{X}} \beta_{t,\mathbf{x}}^{\text{p-sum},*}.$$

The left-hand side equals $\text{OPT}(\nu, U_{t-1}, B_T)$ since $\mathbf{p}_t(h_{t-1}, \cdot)$ is the solution of the primal problem. Thus, the equality above entails, by $\boldsymbol{\beta}_t^{\text{p-sum},*} \geq \mathbf{0}$, that

$$\text{OPT}(\nu, U_{t-1}, B_T) \geq B_T (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{1}. \quad (25)$$

B.3 Third Step: Various High-Probability Controls

We prove below that on the intersection between the event $\mathcal{E}_{\text{prob},\delta}$ of Lemma 1 and another event $\mathcal{E}_{\text{HAz},\delta}$, also of probability at least $1 - \delta$, we have simultaneously that for all rounds $t \geq 2$, the policy $\mathbf{p}_t(h_{t-1}, \cdot)$ is obtained by Phase 2, i.e., by solving $\text{OPT}(\nu, U_{t-1}, B_T)$, and that

$$\begin{aligned} \sum_{t=1}^T r(a_t, \mathbf{x}_t) y_t &\geq \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} - \sqrt{2T \ln \frac{4}{\delta}} - E_T \quad (26) \\ \text{and} \quad \sum_{t=1}^T \mathbf{c}(a_t, \mathbf{x}_t) y_t &\leq \left(B_T + 1 + \sqrt{2T \ln \frac{4d}{\delta}} \right) \mathbf{1} = (B - 1) \mathbf{1}, \end{aligned}$$

where the bound E_T was defined in (12) and where we used the definition of B_T , namely,

$$B_T = B - 2 - \sqrt{2T \ln \frac{4d}{\delta}}. \quad (27)$$

This definition requires that

$$B > c \left(2 + \sqrt{2T \ln \frac{4d}{\delta}} \right) \quad (28)$$

for $c = 1$, but we will rather assume that the inequality holds with $c = 2$.

Applications of the Hoeffding-Azuma inequality to handle the conversions y_t . We recall that we defined h_0 as the empty vector and $h_t = (\mathbf{x}_s, a_s, y_s)_{s \leq t}$ for $t \geq 1$. We introduce the filtration $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$, with $\mathcal{F}_0 = \sigma(a_1, \mathbf{x}_1)$ and for $t \geq 1$,

$$\mathcal{F}_t = \sigma(h_t, a_{t+1}, \mathbf{x}_{t+1}).$$

Then, for all $t \geq 1$, the variables $r(a_t, \mathbf{x}_t) y_t$ and $\mathbf{c}(a_t, \mathbf{x}_t) y_t$ are \mathcal{F}_t -measurable, with conditional expectations with respect to \mathcal{F}_{t-1} equal to

$$\begin{aligned} \mathbb{E}[r(a_t, \mathbf{x}_t) y_t \mid \mathcal{F}_{t-1}] &= r(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t) \\ \text{and} \quad \mathbb{E}[\mathbf{c}(a_t, \mathbf{x}_t) y_t \mid \mathcal{F}_{t-1}] &= \mathbf{c}(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t). \end{aligned}$$

Indeed, the conditioning by \mathcal{F}_{t-1} fixes a_t and \mathbf{x}_t , but not y_t , and exactly means, when $a_t \neq a_{\text{null}}$, integrating over $y_t \sim P(a_t, \mathbf{x}_t)$. When $a_t = a_{\text{null}}$, all equalities above remain valid thanks to the abuse of notation discussed in Section 2.1. Given that r takes values in $[0, 1]$ and \mathbf{c} in $[0, 1]^d$, we may apply $d + 1$ times the Hoeffding-Azuma inequality (once for r and each component of \mathbf{c}) together with a union bound: there exist two events $\mathcal{E}_{r,P,\delta}$ and $\mathcal{E}_{\mathbf{c},P,\delta}$, each of probability at least $1 - \delta/4$, such that

$$\begin{aligned} \text{on } \mathcal{E}_{r,P,\delta}, \quad \sum_{t=1}^T r(a_t, \mathbf{x}_t) y_t &\geq \sum_{t=1}^T r(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t) - \sqrt{\frac{T}{2} \ln \frac{4}{\delta}} \\ \text{and} \quad \text{on } \mathcal{E}_{\mathbf{c},P,\delta}, \quad \sum_{t=1}^T \mathbf{c}(a_t, \mathbf{x}_t) y_t &\leq \sqrt{\frac{T}{2} \ln \frac{4d}{\delta}} \mathbf{1} + \sum_{t=1}^T \mathbf{c}(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t). \end{aligned}$$

Applications of the Hoeffding-Azuma inequality to use the properties of the $p_t(h_{t-1}, \cdot)$. For this sub-step, we rather condition directly by h_{t-1} (instead of \mathcal{F}_{t-1} as in the previous step), where $t \geq 2$. Conditioning by h_{t-1} amounts to integrating first over $a_t \sim p_t(h_{t-1}, \mathbf{x}_t)$ and then over $\mathbf{x}_t \sim \nu$: this is because of the definition of the random draw of a_t according to $p_t(h_{t-1}, \mathbf{x}_t)$ independently from everything else, and because \mathbf{x}_t is drawn independent from the past according to ν . More precisely, for each $t \geq 2$, given that U_{t-1} is h_{t-1} -measurable, we thus have the following equalities:

$$\begin{aligned} \mathbb{E}[r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) \mid h_{t-1}] &= \mathbb{E}\left[\sum_{a \in \mathcal{A}} r(a, \mathbf{x}_t) U_{t-1}(a, \mathbf{x}_t) p_{t,a}(h_{t-1}, \mathbf{x}_t) \mid h_{t-1}\right] \\ &= \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ \text{and} \quad \mathbb{E}[c(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) \mid h_{t-1}] &= \mathbb{E}\left[\sum_{a \in \mathcal{A}} c(a, \mathbf{x}_t) U_{t-1}(a, \mathbf{x}_t) p_{t,a}(h_{t-1}, \mathbf{x}_t) \mid h_{t-1}\right] \\ &= \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right], \end{aligned}$$

where we recall that $\mathbb{E}_{\mathbf{X} \sim \nu}$ denotes an integration solely over $\mathbf{X} \sim \nu$.

By definition of $p_t(h_{t-1}, \cdot)$, whenever the adaptive policy reaches Phase 2 at a given round $t \geq 2$:

$$\begin{aligned} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] &= \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \\ \text{and} \quad \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] &\leq \frac{B_T}{T} \mathbf{1}. \end{aligned}$$

Otherwise, the adaptive policy is stuck in Phase 0, because at least one cost constraint is larger than $B - 1$; in this case, $p_t(h_{t-1}, \mathbf{x}) = \delta_{a_{\text{null}}}$ for all $\mathbf{x} \in \mathcal{X}$ and both expectations above are null. We may summarize all cases by introducing the indicator function that all cost constraints are smaller than $B - 1$,

$$\mathbf{1}_{\{C_{t-1} \leq (B-1)\mathbf{1}\}}, \quad \text{where} \quad C_{t-1} = \sum_{s=1}^{t-1} c(a_s, \mathbf{x}_s) y_s,$$

and stating that

$$\begin{aligned} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] &\geq \mathbf{1}_{\{C_{t-1} \leq (B-1)\mathbf{1}\}} \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \\ \text{and} \quad \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] &\leq \frac{B_T}{T} \mathbf{1}. \end{aligned}$$

Therefore, a second series of applications of the Hoeffding-Azuma inequality (at times $2 \leq t \leq T$, i.e., excluding the first round) entails the existence of two events $\mathcal{E}_{r,U,\delta}$ and $\mathcal{E}_{c,U,\delta}$, each of probability at least $1 - \delta/4$, such that

$$\begin{aligned} \text{on } \mathcal{E}_{r,U,\delta}, \quad \sum_{t=2}^T r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) &\geq \sum_{t=2}^T \mathbf{1}_{\{C_{t-1} \leq (B-1)\mathbf{1}\}} \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} - \sqrt{\frac{T}{2} \ln \frac{4}{\delta}} \\ \text{and on } \mathcal{E}_{c,U,\delta}, \quad \sum_{t=2}^T c(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) &\leq \left(B_T + \sqrt{\frac{T}{2} \ln \frac{4d}{\delta}} \right) \mathbf{1}. \end{aligned}$$

Appeal to results of Appendix B.1. The inequalities (10) and (11) state that on the event $\mathcal{E}_{\text{prob},\delta}$ of Lemma 1, we have

$$\forall t \geq 1, \quad \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \quad \forall \mathbf{x} \in \mathcal{X}, \quad P(a, \mathbf{x}) \leq U_t(a, \mathbf{x}) \leq P(a, \mathbf{x}) + 2\varepsilon_t(a, \mathbf{x}).$$

Thus, on $\mathcal{E}_{\text{prob},\delta}$,

$$\begin{aligned} \sum_{t=1}^T r(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t) &\geq \sum_{t=2}^T r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) - 2 \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \\ &\geq \sum_{t=2}^T r(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t) - E_T \end{aligned}$$

$$\text{and } \sum_{t=1}^T \mathbf{c}(a_t, \mathbf{x}_t) P(a_t, \mathbf{x}_t) \leq \mathbf{1} + \sum_{t=2}^T \mathbf{c}(a_t, \mathbf{x}_t) U_{t-1}(a_t, \mathbf{x}_t),$$

where we recall that the bound E_T was defined in (12).

Conclusion of this step. We define

$$\mathcal{E}_{\text{HAz},\delta} = \mathcal{E}_{r,P,\delta} \cap \mathcal{E}_{c,P,\delta} \cap \mathcal{E}_{r,U,\delta} \cap \mathcal{E}_{c,U,\delta},$$

which is an event of probability at least $1 - \delta$. On the intersection of $\mathcal{E}_{\text{HAz},\delta}$ and $\mathcal{E}_{\text{prob},\delta}$, by collecting all bounds together, we have

$$\sum_{t=1}^T \mathbf{c}(a_t, \mathbf{x}_t) y_t \leq \left(B_T + 1 + \sqrt{2T \ln \frac{4d}{\delta}} \right) \mathbf{1} = (B - 1) \mathbf{1}.$$

This shows that on the intersection of $\mathcal{E}_{\text{HAz},\delta}$ and $\mathcal{E}_{\text{prob},\delta}$, the indicator functions $\mathbb{1}_{\{C_{t-1} \leq (B-1)\mathbf{1}\}}$ all equal 1, and that the policies $\mathbf{p}_t(h_{t-1}, \cdot)$ are always obtained by solving the optimization problems of Phase 2. We conclude this step by collecting the obtained bounds for rewards and by legitimately replacing the indicator functions therein by 1.

B.4 Fourth Step: Conclusion

In this step, we merely combine the bounds exhibited in the first three steps to obtain the closed-form expression of the regret bound. We then propose suitable orders of magnitude for the parameters.

Collecting all bounds to get a closed-form regret bound. By considering $\mathbf{q} = \pi^*(\mathbf{x})$ for each $\mathbf{x} \in \mathcal{X}$ in (23), where we recall that π^* is the optimal static policy, and by the equality (22), we note that, for all $t \geq 2$,

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{X}, \quad \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) p_{t,a}(h_{t-1}, \mathbf{x}) \\ \geq \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ U_{t-1}(a, \mathbf{x}) \pi_a^*(\mathbf{x}). \end{aligned}$$

On the event $\mathcal{E}_{\text{prob},\delta}$ of Lemma 1, the inequality (10) states that $U_{t-1}(a, \mathbf{x}) \geq P(a, \mathbf{x})$, which we may substitute in the inequality above (thanks to the non-negative part in the right-hand side) to get

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{X}, \quad \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) p_{t,a}(h_{t-1}, \mathbf{x}) \\ \geq \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ P(a, \mathbf{x}) \pi_a^*(\mathbf{x}) \\ \geq \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) P(a, \mathbf{x}) \pi_a^*(\mathbf{x}). \end{aligned}$$

We replace the individual \mathbf{x} by a random variable $\mathbf{X} \sim \nu$ and integrate over \mathbf{X} : on $\mathcal{E}_{\text{prob},\delta}$,

$$\begin{aligned} & \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ & \quad - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ \geq & \underbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}_{=\text{OPT}(\nu, P, B)/T} - (\boldsymbol{\beta}_t^{\text{budg},*})^\top \underbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}_{\leq (B/T)\mathbf{1}}. \end{aligned} \quad (29)$$

The equality to $\text{OPT}(\nu, P, B)/T$ and the inequality $\leq (B/T)\mathbf{1}$ above come from the very definition of π^* as the static policy solving $\text{OPT}(\nu, P, B)$. Similarly, Appendix B.3 shows that on the event $\mathcal{E}_{\text{HAz},\delta}$, for all $2 \leq t \leq T$, the policies $p_{t,a}(h_{t-1}, \cdot)$ are obtained by solving $\text{OPT}(\nu, U_{t-1}, B_T)$, so that, by definition,

$$\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) p_{t,a}(h_{t-1}, \mathbf{X}) \right] = \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T}.$$

Substituting this equality and the KKT condition (18) into (29) and rearranging, we see that we proved so far that on the intersection of $\mathcal{E}_{\text{prob},\delta}$ and $\mathcal{E}_{\text{HAz},\delta}$,

$$\forall 2 \leq t \leq T, \quad \text{OPT}(\nu, P, B) - \text{OPT}(\nu, U_{t-1}, B_T) \leq (B - B_T) (\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{1}.$$

We may also substitute (25), i.e., $(\boldsymbol{\beta}_t^{\text{budg},*})^\top \mathbf{1} \leq \text{OPT}(\nu, U_{t-1}, B_T)/B_T$ and finally get that on the intersection of $\mathcal{E}_{\text{prob},\delta}$ and $\mathcal{E}_{\text{HAz},\delta}$,

$$\forall 2 \leq t \leq T, \quad \text{OPT}(\nu, P, B) - \text{OPT}(\nu, U_{t-1}, B_T) \leq \left(\frac{B}{B_T} - 1 \right) \text{OPT}(\nu, U_{t-1}, B_T).$$

Summing over $2 \leq t \leq T$ and using that $\text{OPT}(\nu, P, B)/T \leq 1$ by definition and by the fact that r takes values in $[0, 1]$, we obtain

$$\begin{aligned} \text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} & \leq 1 + \sum_{t=2}^T \frac{\text{OPT}(\nu, P, B) - \text{OPT}(\nu, U_{t-1}, B_T)}{T} \\ & \leq 1 + \left(\frac{B}{B_T} - 1 \right) \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T}. \end{aligned} \quad (30)$$

Distinguishing the cases

$$\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \leq 0 \quad \text{and} \quad \text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \geq 0,$$

we see, based on (30), that in both cases

$$\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\nu, U_{t-1}, B_T)}{T} \leq 1 + \left(\frac{B}{B_T} - 1 \right) \text{OPT}(\nu, P, B).$$

We finally substitute (26) and have proved, as claimed, that on the intersection of $\mathcal{E}_{\text{prob},\delta}$ and $\mathcal{E}_{\text{HAz},\delta}$,

$$\text{OPT}(\nu, P, B) - \sum_{t=1}^T r(a_t, \mathbf{x}_t) y_t \leq 1 + \left(\frac{B}{B_T} - 1 \right) \text{OPT}(\nu, P, B) + E_T + \sqrt{2T \ln \frac{4}{\delta}}, \quad (31)$$

where we recall that E_T was defined in (12).

Improving readability and setting λ . As indicated, we require (28) with $c = 2$. By the definition (27) of B_T and the inequality $1/(1-u) \leq 1+2u$ for $u \in (0, 1/2)$, we have

$$\frac{B}{B_T} - 1 \leq \frac{2(2 + \sqrt{2T \ln(4d/\delta)})}{B}, \quad (32)$$

which we may substitute in (31). It only remains to deal with a bound on E_T to conclude the proof of Theorem 1.

We set below a value of λ larger than 1. Recalling that $\kappa \geq 4$, we may already bound E_T as

$$E_T \leq \gamma_{T,\lambda,\delta} \sqrt{\kappa(4\|\Theta\| + 2)} \sqrt{2mT \ln\left(1 + \frac{T}{4m}\right)} = 2\gamma_{T,\lambda,\delta} \sqrt{\kappa(2\|\Theta\| + 1)} \sqrt{mT \ln\left(1 + \frac{T}{4m}\right)}, \quad (33)$$

where $\gamma_{T,\lambda,\delta}$, defined in the statement of Lemma 1, may itself be bounded by

$$\gamma_{T,\lambda,\delta} \leq \sqrt{\lambda}(\|\Theta\| + 1/2) + \frac{2}{\sqrt{\lambda}} \ln\left(\frac{2^m}{\delta} \left(1 + \frac{T}{4m}\right)^{m/2}\right).$$

For the sake of simplicity, we set the value of λ by only optimizing the orders of magnitude in m and T of (this upper bound on) $\gamma_{T,\lambda,\delta}$, i.e., by considering

$$\sqrt{\lambda} + \frac{m}{\sqrt{\lambda}} \ln T.$$

We take $\lambda = m \ln(1 + T/m)$, which is indeed larger than 1 given that $T \geq 2m$ and $\ln(1 + T/m) \geq 1$. We have

$$\frac{2}{\sqrt{\lambda}} \ln \frac{1}{\delta} \leq \ln \frac{1}{\delta} \quad \text{and} \quad \frac{2}{\sqrt{\lambda}} \ln 2^m \leq (2 \ln 2) \sqrt{m} \leq 2\sqrt{m},$$

as well as

$$\frac{2}{\sqrt{\lambda}} \ln\left(\left(1 + \frac{T}{4m}\right)^{m/2}\right) \leq \frac{\sqrt{m}}{\sqrt{\ln(1 + T/m)}} \ln\left(1 + \frac{T}{4m}\right) \leq \sqrt{m \ln\left(1 + \frac{T}{4m}\right)}.$$

All in all,

$$\gamma_{T,\lambda,\delta} \leq (\sqrt{m} + \|\Theta\| + 1/2) \sqrt{\ln\left(1 + \frac{T}{m}\right)} + 2\sqrt{m} + \ln \frac{1}{\delta}. \quad (34)$$

Combining (33) and (34), we showed:

$$E_T \leq 2\sqrt{\kappa(2\|\Theta\| + 1)} \sqrt{mT \ln\left(1 + \frac{T}{4m}\right)} \left((\sqrt{m} + \|\Theta\| + 1/2) \sqrt{\ln\left(1 + \frac{T}{m}\right)} + 2\sqrt{m} + \ln \frac{1}{\delta} \right). \quad (35)$$

C Adaptation of the Logistic-UCB1 Strategy of Faury et al. [2020]: Proof of Lemma 1

The proof is copied from the proof of Faury et al. [2020, Lemmas 1 and 2], with minor adjustments. We mostly provide it for the sake of self-completeness.

The adjustments are required because the setting of Faury et al. [2020] is slightly different: their action set is a subset $\mathcal{X} \subseteq \mathbb{R}^n$, and when the learner picks an action $\mathbf{x}_t \in \mathcal{X}$ at round t , the obtained reward $r_t \in \{0, 1\}$ is drawn independently at random according to a Bernoulli distribution with parameter $\eta(\mathbf{x}_t^\top \boldsymbol{\theta}_*)$, where $\boldsymbol{\theta}_* \in \mathbb{R}^n$ is unknown. The learner then only observes r_t and not what would have been obtained with a different choice of action.

Therefore, the \mathbf{x}_t and r_t of Faury et al. [2020] correspond to $\varphi(a_t, \mathbf{x}_t)$ and y_t in our setting. The main difference is that while the learner has a full control over the choice of $\mathbf{x}_t \in \mathcal{X}$ in the setting of Faury et al. [2020], in our setting, $\mathbf{x}_t \in \mathcal{X}$ is drawn by the environment and the learner only picks $a_t \in \mathcal{A}$; the learner therefore does not have a full control over $\varphi(a_t, \mathbf{x}_t)$. This is why we carefully check in the present appendix that the results by Faury et al. [2020] (namely, their Lemmas 1 and 2) extend to our setting.

Reminder — A tail inequality for self-normalized martingales. Theorem 1 of Faury et al. [2020] reads as follows. Let $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$ be a filtration, $(\mathbf{U}_t)_{t \geq 1}$ an \mathcal{F} -adapted stochastic vector process in \mathbb{R}^m such that $\|\mathbf{U}_t\| \leq 1$ a.s. for all $t \geq 1$, and $(\Delta_t)_{t \geq 1}$ an \mathcal{F} -martingale difference sequence with $|\Delta_t| \leq 1$ a.s. for all $t \geq 1$; i.e., for all $t \geq 1$, the random variable Δ_t is \mathcal{F}_{t-1} -measurable and

$$\mathbb{E}[\Delta_t | \mathcal{F}_{t-1}] = 0 \quad \text{a.s.}$$

Denote $\sigma_t^2 = \mathbb{E}[\Delta_t^2 | \mathcal{F}_{t-1}]$, let $\lambda > 0$, and define, for $t \geq 1$:

$$S_t = \sum_{s=1}^t \Delta_s \mathbf{U}_s \quad \text{and} \quad M_t = \lambda \mathbf{I}_m + \sum_{s=1}^t \sigma_s^2 \mathbf{U}_s \mathbf{U}_s^\top.$$

Then, for all $\delta \in (0, 1)$, with probability at least $1 - \delta$:

$$\forall t \geq 1, \quad \|S_t\|_{M_t^{-1}} \leq \frac{\sqrt{\lambda}}{2} + \frac{2}{\sqrt{\lambda}} \ln \left(\frac{2^m \det(M_t)^{1/2} \lambda^{-m/2}}{\delta} \right). \quad (36)$$

The result above is proved by Faury et al. [2020] based on Laplace's method on supermartingales, which is a standard argument to provide confidence bounds on self-normalized sums of conditionally centered random vectors and was previously introduced, in the context of linear contextual bandits, by Abbasi-Yadkori et al. [2011, Theorem 2]; see also the monograph by Lattimore and Szepesvári [2020, Theorem 20.2].

Step 1 — A martingale control. We apply (36) to the following elements. We take as filtration $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$, with $\mathcal{F}_0 = \sigma(a_1, \mathbf{x}_1)$ and for $t \geq 1$,

$$\mathcal{F}_t = \sigma(h_t, a_{t+1}, \mathbf{x}_{t+1}),$$

where we recall $h_t = (\mathbf{x}_s, a_s, y_s)_{s \leq t}$. We set $\mathbf{U}_t = \varphi(a_t, \mathbf{x}_t)$, which is indeed \mathcal{F}_{t-1} -measurable and with Euclidean norm smaller than 1 (thanks to the normalization assumed in Section 2.1). Finally, we set, for $t \geq 1$,

$$\Delta_t = \begin{cases} 0 & \text{if } a_t = a_{\text{null}}, \\ y_t - \eta(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*) & \text{if } a_t \neq a_{\text{null}}, \end{cases}$$

which we rewrite, by the abuses of notation indicated in Section 2.1,

$$\Delta_t = \left(y_t - \eta(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*) \right) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} = \left(y_t - P(a_t, \mathbf{x}_t) \right) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}.$$

The conditioning by \mathcal{F}_{t-1} fixes a_t and \mathbf{x}_t , but not y_t , and exactly means, when $a_t \neq a_{\text{null}}$, integrating over $y_t \sim P(a_t, \mathbf{x}_t)$. We therefore have that Δ_t is \mathcal{F}_{t-1} -measurable, with

$$\mathbb{E}[\Delta_t | \mathcal{F}_{t-1}] = 0 \quad \text{a.s.};$$

that is, $(\Delta_t)_{t \geq 1}$ appears as an \mathcal{F} -martingale difference sequence, satisfying the boundedness-by-1 constraint. We may therefore apply the result (36). To do so, we first compute the conditional variances of the Δ_t : for $t \geq 1$,

$$\begin{aligned} \mathbb{E}[\Delta_t^2 | \mathcal{F}_{t-1}] &= \mathbb{E}\left[\left(y_t - \eta(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*)\right)^2 \middle| \mathcal{F}_{t-1}\right] \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \\ &= \eta(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*) \left(1 - \eta(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*)\right) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} = \dot{\eta}(\varphi(a_t, \mathbf{x}_t)^\top \boldsymbol{\theta}_*) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}}, \end{aligned}$$

where we used the fact that $\eta(1 - \eta) = \dot{\eta}$. We rewrite

$$S_t = \sum_{s=1}^t \Delta_s \varphi(a_s, \mathbf{x}_s) = \sum_{s=1}^t \left(y_s - \eta(\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_*)\right) \varphi(a_s, \mathbf{x}_s) \mathbb{1}_{\{a_s \neq a_{\text{null}}\}}$$

and note that $M_t = W_t(\boldsymbol{\theta})$ with the definition (5). The control (36) may be rewritten as follows: for all $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\forall t \geq 1, \quad \|S_t\|_{W_t(\boldsymbol{\theta}_*)^{-1}} \leq \frac{\sqrt{\lambda}}{2} + \frac{2}{\sqrt{\lambda}} \ln\left(\frac{2^m \det(W_t(\boldsymbol{\theta}_*))^{1/2} \lambda^{-m/2}}{\delta}\right).$$

As $\dot{\eta} = \eta(1 - \eta) \in [0, 1/4]$ and as $\|\varphi\| \leq 1$, we have, by a standard trace-determinant inequality (see, e.g., Abbasi-Yadkori et al. [2011, Lemma 10]),

$$\begin{aligned} \det(W_t(\boldsymbol{\theta}_*)) &\leq \left(\lambda + \frac{1}{m} \sum_{s=1}^t \dot{\eta}(\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_*) \|\varphi(a_s, \mathbf{x}_s)\|^2 \mathbb{1}_{\{a_s \neq a_{\text{null}}\}}\right)^m \\ &\leq \left(\lambda + \frac{1}{4m} \sum_{s=1}^t \mathbb{1}_{\{a_s \neq a_{\text{null}}\}}\right)^m \leq \left(\lambda + \frac{t}{4m}\right)^m. \end{aligned}$$

Combining the two inequalities, we have proved that for all $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\forall t \geq 1, \quad \|S_t\|_{W_t(\boldsymbol{\theta}_*)^{-1}} \leq \gamma_{t,\lambda,\delta} - \sqrt{\lambda} \|\Theta\|, \quad (37)$$

where $\gamma_{t,\lambda,\delta}$ was defined in Lemma 1.

Step 2 — Application of the martingale control. The martingale control (37) is applied as follows. We show below that the definition (3) of $\tilde{\boldsymbol{\theta}}_t$ entails that

$$S_t - \lambda \boldsymbol{\theta}_* = \Psi_t(\tilde{\boldsymbol{\theta}}_t) - \Psi_t(\boldsymbol{\theta}_*), \quad (38)$$

where Ψ_t was defined in (4). Taking the $\|\cdot\|_{W_t(\boldsymbol{\theta}_*)^{-1}}$ norms of both sides, together with a triangle inequality (keeping in mind the boundedness of Θ) and noting that

$$\|\boldsymbol{\theta}\|_{W_t(\boldsymbol{\theta}_*)^{-1}} \leq \|\boldsymbol{\theta}\|_{(\lambda \mathbf{I}_m)^{-1}} = \|\boldsymbol{\theta}\| / \sqrt{\lambda},$$

finally yields that for all $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\forall t \geq 1, \quad \left\| \Psi_t(\tilde{\boldsymbol{\theta}}_t) - \Psi_t(\boldsymbol{\theta}_*) \right\|_{W_t(\boldsymbol{\theta}_*)^{-1}} \leq \|S_t\|_{W_t(\boldsymbol{\theta}_*)^{-1}} + \sqrt{\lambda} \|\boldsymbol{\theta}\| \leq \gamma_{t,\lambda,\delta}. \quad (39)$$

We now show (38). The gradient of the continuously differentiable function

$$\boldsymbol{\theta} \in \mathbb{R}^m \mapsto \sum_{s=1}^t \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \left(y_s \ln \eta(\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}) + (1 - y_s) \ln \left(1 - \eta(\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta})\right) \right) - \frac{\lambda}{2} \|\boldsymbol{\theta}\|$$

vanishes at the point $\tilde{\boldsymbol{\theta}}_t$ where it achieves its maximum, i.e., $\tilde{\boldsymbol{\theta}}_t$ defined in (3) satisfies

$$\begin{aligned} &\sum_{s=1}^t \left(\eta(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t) - y_s \right) \varphi(a_s, \mathbf{x}_s) \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \\ &= \sum_{s=1}^t \left(y_s \frac{\dot{\eta}(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t) \varphi(a_s, \mathbf{x}_s)}{\eta(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t)} + (1 - y_s) \frac{\dot{\eta}(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t) \varphi(a_s, \mathbf{x}_s)}{1 - \eta(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t)} \right) \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} = \lambda \tilde{\boldsymbol{\theta}}_t, \end{aligned}$$

where we used $\dot{\eta} = \eta(1 - \eta)$ to get the first equality. In particular,

$$\Psi_t(\tilde{\boldsymbol{\theta}}_t) = \sum_{s=1}^t \eta(\varphi(a_s, \mathbf{x}_s)^\top \tilde{\boldsymbol{\theta}}_t) \varphi(a_s, \mathbf{x}_s) \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} + \lambda \tilde{\boldsymbol{\theta}}_t = \sum_{s=1}^t y_s \varphi(a_s, \mathbf{x}_s) \mathbb{1}_{\{a_s \neq a_{\text{null}}\}},$$

hence the stated rewriting (38).

Step 3 — Bound on prediction error. We now proceed to bounding, for all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$ and \mathbf{x} :

$$|P(a, \mathbf{x}) - \widehat{P}_t(a, \mathbf{x})| = \left| \eta(\varphi(a, \mathbf{x})^\top \boldsymbol{\theta}_*) - \eta(\varphi(a, \mathbf{x})^\top \widehat{\boldsymbol{\theta}}_t) \right|.$$

As η is $1/4$ -Lipschitz (given $\dot{\eta} = \eta(1 - \eta) \in [0, 1/4]$),

$$|\widehat{P}_t(a, \mathbf{x}) - P(a, \mathbf{x})| \leq \frac{1}{4} \left| \varphi(a, \mathbf{x})^\top (\boldsymbol{\theta}_* - \widehat{\boldsymbol{\theta}}_t) \right|. \quad (40)$$

For two $m \times m$ symmetric definite positive matrices M and M' , we write $M \succeq M'$ whenever $\|\mathbf{v}\|_M \geq \|\mathbf{v}\|_{M'}$ for all $\mathbf{v} \in \mathbb{R}^m$. This inequality entails $(M')^{-1} \succeq M^{-1}$. We introduce below a symmetric definite positive matrix G_t such that

$$G_t \succeq \kappa^{-1} V_t, \quad G_t \succeq (1 + 2\|\Theta\|)^{-1} W_t(\widehat{\boldsymbol{\theta}}_t), \quad G_t \succeq (1 + 2\|\Theta\|)^{-1} W_t(\boldsymbol{\theta}_*) \quad (41)$$

and

$$\Psi_t(\widehat{\boldsymbol{\theta}}_t) - \Psi_t(\boldsymbol{\theta}_*) = G_t(\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_*). \quad (42)$$

Based on all these properties, we get, by a Cauchy-Schwarz inequality in the norms $\|\cdot\|_{G_t}$ and $\|\cdot\|_{G_t^{-1}}$:

$$\begin{aligned} & \left| \varphi(a, \mathbf{x})^\top (\boldsymbol{\theta}_* - \widehat{\boldsymbol{\theta}}_t) \right| \\ & \leq \|\varphi(a, \mathbf{x})\|_{G_t^{-1}} \|\boldsymbol{\theta}_* - \widehat{\boldsymbol{\theta}}_t\|_{G_t} \\ & = \|\varphi(a, \mathbf{x})\|_{G_t^{-1}} \left\| G_t(\boldsymbol{\theta}_* - \widehat{\boldsymbol{\theta}}_t) \right\|_{G_t^{-1}} \\ & = \|\varphi(a, \mathbf{x})\|_{G_t^{-1}} \left\| \Psi(\widehat{\boldsymbol{\theta}}_t) - \Psi(\boldsymbol{\theta}_*) \right\|_{G_t^{-1}} \leq \sqrt{\kappa} \|\varphi(a, \mathbf{x})\|_{V_t^{-1}} \left\| \Psi(\widehat{\boldsymbol{\theta}}_t) - \Psi(\boldsymbol{\theta}_*) \right\|_{G_t^{-1}}. \end{aligned}$$

A triangle inequality for the first inequality, followed by applying (41) for the second inequality, and applying the definition of (4) as a projection for the third inequality, shows that

$$\begin{aligned} \left\| \Psi(\widehat{\boldsymbol{\theta}}_t) - \Psi(\boldsymbol{\theta}_*) \right\|_{G_t^{-1}} & \leq \left\| \Psi(\widehat{\boldsymbol{\theta}}_t) - \Psi(\tilde{\boldsymbol{\theta}}_t) \right\|_{G_t^{-1}} + \left\| \Psi(\boldsymbol{\theta}_*) - \Psi(\tilde{\boldsymbol{\theta}}_t) \right\|_{G_t^{-1}} \\ & \leq \sqrt{1 + 2\|\Theta\|} \left(\left\| \Psi(\widehat{\boldsymbol{\theta}}_t) - \Psi(\tilde{\boldsymbol{\theta}}_t) \right\|_{W_t(\widehat{\boldsymbol{\theta}}_t)^{-1}} + \left\| \Psi(\boldsymbol{\theta}_*) - \Psi(\tilde{\boldsymbol{\theta}}_t) \right\|_{W_t(\boldsymbol{\theta}_*)^{-1}} \right) \\ & \leq 2\sqrt{1 + 2\|\Theta\|} \left\| \Psi(\boldsymbol{\theta}_*) - \Psi(\tilde{\boldsymbol{\theta}}_t) \right\|_{W_t(\boldsymbol{\theta}_*)^{-1}}. \end{aligned}$$

Substituting the martingale control (39) and collecting all bounds together finally yields

$$|\widehat{P}_t(a, \mathbf{x}) - P(a, \mathbf{x})| \leq \frac{\sqrt{\kappa}}{2} \|\varphi(a, \mathbf{x})\|_{V_t^{-1}} \sqrt{1 + 2\|\Theta\|} \gamma_{t, \lambda, \delta},$$

as desired.

Note in particular that (39) holds for all $t \geq 1$ with probability $1 - \delta$, and that we took care of the dependencies on a and \mathbf{x} through the $\|\varphi(a, \mathbf{x})\|_{V_t^{-1}}$ term. This explains why the result of Lemma 1 holds with probability $1 - \delta$ for all $t \geq 1$, all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$ and all $\mathbf{x} \in \mathcal{X}$.

Step 4 — Construction of the matrix G_t . It only remains to show that there exists a symmetric definite positive matrix G_t such that (41) and (42) hold. We define

$$\begin{aligned} G_t & = \int_{[0,1]} W_t(v \widehat{\boldsymbol{\theta}}_t + (1-v)\boldsymbol{\theta}_*) \, dv \\ & = \lambda \mathbf{I}_m + \sum_{s=1}^t \left(\int_{[0,1]} \dot{\eta}(v \varphi(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t + (1-v)\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_*) \, dv \right) \varphi(a_s, \mathbf{x}_s) \varphi(a_s, \mathbf{x}_s)^\top \mathbb{1}_{\{a_s \neq a_{\text{null}}\}}. \end{aligned}$$

The thus defined matrix G_t is indeed symmetric definite positive matrix. By definition of κ and the fact that Θ is convex, we have, for all $v \in [0, 1]$,

$$\dot{\eta}(v \varphi(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t + (1-v)\varphi(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_*) = \dot{\eta}(\varphi(a_s, \mathbf{x}_s)^\top (v \widehat{\boldsymbol{\theta}}_t + (1-v)\boldsymbol{\theta}_*)) \geq \kappa^{-1},$$

which also immediately entails the first inequality of (41). To prove (42), we introduce, for $s \geq 1$ such that $a_s \neq a_{\text{null}}$, the continuously differentiable function

$$f_{s,t} : v \in [0, 1] \mapsto f_{s,t}(v) = \eta(v \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t + (1-v) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star),$$

with derivative

$$\dot{f}_{s,t}(v) = \dot{\eta}(v \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t + (1-v) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top (\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_\star),$$

and we have

$$\eta(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t) - \eta(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) = f_{s,t}(1) - f_{s,t}(0) = \int_{[0,1]} \dot{f}_{s,t}(v) dv,$$

These facts, combined with the definition of G_t , immediately entail (42).

It only remains to prove the third inequality of (41), namely,

$$G_t \succeq (1 + 2\|\Theta\|)^{-1} W_t(\boldsymbol{\theta}_\star),$$

as the second one is obtained by symmetry from there, by exchanging the roles of $\boldsymbol{\theta}_\star$ and $\widehat{\boldsymbol{\theta}}_t$. To do so, we rely on the following inequality: for all $z_1, z_2 \in \mathbb{R}$,

$$\int_{[0,1]} \dot{\eta}(z_1 + v(z_2 - z_1)) dv \geq \frac{\dot{\eta}(z_1)}{1 + |z_1 - z_2|}. \quad (43)$$

This inequality is proved, in the case $z_1 \neq z_2$, by noting that the second-order derivative of η equals

$$\ddot{\eta}(x) = \frac{e^{-x} - 1}{1 + e^x} \dot{\eta}(x), \quad \text{where} \quad \frac{e^{-x} - 1}{1 + e^x} \in [-1, 1],$$

so that for all $z, z' \in \mathbb{R}$,

$$\ln \dot{\eta}(z') - \ln \dot{\eta}(z) = \int_z^{z'} \frac{\ddot{\eta}(\tau)}{\dot{\eta}(\tau)} d\tau \geq -|z' - z|, \quad \text{i.e.,} \quad \dot{\eta}(z') \geq \dot{\eta}(z) e^{-|z' - z|}.$$

Therefore,

$$\int_{[0,1]} \dot{\eta}(z_1 + v(z_2 - z_1)) dv \geq \int_{[0,1]} \dot{\eta}(z_1) e^{-v|z_1 - z_2|} dv = \dot{\eta}(z_1) \frac{1 - e^{-|z_1 - z_2|}}{|z_1 - z_2|} \geq \frac{\dot{\eta}(z_1)}{1 + |z_1 - z_2|},$$

where we applied the inequality $(1 - e^{-x})/x \geq 1/(1+x)$, which holds for all $x \geq 0$.

We go back to proving the third inequality of (41). With (43) for the first inequality, followed by an application of the Cauchy-Schwarz inequality for the second inequality, and the fact that $\boldsymbol{\theta}_\star, \widehat{\boldsymbol{\theta}}_t \in \Theta$ have Euclidean norms smaller than $\|\Theta\|$, together with $\|\boldsymbol{\varphi}\| \leq 1$, we have, for all $s \geq 1$ with $a_s \neq a_{\text{null}}$,

$$\begin{aligned} & \int_{[0,1]} \dot{\eta}(v \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \widehat{\boldsymbol{\theta}}_t + (1-v) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) dv \\ & \geq \dot{\eta}(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) \left(1 + |\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top (\boldsymbol{\theta}_\star - \widehat{\boldsymbol{\theta}}_t)|\right)^{-1} \\ & \geq \dot{\eta}(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) \left(1 + \|\boldsymbol{\varphi}(a_s, \mathbf{x}_s)\| \|\boldsymbol{\theta}_\star - \widehat{\boldsymbol{\theta}}_t\|\right)^{-1} \\ & \geq \dot{\eta}(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) (1 + 2\|\Theta\|)^{-1}. \end{aligned}$$

As $(1 + 2\|\Theta\|)^{-1} \leq 1$, we can then conclude, from the definition of G_t , that

$$\begin{aligned} G_t & \succeq (1 + 2\|\Theta\|)^{-1} \left(\lambda I_m + \sum_{s=1}^t \dot{\eta}(\boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \boldsymbol{\theta}_\star) \boldsymbol{\varphi}(a_s, \mathbf{x}_s) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \right) \\ & = (1 + 2\|\Theta\|)^{-1} W_t(\boldsymbol{\theta}_\star), \end{aligned}$$

as announced. This concludes the proof.

D Proof of the Regret Bound in Case of an Unknown Distribution ν : Proof of Theorem 2

We rather explain the differences to and modifications with respect to the proof of Appendix B.

To best do so, we use $\hat{\cdot}$ superscripts to index various quantities defined based on the estimations $\hat{\nu}_t$, even though these quantities are not themselves estimators. In particular, the budget parameter is denoted by \hat{B}_T ; we refer to the policies computed at rounds $t \geq 2$ by $\hat{\mathbf{p}}_t(h_{t-1}, \cdot)$, and by $\hat{\mathcal{E}}_\delta$ the event of Lemma 1 for the sampling strategy pulling actions $a_t \sim \hat{\mathbf{p}}_t(\mathbf{x}_t)$, which is exactly the strategy that we are analyzing here; and finally, the optimal dual variables linked to the budget are denoted by $\hat{\beta}_t^{\text{budg}, \star}$.

D.1 Key New Building Block: Uniform Deviation Inequality

Throughout this appendix, we will need to relate quantities defined based on $\hat{\nu}_t$ to the target quantities defined based on ν . All these quantities will be of the form: for $1 \leq t \leq T$ and for various functions $f : \mathcal{X} \rightarrow [0, 1]$,

$$\mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} [f(\mathbf{X})] \quad \text{and} \quad \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} [f(\mathbf{X})].$$

A simple (but probably slightly suboptimal) way to do so is to apply $T|\mathcal{X}|$ times the Hoeffding-Azuma inequality together with a union bound. We get that on an event $\hat{\mathcal{E}}_{\text{unif}, \delta}$ of probability at least $1 - \delta$,

$$\forall 1 \leq t \leq T, \quad \forall \mathbf{x} \in \mathcal{X}, \quad |\hat{\nu}_t(\mathbf{x}) - \nu(\mathbf{x})| \leq \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}.$$

In particular, with probability at least $1 - \delta$, for all functions $f : \mathcal{X} \rightarrow [0, 1]$ and all $1 \leq t \leq T$,

$$\left| \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} [f(\mathbf{X})] - \mathbb{E}_{\mathbf{X} \sim \nu} [f(\mathbf{X})] \right| \leq \sum_{\mathbf{x} \in \mathcal{X}} |\hat{\nu}_t(\mathbf{x}) - \nu(\mathbf{x})| \leq |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}. \quad (44)$$

We now explain the adaptations required (or not required) for each of the four steps of the proof provided in Appendix B.

D.2 First and Second Steps: No Adaptation Required

These two steps do not require any adaptation; we merely re-state the useful results extracted therein, with the corresponding adapted notation.

The first step (Appendix B.1) held for any sampling strategy. Therefore, the same upper-confidence bonuses (9) and Lemma 1 entail that on an event $\hat{\mathcal{E}}_\delta$ of probability at least $1 - \delta$, for all $t \geq 1$, all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$, and all $\mathbf{x} \in \mathcal{X}$:

$$P(a, \mathbf{x}) \leq U_t(a, \mathbf{x}) \leq P(a, \mathbf{x}) + 2\varepsilon_t(a, \mathbf{x}). \quad (45)$$

The bound (12) also still holds, as it was obtained in a deterministic manner not using any specific feature of the sampling strategy; namely,

$$2 \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \leq E_T. \quad (46)$$

Similarly, the second step (Appendix B.2) actually yielded general results between the primal and the dual formulations of the OPT problems considered. The equality (18), the characterizations (22) and (23), as well as the inequality (25) may be instantiated with $\hat{\nu}_t$ (in lieu of ν) and \hat{B}_T (in lieu of B_T) as follows. For each $t \geq 2$ such that the cost constraints of Phase 0 of the adaptive policy are not violated and the optimization problem $\text{OPT}(\hat{\nu}_t, U_{t-1}, \hat{B}_T)$ is to be solved, there exists a vector $\hat{\beta}_t^{\text{budg}, \star} \geq \mathbf{0}$ such that first,

$$\left(\hat{\beta}_t^{\text{budg}, \star} \right)^\top \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \hat{\mathbf{p}}_{t,a}(\mathbf{X}) \right] = \frac{\hat{B}_T}{T} \left(\hat{\beta}_t^{\text{budg}, \star} \right)^\top \mathbf{1}. \quad (47)$$

By (44) and the fact that the sums at stake below lie in $[0, 1]$, we have, on $\widehat{\mathcal{E}}_{\text{unif},\delta}$, that for all $t \geq 2$,

$$\begin{aligned} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ \leq \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} c(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] + |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}} \mathbf{1}. \end{aligned}$$

Combining the two inequalities above with (53) and (55), as well as with the bounds $P \leq U_{t-1}$ of (45), we proved so far that

$$\text{on } \widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}, \quad \sum_{t=1}^T c(a_t, \mathbf{x}_t) y_t \leq \left(\widehat{B}_T + 1 + \sqrt{2T \ln \frac{4d}{\delta}} + \sum_{t=2}^T |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}} \right) \mathbf{1}.$$

Since $\sum_{t \leq T} 1/\sqrt{t} \leq 2\sqrt{T}$ and by the definition of \widehat{B}_T (see the statement of Theorem 2), we proved that

$$\text{on } \widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}, \quad \sum_{t=1}^T c(a_t, \mathbf{x}_t) y_t \leq (B-1)\mathbf{1},$$

as claimed. Therefore, on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$, the adaptive policy of Box B of Section 3 never stays in Phase 0 and instead solves, at each round $t \geq 2$, the Phase-2 problem $\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)$.

Dealing with the rewards. By (52) and (54), by the bound $U_{t-1} \leq P + 2\varepsilon_{t-1}$ of (45) together with the bound E_T of (46), and by the uniform control (44), we have similarly that on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$,

$$\begin{aligned} \sum_{t=1}^T r(a_t, \mathbf{x}_t) y_t \geq \sum_{t=2}^T \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ - \left(E_T + \sqrt{2T \ln \frac{4d}{\delta}} + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} \right). \end{aligned}$$

Given that on the intersection of events considered, the adaptive policy solves $\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)$ for all $t \geq 2$, we have

$$\mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] = \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T}. \quad (56)$$

This concludes the proof of (51), and hence, the adaptation of the third step.

D.4 Fourth Step: Some Adaptations are Also Required

In this final step, we collect the bounds from the previous three steps. Some (rather minor) adaptations are required, e.g., we would like to integrate (48) and (49) over $\widehat{\nu}_t$ in the left-hand sides and ν in the right-hand sides.

Main modification. We consider some $t \geq 2$. For each $\mathbf{x} \in \mathcal{X}$, we apply (48) and (49) with $\mathbf{q} = \pi^*(\mathbf{x})$ and get

$$\begin{aligned} \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\widehat{\beta}_t^{\text{budg},*})^\top c(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{x}) \\ \geq \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\widehat{\beta}_t^{\text{budg},*})^\top c(a, \mathbf{x}) \right)_+ U_{t-1}(a, \mathbf{x}) \pi_a^*(\mathbf{x}). \end{aligned}$$

The bound $U_{t-1} \geq P$ of (45) and the non-negative parts taken in the right-hand side then entail that

$$\begin{aligned} \text{on } \widehat{\mathcal{E}}_{\text{prob},\delta}, \quad \forall \mathbf{x} \in \mathcal{X}, \quad & \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right) U_{t-1}(a, \mathbf{x}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{x}) \\ & \geq \sum_{a \in \mathcal{A}} \left(r(a, \mathbf{x}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{x}) \right)_+ P(a, \mathbf{x}) \pi_a^*(\mathbf{x}). \end{aligned}$$

We replace the individual \mathbf{x} by a random variable $\mathbf{X} \sim \widehat{\nu}_t$ and take expectations with respect to \mathbf{X} :

$$\begin{aligned} \text{on } \widehat{\mathcal{E}}_{\text{prob},\delta}, \quad & \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ & \geq \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right)_+ P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]. \end{aligned}$$

Thanks to the non-negative parts in the right-hand side, we identify some function $f(\mathbf{X})$ where f takes values in $[0, 1]$, so that we may apply the uniform control (44) and get

$$\begin{aligned} \text{on } \widehat{\mathcal{E}}_{\text{unif},\delta}, \quad & \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right)_+ P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right] \\ & \geq \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right)_+ P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right] - |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}} \\ & \geq \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \left(r(a, \mathbf{X}) - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{c}(a, \mathbf{X}) \right) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right] - |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}, \end{aligned}$$

where for the second inequality, we simply dropped the non-negative parts.

The rest of the proof for this final step is basically unchanged. Combining all the bounds exhibited so far in this updated fourth step, we have, for each $t \geq 2$,

$$\begin{aligned} \text{on } \widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}, \quad & \overbrace{\mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right]}^{=\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)/T} \\ & - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \\ & \geq \overbrace{\mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}^{=\text{OPT}(\nu, P, B)/T} \\ & - \underbrace{(\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) P(a, \mathbf{X}) \pi_a^*(\mathbf{X}) \right]}_{\leq (B/T)\mathbf{1}} - |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}, \end{aligned}$$

where we substituted the inequalities stemming from the definition of π^* as well as the rewriting (56).

Rearranging the inequality above and substituting (47), we get that on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$,

$$\begin{aligned} & \frac{\text{OPT}(\nu, P, B)}{T} - \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} \\ & \leq \frac{B(\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{1}}{T} - (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(a, \mathbf{X}) U_{t-1}(a, \mathbf{X}) \widehat{p}_{t,a}(\mathbf{X}) \right] + |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}} \\ & = \frac{B - \widehat{B}_T}{T} (\widehat{\boldsymbol{\beta}}_t^{\text{budg},*})^\top \mathbf{1} + |\mathcal{X}| \sqrt{\frac{1}{2t} \ln \frac{2T|\mathcal{X}|}{\delta}}. \end{aligned} \tag{57}$$

Summing this bound over $2 \leq t \leq T$ and combining it with (50), we obtain that on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$,

$$\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} \leq 1 + \frac{B - \widehat{B}_T}{\widehat{B}_T} \sum_{t=2}^T \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}}.$$

By distinguishing the cases

$$\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} \leq 0 \quad \text{and} \quad \text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} \geq 0,$$

the inequality above entails that on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$,

$$\text{OPT}(\nu, P, B) - \sum_{t=2}^T \frac{\text{OPT}(\widehat{\nu}_t, U_{t-1}, \widehat{B}_T)}{T} \leq \frac{B - \widehat{B}_T}{\widehat{B}_T} \text{OPT}(\nu, P, B) + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} + 1.$$

Substituting this upper bound into (51), we finally obtain that on $\widehat{\mathcal{E}}_{\text{prob},\delta} \cap \widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$,

$$\begin{aligned} \text{OPT}(\nu, P, B) - \sum_{t=1}^T r(a_t, \mathbf{x}_t) y_t &\leq \frac{B - \widehat{B}_T}{\widehat{B}_T} \text{OPT}(\nu, P, B) \\ &\quad + E_T + \underbrace{\sqrt{2T \ln \frac{4d}{\delta}} + 2|\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}}}_{\leq 2b_T} + 1. \end{aligned}$$

We conclude the proof by the same modifications to improve readability as at the end of the proof of Theorem 1: namely, since the definition of E_T did not change, the bound (35) is still applicable, while

$$\frac{B - \widehat{B}_T}{\widehat{B}_T} \leq \frac{2b_T}{B} = \frac{1}{B} \left(4 + 2\sqrt{2T \ln \frac{4d}{\delta}} + 2|\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} \right)$$

is obtained with the same techniques and similar conditions as for (32).

E Details for the Material of Section 6

In this section, we first recall (Appendix E.1) the setting of linear CBwK introduced by Agrawal and Devanur [2016]—and actually, slightly generalize it to match the setting of CBwK for a logistic-regression conversion model. We then state the adaptive policy considered (Appendix E.2), which relies on upper-confidence estimates of the rewards and lower-confidence estimates of the cost vectors. We also state the corresponding estimation guarantees in a key lemma (Lemma 3). The heart of this section is to state, discuss (Appendix E.3) and prove (Appendix E.4) a regret bound, matching the one of Agrawal and Devanur [2016, Theorem 3], with a slight improvement consisting of a relaxation of the budget constraints. For the sake of completeness, we finally recall (Appendix E.5) the statement of the adaptive policy of Agrawal and Devanur [2016].

E.1 Setting

The setting is the following. We consider a finite action set \mathcal{A} including a no-op action a_{null} , a finite context set $\mathcal{X} \subseteq \mathbb{R}^n$, a number T of rounds and a total budget constraint $B > 0$. All these parameters are known. Contexts $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ are drawn i.i.d. according to some distribution ν . At round $t \geq 1$, the learner observes the context \mathbf{x}_t , picks an action a_t and, conditionally to \mathbf{x}_t and a_t , when $a_t \neq a_{\text{null}}$, obtains a reward $r_t \in [0, 1]$ drawn independently at random according to a distribution with expectation $\bar{r}(a_t, \mathbf{x}_t)$, where

$$\forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \bar{r}(a, \mathbf{x}) = \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\mu}_*,$$

and suffers a vector cost $\mathbf{c}_t \in [0, 1]^d$ drawn independently at random according to a distribution with vector of expectations $\bar{\mathbf{c}}(a_t, \mathbf{x}_t)$, where each component \bar{c}_i of $\bar{\mathbf{c}}$, for $i \in \{1, \dots, d\}$, is given by

$$\forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \bar{c}_i(a, \mathbf{x}) = \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta}_{*,i}.$$

In the definitions above, $\boldsymbol{\varphi} : \mathcal{A} \setminus \{a_{\text{null}}\} \times \mathcal{X} \rightarrow \mathbb{R}^m$ is a known transfer function, with $\|\boldsymbol{\varphi}\| \leq 1$, and $\boldsymbol{\mu}_*$ and the $\boldsymbol{\theta}_{*,i}$ are unknown parameters in \mathbb{R}^m . We assume that these unknown parameters lie in some bounded set Θ , with maximal norm still denoted by $\|\Theta\|$. When $a_t = a_{\text{null}}$, the obtained reward and suffered costs are null: $r_t = 0$ and $\mathbf{c}_t = \mathbf{0}$.

Comparison to the canonical setting of linear CBwK. Note that in the original formulation of Agrawal and Devanur [2016], we have (where \mathbf{x}_a also denote vectors):

$$\mathbf{x} = (\mathbf{x}_a)_{a \in \mathcal{A} \setminus \{a_{\text{null}}\}} \quad \text{and} \quad \boldsymbol{\varphi}(a, \mathbf{x}) = \mathbf{x}_a.$$

Benchmark and regret. The goal is still to maximize the accumulated rewards while controlling the costs:

$$\text{maximize} \quad \sum_{t \leq T} r_t \quad \text{while controlling} \quad \sum_{t \leq T} \mathbf{c}_t \leq B \mathbf{1}.$$

The goal can be equivalently defined as the minimization of the regret while controlling the costs, where the regret equals

$$R_T = \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) - \sum_{t \leq T} r_t$$

for the benchmark given by the static policy π^* achieving the largest expected cumulative rewards under the condition that its cumulative vector costs abide by the budget constraints in expectation, i.e.,

$$\begin{aligned} \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) &= \max_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \bar{r}(a, \mathbf{X}) \right] \\ &\text{under} \quad T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \bar{\mathbf{c}}(a, \mathbf{x}) \right] \leq B \mathbf{1}. \end{aligned} \tag{58}$$

In the sequel, we will use the definition (58) of OPT with different quadruplets of parameters; see, for instance, the definition of the adaptive policy of Box C.

E.2 Statement of an Adaptive Policy

The considered adaptive policy is stated in Box C. It is adapted from the adaptive policy of Section 3. The (almost only) changes lie in Phase 1, which depends heavily of the model. Here, we resort (as Agrawal and Devanur [2016]) to a LinUCB-type estimation of the parameters of the stochastic linear bandits yielding rewards and costs. Based on these estimated parameters, we are able to issue, at each round $t \geq 2$, an upper-confidence expected reward function U_{t-1} and a lower-confidence expected vector-cost function L_{t-1} . We also use the empirical estimate of the context distribution. In Phase 2, we solve the OPT problem on these estimates and with the conservative budget \hat{B}_T .

The adaptive policy of Box C bears links, and actually generalizes, the one by Xu and Truong [2019]. The setting of the latter reference is more limited as more information is provided to the learner, such as the costs for taking each action and the distribution ν of contexts (clients in their case).

For the sake of completeness, we state in Appendix E.5 the adaptive policy introduced by Agrawal and Devanur [2016].

The main additional ingredient in the analysis of the policy of Box C, compared to the analyses of the adaptive policy of Section 3, is a guarantee on the outcomes of Phase 1. We recall that we assumed that the parameters $\boldsymbol{\mu}_*$ and $\boldsymbol{\theta}_{*,i}$ lie in a bounded set Θ with maximal norm denoted by $\|\Theta\|$.

Lemma 3 (direct adaptation from Abbasi-Yadkori et al. [2011, Theorem 2]). *Fix any sampling strategy and consider the version of LinUCB given by Phase 1 of Box C. For all $\delta \in (0, 1)$, there exists an event $\mathcal{E}_{\text{lin},\delta}$ with probability at least $1 - \delta$ and such that over $\mathcal{E}_{\text{lin},\delta}$:*

$$\forall t \geq 1, \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad |\hat{r}_t(a, \mathbf{x}) - \bar{r}(a, \mathbf{x})| \leq \gamma_{t,\lambda,\delta} \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{X_t^{-1}}$$

$$\text{and} \quad |\hat{c}_t(a, \mathbf{x}) - \bar{c}(a, \mathbf{x})| \leq \gamma_{t,\lambda,\delta} \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{X_t^{-1}} \mathbf{1},$$

where

$$\gamma_{t,\lambda,\delta} = \frac{1}{4} \sqrt{m \ln \frac{1+t/(\lambda m)}{\delta/(d+1)}} + \sqrt{\lambda} \|\Theta\|.$$

Proof sketch. We explain why the bound for \bar{r} holds with probability at least $1 - \delta/(d+1)$. The lemma follows by repeating the argument for the components of \bar{c} and resorting to a union bound.

Given that rewards lie in $[0, 1]$ and are thus $1/4$ -sub-Gaussian, the martingale analysis by Abbasi-Yadkori et al. [2011, Theorem 2], with the same adaptations as the ones carried out in Appendix C to take into account the rounds when $a_t = a_{\text{null}}$, shows that with probability at least $1 - \delta/(d+1)$,

$$\forall t \geq 1, \quad \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_*\|_{X_t} \leq \frac{1}{4} \sqrt{m \ln \frac{1+t/(\lambda m)}{\delta/(d+1)}} + \sqrt{\lambda} \|\Theta\| = \gamma_{t,\lambda,\delta}.$$

We then proceed similarly to the Cauchy-Schwarz inequalities following (42): for all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$ and $\mathbf{x} \in \mathcal{X}$,

$$|\hat{r}_t(a, \mathbf{x}) - \bar{r}(a, \mathbf{x})| = \left| \boldsymbol{\varphi}(a, \mathbf{x})^\top (\boldsymbol{\mu}_{t-1} - \boldsymbol{\mu}_*) \right| \leq \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{X_t^{-1}} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_*\|_{X_t}.$$

This concludes the proof. \square

As a consequence, we set, when defining the adaptive policy of Box C,

$$\varepsilon_s(a, \mathbf{x}) = \gamma_{t,\lambda,\delta} \|\boldsymbol{\varphi}(a, \mathbf{x})\|_{X_s^{-1}}$$

for all $s \geq 1$, and denote by $\hat{\mathcal{E}}_{\text{lin},\delta}$ the event of Lemma 3 for the sampling policy of Box C. This event is of probability at least $1 - \delta$. We have:

$$\text{on } \hat{\mathcal{E}}_{\text{lin},\delta}, \quad \forall t \geq 1, \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X},$$

$$\bar{r}(a, \mathbf{x}) \leq U_t(a, \mathbf{x}) \leq \bar{r}(a, \mathbf{x}) + 2\varepsilon_t(a, \mathbf{x}) \quad (59)$$

$$\text{and} \quad L_t(a, \mathbf{x}) \leq \bar{c}(a, \mathbf{x}) \leq L_t(a, \mathbf{x}) + 2\varepsilon_t(a, \mathbf{x}) \mathbf{1}. \quad (60)$$

BOX C: LINUCB FOR DIRECT SOLUTIONS TO OPT PROBLEMS

Parameters: regularization parameter $\lambda > 0$; conservative-budget parameter \widehat{B}_T ; upper-confidence bonuses $\varepsilon_s(a, \mathbf{x}) > 0$, for $s \geq 1$ and $(a, \mathbf{x}) \in (\mathcal{A} \setminus \{a_{\text{null}}\}) \times \mathcal{X}$.

Round $t = 1$: play an arbitrary action $a_1 \in \mathcal{A} \setminus \{a_{\text{null}}\}$

At rounds $t \geq 2$:

Phase 0 If $\sum_{s \leq t-1} \mathbf{c}_s \leq (B-1)\mathbf{1}$ is violated, then $\widehat{\mathbf{p}}_t(h_{t-1}, \mathbf{x}) = \delta_{a_{\text{null}}}$ for all \mathbf{x}

Phase 1 Otherwise, estimate the parameters by

$$\begin{aligned} \boldsymbol{\mu}_{t-1} &= X_{t-1}^{-1} \sum_{s=1}^{t-1} \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \boldsymbol{\varphi}(a_s, \mathbf{x}_s) r_s \\ \text{and } \widehat{\boldsymbol{\theta}}_{t-1,i} &= X_{t-1}^{-1} \sum_{s=1}^{t-1} \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \boldsymbol{\varphi}(a_s, \mathbf{x}_s) c_{s,i} \\ \text{where } X_t &= \sum_{s=1}^t \mathbb{1}_{\{a_s \neq a_{\text{null}}\}} \boldsymbol{\varphi}(a_s, \mathbf{x}_s) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top + \lambda \mathbf{I}_m \end{aligned}$$

Define the expected reward function \widehat{r} and cost function $\widehat{\mathbf{c}}_{t-1} = (\widehat{c}_{t-1,i})_{1 \leq i \leq d}$ as

$$\begin{aligned} \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \widehat{r}_{t-1}(a, \mathbf{x}) &= \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\mu}_{t-1} \\ \text{and } \forall 1 \leq i \leq d, \quad \widehat{c}_{t-1,i} &= \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta}_{t-1,i} \end{aligned}$$

Build the upper-confidence expected reward function U_{t-1} and the lower-confidence expected vector-cost function \mathbf{L}_{t-1} as

$$\begin{aligned} \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \\ U_{t-1}(a, \mathbf{x}) &= \max \left\{ \min \{ \widehat{r}_{t-1}(a, \mathbf{x}) + \varepsilon_{t-1}(a, \mathbf{x}), 1 \}, 0 \right\} \\ \mathbf{L}_{t-1}(a, \mathbf{x}) &= \max \left\{ \min \{ \widehat{\mathbf{c}}_{t-1} - \varepsilon_{t-1}(a, \mathbf{x}) \mathbf{1}, \mathbf{1} \}, \mathbf{0} \right\} \end{aligned}$$

where the maximum and minimum are taken pointwise in the definition of \mathbf{L}_{t-1}

Set $U_{t-1}(a_{\text{null}}, \mathbf{x}) = 0$ and $\mathbf{L}_{t-1}(a_{\text{null}}, \mathbf{x}) = \mathbf{0}$ for all $\mathbf{x} \in \mathcal{X}$

Also estimate the context distribution by $\widehat{\nu}_t = \frac{1}{t} \sum_{s=1}^t \delta_{\mathbf{x}_s}$

Phase 2 Compute the solution $\widehat{\mathbf{p}}_t(h_{t-1}, \cdot)$ of

$$\begin{aligned} \text{OPT}(\widehat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, B_T) &= \max_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} T \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} U_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] \\ \text{under } T \mathbb{E}_{\mathbf{X} \sim \widehat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{L}_{t-1}(a, \mathbf{X}) \pi_a(\mathbf{X}) \right] &\leq B_T \mathbf{1} \end{aligned}$$

Draw an arm $a_t \sim \widehat{\mathbf{p}}_t(h_{t-1}, \mathbf{x}_t)$

E.3 Regret Bound

We sketch in Appendix E.4 below the proof of the following result. We use the $\varepsilon_s(a, \mathbf{x})$ indicated by Lemma 3.

Theorem 3. *In the setting of Appendix E.1, we consider the adaptive policy of Box C of Appendix E.2. We set a confidence level $1 - \delta \in (0, 1)$ and use parameters $\lambda = m \ln(1 + T/m)$, a working budget of $\hat{B}_T = B - b_T$, where*

$$b_T = 2 + m(2\sqrt{2}\|\Theta\| + 1)\sqrt{T} \ln \frac{1 + T/m}{\delta/(d+1)} + \sqrt{2T \ln \frac{4d}{\delta}} + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}},$$

and $\varepsilon_t(a, \mathbf{x}) = \gamma_{t,\lambda,\delta} \|\varphi(a, \mathbf{x})\|_{X_t^{-1}}$. Then, provided that $T \geq 2m$ and $B > 2b_T$, we have, with probability at least $1 - 3\delta$,

$$\text{OPT}(\nu, \bar{r}, \bar{c}, B) - \sum_{t=1}^T r_t \leq 2b_T \left(1 + \frac{\text{OPT}(\nu, \bar{r}, \bar{c}, B)}{B} \right).$$

The order of magnitude of the regret bound, in terms of T , m , and B is

$$\frac{\text{OPT}(\nu, \bar{r}, \bar{c}, B)}{B} m \sqrt{T} \ln T.$$

This matches the regret bound achieved by Agrawal and Devanur [2016, Theorem 3], except that the latter reference required a budget B of order $T^{3/4}$ up to logarithmic terms, while we relax the budget amount to $B \geq 2b_T$, i.e., B of order \sqrt{T} up to logarithmic terms.

Also, and more importantly, we provide a natural strategy in Box C, whose parameters are easy to tune, while the fully adaptive algorithm underlying Agrawal and Devanur [2016, Theorem 3] has to estimate a critical parameter Z to trade off rewards and costs (the equivalent of our $\hat{\beta}_t^{\text{budg},*}$ dual optimal variable). This Z should be of order $\text{OPT}(\nu, \bar{r}, \bar{c}, B)/B$ and \sqrt{T} initial rounds of the strategy underlying Agrawal and Devanur [2016, Theorem 3] are devoted to computing a suitable value of Z .

We also provide, in the analysis of Appendix E.4, a rigorous treatment of the use of the no-op action a_{null} .

However, the main advantage of Agrawal and Devanur [2016, Theorem 3] over Theorem 3 above lies in the absence of finiteness restriction on the context set \mathcal{X} , which we have to (somewhat artificially) introduce to ensure that the linear program of Phase 2 of the adaptive policy of Box C is tractable.

E.4 Proof Sketch of Theorem 3

We use the same $\hat{\cdot}$ conventions as in Appendix D. The main (but rather minor) changes with respect to the proofs of Appendices B and D are specifically underlined below. The reason why it is handy to consider instead a lower-confidence bound on the vector costs is to be found in Step 4 below.

Step 1. The first step corresponds to Lemma 3 above, together with the introduction of the bound E_T . Given that we pick $\lambda = m \ln(1 + T/m) \geq 1 \geq 1$, we get the following counterpart of (12), by replacing κ by 1:

$$2 \sum_{t=2}^T \varepsilon_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} \leq 2\gamma_{T,\lambda,\delta} \sqrt{2mT \ln \left(1 + \frac{T}{\lambda m} \right)} \stackrel{\text{def}}{=} E_T.$$

Substituting the value of λ and upper bounding $\gamma_{T,\lambda,\delta}$ by

$$\gamma_{T,\lambda,\delta} \leq \left(\|\Theta\| + \frac{1}{4} \right) \sqrt{m \ln \frac{1 + T/m}{\delta/(d+1)}},$$

we get

$$E_T \leq m(2\sqrt{2}\|\Theta\| + 1)\sqrt{T} \ln \frac{1 + T/m}{\delta/(d+1)}.$$

Step 2. There are three main outcomes of Step 2 (see the summary in Appendix D.2). Up to considering the new $U_t(a, \mathbf{x})$ and $\mathbf{L}_t(a, \mathbf{x})$ in lieu of the $r(a, \mathbf{x}) U_t(a, \mathbf{x})$ and $\mathbf{c}(a, \mathbf{x}) U_t(a, \mathbf{x})$, respectively, we have the following counterparts to (18), (20) and (25). For each $t \geq 2$ such that the cost constraints of Phase 0 of the adaptive policy of Box C are not violated and the optimization problem $\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)$ is to be solved, there exists a vector $\hat{\beta}_t^{\text{budg},*} \geq \mathbf{0}$ such that the complementary slackness condition of KKT reads

$$(\hat{\beta}_t^{\text{budg},*})^\top \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{L}_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] = \frac{\hat{B}_T}{T} (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{1}, \quad (61)$$

the policy $\hat{p}_t(h_{t-1}, \cdot)$ satisfies

$$\hat{p}_t(h_{t-1}, \mathbf{x}) \in \operatorname{argmax}_{\mathbf{q} \in \mathcal{P}(\mathcal{A})} \sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right) q_a, \quad (62)$$

and the value of the optimization problem is larger than

$$\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T) \geq \hat{B}_T (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{1}. \quad (63)$$

Step 3. The uniform deviation argument (8), formulated equivalently as (44), still holds, on an event referred to as $\hat{\mathcal{E}}_{\text{unif},\delta}$. Also, we assumed that rewards r_t and cost vectors \mathbf{c}_t are bounded in $[0, 1]$ and $[0, 1]^d$, respectively. Several applications of the Hoeffding-Azuma inequality, together with a union bound, show that there exist an event $\hat{\mathcal{E}}_{\text{HAz},\delta}$ of probability at least $1 - \delta$ such that, simultaneously, various high-probability controls similar to (52)–(55) hold. We do not rewrite them explicitly.

On the intersection $\hat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \hat{\mathcal{E}}_{\text{unif},\delta}$, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbf{c}_t \\ & \leq \sum_{t=1}^T \bar{\mathbf{c}}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} + \sqrt{\frac{T}{2} \ln \frac{4d}{\delta}} \mathbf{1} \\ & \leq 1 + \sum_{t=2}^T \mathbf{L}_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} + \left(E_T + \sqrt{\frac{T}{2} \ln \frac{4d}{\delta}} \right) \mathbf{1} \\ & \leq 1 + \sum_{t=2}^T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{L}_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] + \left(E_T + \sqrt{2T \ln \frac{4d}{\delta}} \right) \mathbf{1} \\ & \leq 1 + \sum_{t=2}^T \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{L}_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] + \left(E_T + \sqrt{2T \ln \frac{4d}{\delta}} + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} \right) \mathbf{1}, \end{aligned}$$

where, among others, we used (60) and the definition of E_T for the second inequality. Note that the E_T term was not necessary in Appendices B and D as we were then using an upper-confidence bound on the vector costs, obtained thanks to an upper-confidence bound on the conversion rate. At each round $t \geq 2$, whether the strategy picks $\hat{p}_{t,a}(h_{t-1}, \cdot)$ in Phase 0 (in which case the left-hand side in the display below equals $\mathbf{0}$) or Phase 2 (in which case we have an equality in the display below), it holds that

$$\mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \mathbf{L}_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] \leq \frac{\hat{B}_T}{T} \mathbf{1}.$$

Substituting the value of \hat{B}_T , we proved that on $\hat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \hat{\mathcal{E}}_{\text{unif},\delta}$, which is an event of probability at least $1 - 3\delta$,

$$\sum_{t=1}^T \mathbf{c}_t \leq (B - 1) \mathbf{1}.$$

This shows that on this intersection of events, the adaptive policy of Box C always resorts to Phase 2. We will consider this event for the rest of the proof, so that the results of Step 2 may be applied.

We may proceed similarly for rewards, based on (59): on $\widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$

$$\begin{aligned}
& \sum_{t=1}^T r_t \\
& \geq \sum_{t=1}^T \bar{r}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} - \sqrt{\frac{T}{2} \ln \frac{4}{\delta}} \\
& \geq \sum_{t=2}^T U_{t-1}(a_t, \mathbf{x}_t) \mathbb{1}_{\{a_t \neq a_{\text{null}}\}} - \left(E_T + \sqrt{\frac{T}{2} \ln \frac{4}{\delta}} \right) \\
& \geq \sum_{t=2}^T \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} U_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right] - \left(E_T + \sqrt{2T \ln \frac{4}{\delta}} \right) \\
& \geq \underbrace{\sum_{t=2}^T \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} U_{t-1}(a, \mathbf{X}) \hat{p}_{t,a}(h_{t-1}, \mathbf{X}) \right]}_{=\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)/T} - \underbrace{\left(E_T + \sqrt{2T \ln \frac{4}{\delta}} + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} \right)}_{\leq b_T},
\end{aligned} \tag{64}$$

where the indicated equality to $\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)$ follows from the fact that on the considered intersection of events, the adaptive policy always resorts to Phase 2. We also used the piece of notation b_T introduced in the statement of Lemma 3.

Step 4. We build on (62) as follows. By the existence of a_{null} , the maximum in (62) can be taken with non-negative parts. We also substitute the upper confidence control (59) and the lower confidence control (60)—this piece of the proof is the very reason why such upper and lower confidence estimates were picked. We get: on $\widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$, for all $t \geq 2$, for all $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned}
& \sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right) \hat{p}_{t,a}(h_{t-1}, \mathbf{x}) \\
& = \sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right)_+ \hat{p}_{t,a}(h_{t-1}, \mathbf{x}) \\
& \geq \sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right)_+ \pi_a^*(\mathbf{x}) \\
& \geq \sum_{a \in \mathcal{A}} \left(\bar{r}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \bar{\mathbf{c}}(a, \mathbf{x}) \right)_+ \pi_a^*(\mathbf{x}).
\end{aligned}$$

The rest of the proof follows the exact same logic as in Appendices B and D. By replacing the \mathbf{x} by a random variable $\mathbf{X} \sim \hat{\nu}_t$ and integrating, we have, on $\widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$, that for all $t \geq 2$,

$$\begin{aligned}
& \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right) \hat{p}_{t,a}(h_{t-1}, \mathbf{x}) \right] \\
& \geq \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(\bar{r}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \bar{\mathbf{c}}(a, \mathbf{x}) \right)_+ \pi_a^*(\mathbf{x}) \right] \\
& \geq \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \left(\bar{r}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \bar{\mathbf{c}}(a, \mathbf{x}) \right)_+ \pi_a^*(\mathbf{x}) \right] - |\mathcal{X}| \sqrt{\frac{t}{2} \ln \frac{2T|\mathcal{X}|}{\delta}},
\end{aligned}$$

where the second inequality follows by (44), which is legitimately applied thanks the fact that the sum over a in the right-hand side takes values in $[0, 1]$, given the non-negative parts and the fact that $\bar{r}(a, \mathbf{x}) \leq 1$ by definition. Now, with (61) and the definition of π^* : on $\widehat{\mathcal{E}}_{\text{HAz},\delta} \cap \mathcal{E}_{\text{lin},\delta} \cap \widehat{\mathcal{E}}_{\text{unif},\delta}$, for all

$t \geq 2$,

$$\begin{aligned}
& \frac{\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)}{T} - \frac{\hat{B}_T}{T} (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{1} \\
&= \mathbb{E}_{\mathbf{X} \sim \hat{\nu}_t} \left[\sum_{a \in \mathcal{A}} \left(U_{t-1}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{L}_{t-1}(a, \mathbf{x}) \right) \hat{p}_{t,a}(h_{t-1}, \mathbf{x}) \right] \\
&\geq \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \left(\bar{r}(a, \mathbf{x}) - (\hat{\beta}_t^{\text{budg},*})^\top \bar{\mathbf{c}}(a, \mathbf{x}) \right)_+ \pi_a^*(\mathbf{x}) \right] - |\mathcal{X}| \sqrt{\frac{t}{2} \ln \frac{2T|\mathcal{X}|}{\delta}} \\
&\geq \frac{\text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B)}{T} - \frac{B}{T} (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{1} - |\mathcal{X}| \sqrt{\frac{t}{2} \ln \frac{2T|\mathcal{X}|}{\delta}}.
\end{aligned}$$

Based on these inequalities, we have

$$\begin{aligned}
& \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) - \sum_{t=2}^T \frac{\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)}{T} \\
&\leq |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} + 1 + \sum_{t=2}^T \frac{B - \hat{B}_T}{T} (\hat{\beta}_t^{\text{budg},*})^\top \mathbf{1} \\
&\leq |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}} + 1 + \frac{B - \hat{B}_T}{\hat{B}_T} \sum_{t=2}^T \frac{\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)}{T},
\end{aligned}$$

where we substituted (63) for the second inequality. By a case analysis, we finally proved that on $\hat{\mathcal{E}}_{\text{HAZ}, \delta} \cap \mathcal{E}_{\text{lin}, \delta} \cap \hat{\mathcal{E}}_{\text{unif}, \delta}$,

$$\begin{aligned}
& \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) - \sum_{t=2}^T \frac{\text{OPT}(\hat{\nu}_t, U_{t-1}, \mathbf{L}_{t-1}, \hat{B}_T)}{T} \\
&\leq \left(\frac{B}{\hat{B}_T} - 1 \right) \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) + \underbrace{1 + |\mathcal{X}| \sqrt{2T \ln \frac{2T|\mathcal{X}|}{\delta}}}_{\leq b_T}.
\end{aligned}$$

The proof is concluded by combining the inequality above with the bound (64) on cumulative rewards:

$$\text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B) - \sum_{t=1}^T r_t \leq 2b_T + \underbrace{\left(\frac{B}{\hat{B}_T} - 1 \right)}_{\leq 2b_T} \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B),$$

where the bound on $B/\hat{B}_T - 1$ is obtained with the same techniques and similar conditions as for (32).

E.5 Reminder: Algorithm 1 from Agrawal and Devanur [2016]

We recall (and actually slightly adapt to our setting) the adaptive policy of Agrawal and Devanur [2016] titled Algorithm 1 therein. We describe it in Box D. One of the adaptations is to state it with general upper-confidence bonuses $\varepsilon_s(a, \mathbf{x}) > 0$. As in Agrawal and Devanur [2016], who proceed as in the proof of Lemma 3, we will use the same values for $\varepsilon_s(a, \mathbf{x})$ as in Theorem 3. The same comment applies to λ . Another adaptation is that we specified the online convex optimization algorithm to be used and picked a simple strategy (instead of other possible choices discussed in Agrawal and Devanur [2016]), namely, the projected gradient descent introduced by Zinkevich [2003]. The latter relies on a learning rate $\eta > 0$. The drawback of the projected gradient descent is however that its dependency in the ambient dimension is suboptimal.

The final parameter of the adaptive policy of Box D is a parameter Z to trade off between rewards and costs. A recommended choice is, for instance, $Z = \text{OPT}(\nu, \bar{r}, \bar{\mathbf{c}}, B)/B$, the issue being that, of course, the latter value is unknown. In the simulation study of Appendix F, we will provide a good value of Z to the adaptive policy, even though Agrawal and Devanur [2016] introduce a variant with a preliminary exploration phase meant to provide in an automatic way such a good value for Z .

BOX D: ADAPTATION OF ALGORITHM 1 FROM AGRAWAL AND DEVANUR [2016]

Parameters: regularization parameter $\lambda > 0$; trade-off parameter Z between reward and costs; upper-confidence bonuses $\varepsilon_s(a, \mathbf{x}) > 0$, for $s \geq 1$ and $(a, \mathbf{x}) \in (\mathcal{A} \setminus \{a_{\text{null}}\}) \times \mathcal{X}$; learning rate $\eta > 0$.

Round $t = 1$: play an arbitrary action $a_1 \in \mathcal{A} \setminus \{a_{\text{null}}\}$; pick $\zeta_1 = \mathbf{0}$

At rounds $t \geq 2$:

Phase 0 If $\sum_{s \leq t-1} \mathbf{c}_s \leq (B-1)\mathbf{1}$ is violated, then $a_t = a_{\text{null}}$

Phase 1 Otherwise, estimate the parameters by

$$\begin{aligned} \boldsymbol{\mu}_{t-1} &= X_{t-1}^{-1} \sum_{s=1}^{t-1} \boldsymbol{\varphi}(a_s, \mathbf{x}_s) r_s \\ \text{and } \hat{\boldsymbol{\theta}}_{t-1,i} &= X_{t-1}^{-1} \sum_{s=1}^{t-1} \boldsymbol{\varphi}(a_s, \mathbf{x}_s) c_{s,i} \\ \text{where } X_t &= \sum_{s=1}^t \boldsymbol{\varphi}(a_s, \mathbf{x}_s) \boldsymbol{\varphi}(a_s, \mathbf{x}_s)^\top + \lambda \mathbf{I}_m \end{aligned}$$

Define the expected reward function \hat{r} and cost function $\hat{\mathbf{c}}_{t-1} = (\hat{c}_{t-1,i})_{1 \leq i \leq d}$ as

$$\begin{aligned} \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \quad \hat{r}_{t-1}(a, \mathbf{x}) &= \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\mu}_{t-1} \\ \text{and } \forall 1 \leq i \leq d, \quad \hat{c}_{t-1,i} &= \boldsymbol{\varphi}(a, \mathbf{x})^\top \boldsymbol{\theta}_{t-1,i} \end{aligned}$$

Build the upper-confidence expected reward function U_{t-1} and the lower-confidence expected vector-cost function \mathbf{L}_{t-1} as

$$\begin{aligned} \forall a \in \mathcal{A} \setminus \{a_{\text{null}}\}, \forall \mathbf{x} \in \mathcal{X}, \\ U_{t-1}(a, \mathbf{x}) &= \hat{r}_{t-1}(a, \mathbf{x}) + \varepsilon_{t-1}(a, \mathbf{x}) \\ \mathbf{L}_{t-1}(a, \mathbf{x}) &= \hat{\mathbf{c}}_{t-1} - \varepsilon_{t-1}(a, \mathbf{x}) \mathbf{1} \end{aligned}$$

Phase 2 Play

$$a_t \in \operatorname{argmax}_{a \in \mathcal{A} \setminus \{a_{\text{null}}\}} U_{t-1}(a, \mathbf{x}) - Z (\boldsymbol{\zeta}_{t-1}^\top \mathbf{L}_{t-1}(a, \mathbf{x}))$$

Compute

$$\boldsymbol{\zeta}_t = \Pi_{\text{unit}} \left(\boldsymbol{\zeta}_{t-1} + \eta (\mathbf{c}_{t-1} - (B/T)\mathbf{1}) \right)$$

where Π_{unit} denotes the Euclidean projection onto the unit ℓ_1 -ball

$$\{\boldsymbol{\zeta} \in \mathbb{R}^d : \boldsymbol{\zeta} \geq \mathbf{0} \text{ and } \zeta_1 + \dots + \zeta_d \leq 1\}$$

F Simulation Study

This appendix reports numerical simulations performed on partially simulated but realistic data (Appendix F.1), for the sake of illustration only. We describe (Appendix F.2) the specific experimental setting of CBwK for a conversion model considered—i.e., the features available, the parameters of the logistic regression, the reward and cost functions. A key point is that continuous variables are used to define rewards, costs, and even the conversion rate, while the adaptive policy of Box C must discretize these variables to solve its Phase 2 linear program. Though the experimental setting introduced is not a setting of linear CBwK, we may still apply the Box D adaptive policy (Appendix F.3), with the underlying idea that it provides linear approximations to non-linear reward and cost functions. We carefully explain how the hyper-parameters picked (Appendix F.4) before providing and discussing the outcomes of the simulations (Appendix F.5). The main outcome is that, as expected, the ad hoc adaptive policy of Box C outperforms the adaptive policy of Box D, which essentially linearly approximates non-linear rewards and costs. We end with a note (Appendix F.6) on the computation environment and time.

F.1 Data Preparation and Available Contexts

The underlying dataset for the simulations is the standard “default of credit card clients” dataset of UCI [2016], initially provided by Yeh and Lien [2009]. (It may be used under a Creative Commons Attribution 4.0 International [CC BY 4.0] license.) This dataset is originally for comparing algorithms predicting default probability of credit card clients. It includes some socio-demographic data, debt level, and payment/default history of the clients. It also includes a target measuring whether the client will default in the future (1-month ahead). We transform it to match our motivating application of market share expansion for loans, described in Appendix A. To do so, we consider each line of the dataset as a loan application. We then discard some variables (e.g., the target) and create new ones (requested amount, standard interest rate offered, risk score). Below, we begin with describing the variables that we keep as they are and explain next how we created the additional variables, based on existing ones.

Variables retrieved. The variable *Age* provides the age of a given client at the time of the loan application, in years. We discretize it into 5 levels with similar numbers of loan requests in each level. The cutoffs for each level are 27, 31, 37, and 43, respectively. This gives rise to a variable referred to as *Age-discrete*.

The variable *Education* reports the education level of a client; in the data there are 4 levels: “others” (level 1, representing 2% of clients), “high school” (level 2, with a share of 16%), “university degree” (level 3, with 47%), and “graduate school degree” (level 4, with 35%).

Finally, the variable *Marital status* provides the marital status of a client: “others” (level 1, accounting for 1.3% of clients), “single” (level 2, for 53.3% of clients), and “married” (level 3, for 45.4%).

Variables created based on existing ones. We create a variable *Requested amount*, in dollars (\$), based on a variable provided in the data set that measures the current debt level, in dollars, of the clients: we do so by multiplying the debt level by 0.2. We cap the value of *Requested amount* to 100K\$. We then discretize *Requested amount* into 5 levels with similar numbers of loan requests in each level; the obtained variable is referred to as *Requested amount-discrete*. The cutoffs for each level are 10K\$, 20K\$, 36K\$, and 54K\$, respectively.

For the final two variables, *Standard interest rate* and *Risk score*, we first build a probability-of-default model with the variables from the raw database as predictors and the occurrence of a default within the next month as a target. This probability-of-default model is based on XGBoost (Chen and Guestrin [2016]), run with no penalization, depth 3, learning rate 0.01, subsample parameter 0.8, min child weight 10, and number of trees 1,176. We only set the number of trees by cross validation, while the rest of the hyper-parameters were set arbitrarily. As the default target is on a credit card, the predicted default rate seems high compared to what we deem as typical default rates on loans. We therefore divide the predicted probability of default by factor of 4 and cap this probability to 20%. This gives us a working variable called *PD*, for probability of default.

We build a *Risk score* rating the risk of a client’s default, with 5 levels, coded from A (level 1) to E (level 5), where E represents the highest risk. It is created based on the 20% - 40% - 60% - 80% quantiles of *PD*.

We finally set the *Standard interest rate* variable as 0.9 times the *PD*, with a maximal value of 18% and a minimum one of 1%. This constitutes an oversimplification of risk-based pricing, since we do not take into account any loss given default; but we do not have enough information in the dataset to do so, which is why we basically assume that the loss given default is constant. Then, theory has it that the *Standard interest rate* can be considered proportional to *PD*, and we carefully picked the factor 0.9 to get realistic values. In the dataset, the thus created *Standard interest rate* variable exhibits an average of 4.9%, with median 3.3%. An important note is that this variable is continuous and does not take finitely many values, while the setting described in Section 2.1 imposes such a restriction (the linear setting of Section 6 does not require it). We discuss below how we take this fact into account: by only using this variable for the conversion model (i.e., in Phase 1 of the adaptive policy of Box C), but not to pick actions (i.e., not in Phase 2 of the adaptive policy of Box C).

Lastly, we filter the database to remove the outliers: the lines for which *Standard interest rate* times *Requested amount* is larger than 10K, which happens for 133 clients. Our final database contains 29,867 loan applications, out of which we will bootstrap $T = 50,000$ applications.

Additional comments. Note that for *Requested amount* and *Age*, the discretizations performed aim to get five balanced classes; however, as some requests are with some specific boundary values, we do not get exact equal distributions over the classes.

All the parameters, constants, and cutoff/filter thresholds used in this data preparation step were decided arbitrarily and were not based on any real information. The context variables here were also selected somewhat arbitrarily (based on their availability), and solely for illustration purposes. In reality, the variables that can be used for commercial discounts need to comply with relevant laws, regulations and company’s internal compliance rules.

Summary. The context \mathbf{x} for a given client thus contains the following variables: *Age*, *Age-discrete*, *Education*, *Marital status*, *Requested amount*, *Requested amount–discrete*, *Risk score*, and *Standard interest rate*. Categorical variables are hot-one encoded via binary variables.

F.2 Specific Setting of CBwK for Logistic Conversions

We recall that our aim is to provide simulations matching the motivating example of market share expansion for loans described in Appendix A. We take as action set, i.e., as possible discount rates, $\mathcal{A} = \{a_{\text{null}}, 0.1, 0.2, 0.35, 0.55, 0.8\}$.

Features. The feature vectors $\varphi_{\text{conv}}(a, \mathbf{x})$ used are composed of only some of the variables defining the context \mathbf{x} , namely, *Age-discrete*, *Education*, *Marital status*, *Requested amount–discrete*, *Risk score*, and *Standard interest rate* (we recall that this variable is not discrete but will not use it in the linear program of Phase 2), with the addition of a new variable called *Final interest rate* equal to the discounted standard interest rate offered, i.e., $\text{Final interest rate} = \text{Standard interest rate} \times (1 - a)$.

Reward and vector cost functions. We use the following (normalized) reward and cost functions, inspired from Appendix A. We set a common duration for all loans, say, 2 years, so that the requested amount equals the outstanding. For all $a \in \mathcal{A} \setminus \{a_{\text{null}}\}$ and \mathbf{x} ,

$$r(a, \mathbf{x}) = x_{\text{am}}/M_{\text{am}} \quad \text{and} \quad \mathbf{c}(a, \mathbf{x}) = (a/M_{\text{disc}}, x_{\text{ir}}x_{\text{am}}/M_{\text{ir,am}}) \quad (65)$$

where x_{am} and x_{ir} denote the components of the context \mathbf{x} containing the *Requested amount* and *Standard interest rate*, respectively, and the normalization factors equal $M_{\text{am}} = 10^5$, $M_{\text{disc}} = 7$, and $M_{\text{ir,am}} = 9,996$.

Note that the definition of the first cost here is different from that in Appendix A: we use a/M_{disc} here instead of $\mathbb{1}_{\{a \neq 0\}}$ there. We do so not to disfavor the Box D policy, see details in Appendix F.3 below.

Conversion rate function P . We model the conversion rate function P with the logistic-regression model stated in (1), with $\varphi = \varphi_{\text{conv}}$, and only need to provide the numerical value of θ_* , which we do

Table 1: Coefficients picked for the logistic-regression model of P .

Intercept	0.8177				
Continuous variable	Single coefficient				
<i>Final interest rate</i>	−13.1101				
Discrete variables	Coefficients for each level				
	Level 1	Level 2	Level 3	Level 4	Level 5
<i>Risk score</i>	−0.3045	−0.0383	0.0515	0.1261	0.1636
<i>Requested amount–discrete</i>	0.7093	0.4703	0.1113	−0.2748	−1.0179
<i>Age–discrete</i>	−0.1837	−0.1392	−0.0476	0.1096	0.2592
<i>Education</i>	0.1836	0.0126	−0.0896	−0.1084	
<i>Marital status</i>	0.0799	0.0102	−0.0918		

in Table 1. This model, as well as the Phase 1 learning of θ_* described by (3)–(5) in Section 3, holds for possibly non-discrete contexts.

The numerical values picked for θ_* were so in some arbitrary way, to get somewhat realistic outcomes with a simple model structure. We imposed monotonicity constraints, as these are most natural: for instance, the conversion rate increases with the level of *Risk score* and *Age–discrete* increase, and decreases with the level of *Education*, *Requested amount–discrete*, and *Final interest rate*. The coefficients for *Marital status* indicates that conversions are more likely for clients that are single than for married clients.

The average conversion rate in the case $a = 0$ of no discount (i.e, by replacing *Final interest rate* by *Standard interest rate*) is around 50%.

Adaptive policy: based only on the discrete variables. As indicated above, the logistic-regression model and the learning of its parameters apply to continuous variables. The restriction that the context set \mathcal{X} should be finite only came from Phase 2 of the Box C adaptive policy, i.e., the linear program—in particular, for it to be computationally tractable. Here, we thus restrict our attention to policies that map the discrete variables in \mathbf{x} to distributions over \mathcal{A} : policies that ignore the variables *Age*, *Requested amount*, and *Standard interest rate*. For the first two variables, they may use their discretized versions *Age–discrete* and *Requested amount–discrete*. For *Standard interest rate*, given how it was constructed, *Risk Score* appears as its discretized version.

The aim of these simulations is to show, among others, that using discretizations only in Phase 2 is relevant and efficient.

Note that, on the theoretical side, the proof sketches provided in Sections 4 and 5 reveal that the errors $\varepsilon_t(a, \mathbf{x})$ for learning θ_* and P , obtained as outcomes of the first step of the analyses, are carried over in the subsequent steps, where the optimization part is evaluated. Using discretizations only in Phase 2 does therefore not come at the price of loosing theoretical guarantees.

E.3 Consideration of the Box D Adaptive Policy for Linear CBwK

In these experiments, we also consider the Box D adaptive policy of Appendix E.5, which was introduced by Agrawal and Devanur [2016] in a different setting. To be as fair as possible to this adaptive policy, we do so with the extended features $\varphi_{\text{lin}}(a, \mathbf{x})$ consisting of the features $\varphi_{\text{conv}}(a, \mathbf{x})$ described above and three additional components: the discount a , the *Requested amount* x_{am} , and the product $x_{\text{ir}}x_{\text{am}}$ of the *Requested amount* by the *Standard interest rate*. Actually, to ensure that $\varphi_{\text{lin}}(a, \mathbf{x}) \in [0, 1]^m$, the last two components added are normalized: we rather use $x_{\text{am}}/M_{\text{am}}$ and $x_{\text{ir}}x_{\text{am}}/M_{\text{ir,am}}$. The reward and vector cost functions introduced in (65) are linear in $\varphi_{\text{lin}}(a, \mathbf{x})$. Even better, each component of $r(a, \mathbf{x})$ and $\mathbf{c}(a, \mathbf{x})$ is given directly by a component of $\varphi_{\text{lin}}(a, \mathbf{x})$ —an extremely simple linear dependency on $\varphi_{\text{lin}}(a, \mathbf{x})$.

However, the expected reward and cost functions

$$\bar{r}(a, \mathbf{x}) = r(a, \mathbf{x}) P(a, \mathbf{x}) \quad \text{and} \quad \bar{\mathbf{c}}(a, \mathbf{x}) = \mathbf{c}(a, \mathbf{x}) P(a, \mathbf{x}),$$

which are the ones that should be linear in $\varphi_{\text{in}}(a, \mathbf{x})$ according to the setting described in Appendix E.1, are not linear in these features. This is due to the $P(a, \mathbf{x})$ terms, which are given by logistic regressions.

The Box D adaptive policy is therefore disadvantaged. This is even more true as it models rewards and costs independently, while they are coupled through conversions. We nonetheless consider this linear-modeling policy because a typical justification for linear approximations is that they offer a typical and efficient first-stage approach to possibly complex problems. Another reason was the desire to have some competitor to our policy in the simulations, and Box D adaptive policy was an easy-to-implement strategy—unlike the policies by Badanidiyuru et al. [2014] and Agrawal et al. [2016], which rely on considering finitely many benchmark policies.

All in all, we report the performance of the Box D policy as well in our experiments, though, as expected, the ad-hoc Box C policy outperforms it.

F.4 Hyper-parameters Picked

We actually set the hyper-parameters based on the budget B , and therefore, first explain how we set its possible values. It turns out that setting $B \geq 3,650$ is equivalent to not imposing any constraint, while setting $B \geq 2,900$ is equivalent to no second budget constraint. To get meaningful results, we therefore picked $B = 1,600$ and $B = 2,200$ as the two possible values for B . We now describe how we tune each of the two adaptive policies considered.

Hyper-parameters common to the two adaptive policies. We take $T = 50,000$ clients in the experiments, by bootstrapping them from the enriched dataset prepared in Appendix F.1. We set initial 50 rounds as a warm start for the sequential logistic regression and sequential linear regression carried out in Phase 1 of the adaptive policies.

Both adaptive policies use upper-confidence bonuses $\varepsilon_s(a, \mathbf{x})$, which are roughly of the form (considering λ as a constant)

$$C(1 + \ln s) \|\varphi(a, \mathbf{x})\|_{X_s^{-1}},$$

where the matrix X_s was defined in Box C; for simplicity, we set $\kappa = 1$, so that the matrices V_s of Lemma 1 and X_s are equal, which explains the common form of the upper-confidence bonuses $\varepsilon_s(a, \mathbf{x})$. The hyper-parameter C controls the exploration: the higher C , the more exploration. We report in the simulations the results achieved for C in the range $\{0.025, 0.1, 0.3\}$. That range was set so that at round $s = 51$, which is the first round after the warm start, the $\varepsilon_{51}(a, \mathbf{x})$ take values around 0.05, 0.3, and 0.9, respectively.

For simplicity, we set $B_T = B$ as a working budget.

Hyper-parameters for the Box C adaptive policy. We feed this adaptive policy with a good value of λ , namely, $\lambda = 0.0129$. We obtained it by cross-validation on an independent T -sample of data, using the Phase 1 estimation. In the T -sample for estimating λ , at each round s , we take action from the optimal static policy and use the associated conversion y_s as target for estimation. We omit the projection step in Phase 1 by considering that a large enough set Θ was picked.

Hyper-parameters for the Box D adaptive policy. As discussed in Appendix E.5, we set $Z = \text{OPT}(\nu, P, B)/B$, that is, $Z = 5.16$ for $B = 1,600$ and $Z = 3.87$ for $B = 2,200$. We also set $\lambda = 0.2452$ for $B = 1,600$ and $\lambda = 0.2765$ for $B = 2,200$. These values were obtained as weighted averages: the sum of 0.5 times the optimal λ for rewards and 0.25 times the optimal λ for each of the two cost components. These optimal λ s for rewards and cost components were set by cross validation on an independent T -sample; with actions a_s taken at each round s from the optimal static policy and associated rewards r_s and costs c_s as targets.

Finally, the learning parameter η was selected in the range $\{0.005, 0.01, 0.05, 0.1, 0.2\}$. We did so given the other choices, by picking in hindsight the η with best performance; this of course, just like the clever choice of Z , should give an advantage to the Box D adaptive policy. Namely, when $B = 1,600$, for C equal 0.025, 0.1, and 0.3, we selected η equal to 0.05, 0.01, and 0.1, respectively; and when $B = 2,200$, we selected η equal to 0.01, 0.01, and 0.05, respectively. When performing these retrospective choices, we however noted that the performance was not significantly impacted by the choice of η .

F.5 Outcomes of the Simulations

We were limited by the computational power (see Appendix F.6) and could only perform 10 simulations for each pair of $B \in \{1,600, 2,200\}$ and $C \in \{0.025, 0.1, 0.3\}$. We report averages (strong lines) as well as ± 2 times the standard errors (shaded areas).

Figure 1 reports, in the first line of graphs, the regret achieved with respect to what achieves the optimal static policy, i.e.,

$$t \mapsto \frac{t}{T} \text{OPT}(\nu, P, B) - \sum_{s=1}^t r_s,$$

where $\text{OPT}(\nu, P, B)$ is larger than 8,000 for both values of B . This regret can take negative or positive values, but in expectation, it is non-negative. This is not immediately clear from the figure, which reports the empirical averages of the regret over 10 runs: these empirical averages are sometimes negative, but they always lie in confidence intervals containing the value 0.

The figures also reports, in the second and third lines, the difference between a constant linear increase of the costs (between a 0 initial cost and a final B cost) and the costs actually achieved by the adaptive policies. I.e., these graphs report the averages and standard errors of the following quantities: for each cost component $i \in \{1, 2\}$,

$$t \mapsto \sum_{s=1}^t c_{s,i} - \frac{t}{T} B;$$

by design, the difference above must be non-positive. The second line of Figure 1 deals with the first cost component, and its third line reports the results for the second cost component.

The experiments reveal that while both adaptive policies seem to achieve sublinear regret, the Box D adaptive policy, which is suited for the CBwK setting for a conversion model, performs better than the Box C adaptive policy in terms of rewards: it achieves a smaller, sometimes negative, regret. In terms of costs, we globally see the same trend, with, for a given value of C , the Box D adaptive policy suffering smaller costs than the Box C adaptive policy while achieving higher rewards. This hints at a better use of the discounts.

The hyper-parameter C has an interesting impact: the lower C , the lower the regret (the higher the rewards) and the lower the costs. Rewards and costs go hand in hand: for a given adaptive policy, higher rewards are associated with higher costs.

F.6 Computation Time and Environment

As requested by the NeurIPS checklist, we provide details on the computation time and environment. Our experiments were ran on the following hardware environment: no GPU was required, CPU is 2.7 GHz Quad-Core with total of 8 threads and RAM is 16 GB 2133 MHz LPDDR3. We ran 5 simulations with different seeds on parallel each time. In the setting and for the data described above, it took us 8 hours for each such bunch of 5 runs of the adaptive policy of Box C, and 1.5 hours for Box D.

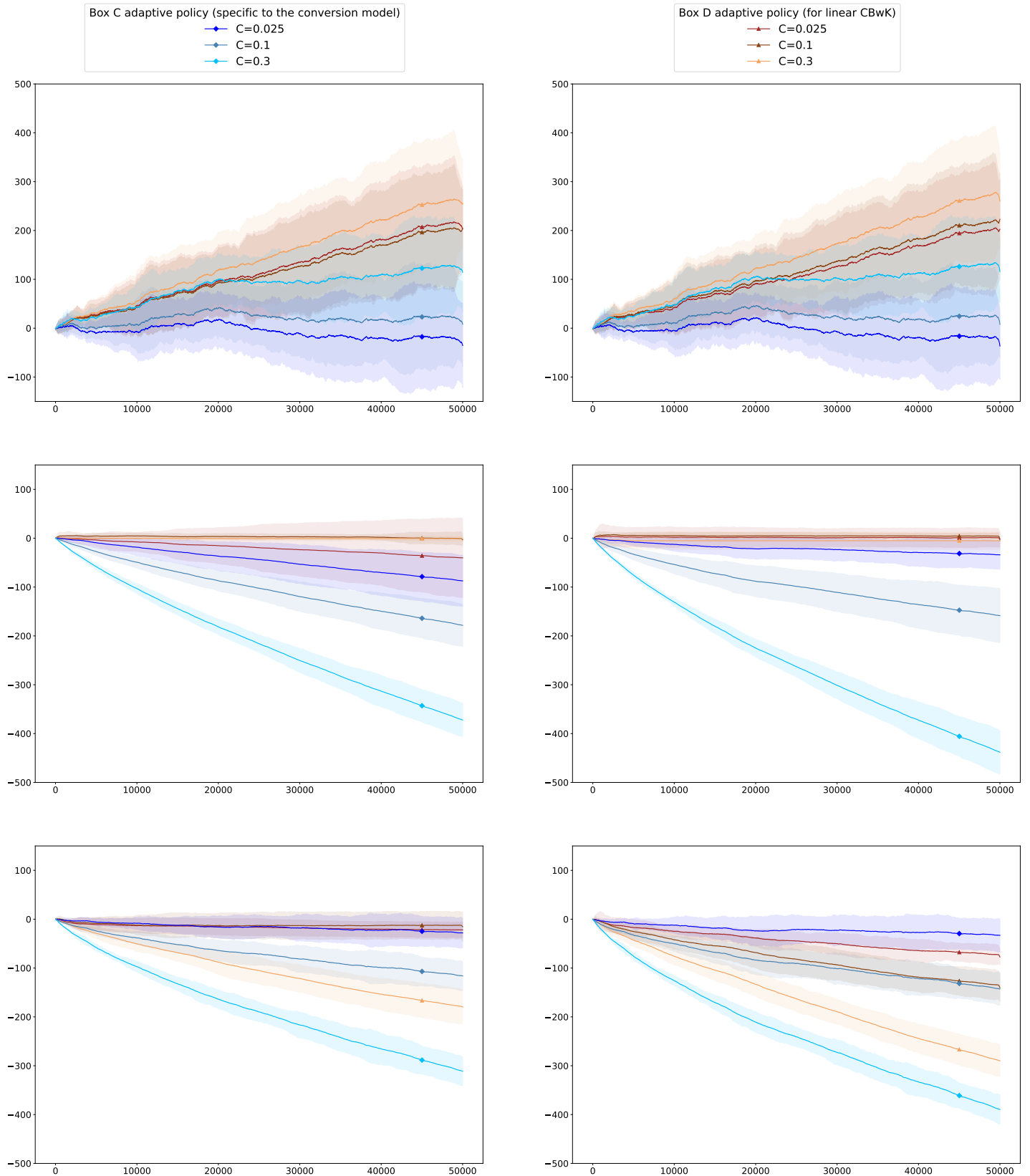


Figure 1: Averages (solid lines) and ± 2 times standard errors (shaded areas), achieved on 10 runs by the Box C (blue) and Box D (orange) adaptive policies: of the regret (first line), of the difference of the first cost component to tB/T (second line), and of the difference of the second cost component to tB/T (third line), by the values of the budget ($B = 1,600$ in the first column, $B = 2,200$ in the second column).