

Small Total-Cost Constraints in Contextual Bandits with Knapsacks [CBwK], with Application to Fairness

Evgenii Chzhen – Christophe Giraud – Gilles Stoltz



Zhen Li



BNP PARIBAS

Neurips 2023

CBwK framework – Novelty is $T\mathbf{B}$ with components of order \sqrt{T}

Known: finite \mathcal{A} , rounds T , average costs $\mathbf{B} \in [0, 1]^d$

Unknowns:

Context distribution ν on \mathcal{X}

Scalar mean-payoff function $r : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$

Vector-valued mean-cost function $\mathbf{c} : \mathcal{X} \times \mathcal{A} \rightarrow [-1, 1]^d$

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

Observe context $\mathbf{x}_t \sim \nu$, and pick $a_t \in \mathcal{A}$

Get payoff r_t and costs \mathbf{c}_t with cond. exp. $r(\mathbf{x}_t, a_t)$ and $\mathbf{c}(\mathbf{x}_t, a_t)$

Goals (cf. fairness costs: $T\mathbf{B}$ as small as possible, possibly \sqrt{T})

Ensure $\sum_{t=1}^T \mathbf{c}_t \leq T\mathbf{B}$ a.s. while maximizing $\sum_{t=1}^T r_t$

First reference for CBwK: Badanidiyuru, Langford, Slivkins [2014]

State of the art = $T\mathbf{B}$ at best $T^{3/4}$: Agrawal and Devanur [2016], Han et al. [2022]

Regret: Minimize $R_T = T \text{opt}(r, \mathbf{c}, \mathbf{B}) - \sum_{t=1}^T r_t$ where

$\text{opt}(r, \mathbf{c}, \mathbf{B})$

$$= \sup_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \left\{ \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(\mathbf{X}, a) \pi_a(\mathbf{X}) \right] : \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} \mathbf{c}(\mathbf{X}, a) \pi_a(\mathbf{X}) \right] \leq \mathbf{B} \right\}$$

$$= \sup_{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})} \inf_{\lambda \geq 0} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\sum_{a \in \mathcal{A}} r(\mathbf{X}, a) \pi_a(\mathbf{X}) + \left\langle \lambda, \mathbf{B} - \sum_{a \in \mathcal{A}} \mathbf{c}(\mathbf{X}, a) \pi_a(\mathbf{X}) \right\rangle \right]$$

$$= \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{X} \sim \nu} \left[\max_{a \in \mathcal{A}} \left\{ r(\mathbf{X}, a) - \langle \mathbf{c}(\mathbf{X}, a) - \mathbf{B}, \lambda \rangle \right\} \right]$$

→ Suffices to learn r and \mathbf{c} , as well as λ^*

Learn r and \mathbf{c} : via structural assumptions; uniform bounds

Linear model: Agrawal and Devanur [2016], based on LinUCB from Abbasi-Yadkori et al. [2011]. **Logistic model:** Li and Stoltz [2022], based on Logistic-UCB1 from Faury et al. [2020].

Target:
$$\text{opt}(r, \mathbf{c}, \mathbf{B}) = \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{x} \sim \nu} \left[\max_{a \in \mathcal{A}} \left\{ r(\mathbf{X}, a) - \langle \mathbf{c}(\mathbf{X}, a) - \mathbf{B}, \lambda \rangle \right\} \right]$$

→ Gradient descent on dual / best response for primal var.

Algorithm with fixed step size γ

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

1. Play $a_t \in \arg \max_{a \in \mathcal{A}} \left\{ \hat{r}_{t-1}(\mathbf{x}_t, a) - \langle \hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - b\mathbf{1}), \lambda_{t-1} \rangle \right\}$
2. Make gradient step $\lambda_t = \left(\lambda_{t-1} + \gamma (\hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - b\mathbf{1})) \right)_+$
3. Update estimates \hat{r}_t and $\hat{\mathbf{c}}_t$ of functions r and \mathbf{c}

Analysis

Cost margin Tb , of order $(1 + \|\lambda^*\|)/\gamma$; adds $\|\lambda^*\| (Tb + \sqrt{T})$ to regret

→ Oracle choice $(1 + \|\lambda^*\|)/\sqrt{T}$ for γ , leads to $(1 + \|\lambda^*\|)\sqrt{T}$ regret

Reminder of the issue: oracle choice $(1 + \|\lambda^*\|)/\sqrt{T}$ for γ

Typical bypass by estimating $\|\lambda^*\|$ on \sqrt{T} preliminary rounds (see, e.g.: Agrawal and Devanur [2016], Han et al. [2022]) leads to $\min \mathbf{B} \geq T^{-1/4}$

Theorem

Algorithm based on a *careful doubling trick* $\gamma_k = 2^k/\sqrt{T}$

W.h.p.: controls cumulative costs & regret of order $(1 + \|\lambda^*\|)\sqrt{T}$

Only requires $\min \mathbf{B}$ to be larger than $1/\sqrt{T}$ up to poly-log terms

Note 1: if null-cost action, $\|\lambda^*\| \leq \frac{2 \text{opt}(r, \mathbf{c}, \mathbf{B})}{\min \mathbf{B}}$ = usual bound

Note 2: explicit, closed-form bounds in the article

Note 3: fairness example from Chohlas-Wood et al. [2021]