

Contextual Stochastic Bandits with Budget Constraints and Fairness Application

Gilles Stoltz

Laboratoire de mathématiques d'Orsay



Joint work with **Evgenii Chzhen**, **Christophe Giraud**, and **Zhen Li**

Contextual Stochastic Bandits with Budget Constraints and Fairness Application

- 1 Stochastic bandits
- 2 **Contextual** stochastic bandits
- 3 Contextual stochastic bandits **with budget constraints**
 - Application to fairness: small budgets

K -armed stochastic bandits

Simplest possible framework

K probability distributions ν_1, \dots, ν_K in a model \mathcal{D}
with expectations μ_1, \dots, μ_K $\longrightarrow \mu^* = \max_{a \in [K]} \mu_a$

At each round $t = 1, 2, \dots$,

1. Statistician picks arm $A_t \in [K]$
2. She gets a reward Y_t drawn according to ν_{A_t}
3. This is the **only feedback** she receives

\longrightarrow **Exploration–exploitation dilemma**
estimate the ν_a **vs.** get high rewards Y_t

Goal:

Maximize expected cumulative rewards \longleftrightarrow Minimize **regret**

$$R_T = T\mu^* - \mathbb{E} \left[\sum_{t=1}^T Y_t \right] = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E} [N_a(T)]$$

\longleftrightarrow Control the $\mathbb{E} [N_a(T)]$ for suboptimal arms a

Setting:

Distributions ν_1, \dots, ν_K with expectations μ_1, \dots, μ_K

At each round $t \geq 1$, pick arm $A_t \in [K]$, get and observe $Y_t \sim \nu_{A_t}$

Proof of the rewriting of regret

Tower rule: $\mathbb{E}[Y_t | A_t] = \mu_{A_t}$ thus $\mathbb{E}[Y_t] = \mathbb{E}[\mu_{A_t}]$

$$\begin{aligned} R_T &= \sum_{t=1}^T (\mu^* - \mathbb{E}[Y_t]) = \sum_{t=1}^T (\mu^* - \mathbb{E}[\mu_{A_t}]) \\ &= \sum_{t=1}^T \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}[\mathbb{I}_{\{A_t=a\}}] = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}[N_a(T)] \end{aligned}$$

where $N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$

Model: ν_1, \dots, ν_K are distributions over $[0, 1]$

A popular strategy: **UCB** [upper confidence bound]

Auer, Cesa-Bianchi and Fisher [2002]

For $t \geq K$, pick $A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + \sqrt{\frac{2 \ln t}{N_a(t)}} \right\}$

Exploitation: cf. empirical mean $\hat{\mu}_a(t) = \frac{1}{N_a(t)} \sum_{s=1}^t Y_s \mathbb{I}_{\{A_s=a\}}$

Exploration: cf. $\sqrt{2 \ln t / N_a(t)}$ favors arms a not pulled often

Regret bounds (suboptimal) of two types

– **Distribution-dependent** bound: $R_T \lesssim \sum_{a: \mu_a < \mu^*} \frac{8 \ln T}{\mu^* - \mu_a}$

– **Distribution-free** bound: $\sup_{\nu_1, \dots, \nu_K} R_T \lesssim \sqrt{8KT \ln T}$

$$\text{Proof of } R_T \lesssim \sum_{a: \mu_a < \mu^*} \frac{8 \ln T}{\mu^* - \mu_a}$$

$$\text{Hoeffding-Azuma: } \mathbb{P} \left\{ |\mu_a - \hat{\mu}_a(t)| \leq \sqrt{\frac{2 \ln t}{N_a(t)}} \right\} \geq 1 - 2t^{-3}$$

Indeed, by optional skipping:

$$\begin{aligned} & \mathbb{P} \left\{ |\mu_a - \hat{\mu}_a(t)| > \sqrt{\frac{2 \ln t}{N_a(t)}} \right\} \\ &= \sum_{n=1}^t \mathbb{P} \left\{ |\mu_a - \hat{\mu}_{a,n}| > \sqrt{\frac{2 \ln t}{n}} \text{ and } N_a(t) = n \right\} \\ &\leq \underbrace{\sum_{n=1}^t \mathbb{P} \left\{ |\mu_a - \hat{\mu}_{a,n}| > \sqrt{\frac{\ln(1/t^{-4})}{2n}} \right\}}_{\leq 2t^{-4}} \end{aligned}$$

where $\hat{\mu}_{a,n}$ denotes the average of n -sample with distribution ν_a

Proof of $R_T \lesssim \sum_{a:\mu_a < \mu^*} \frac{8 \ln T}{\mu^* - \mu_a}$

Hoeffding–Azuma: $\mathbb{P} \left\{ \left| \mu_a - \hat{\mu}_a(t) \right| \leq \sqrt{\frac{2 \ln t}{N_a(t)}} \right\} \geq 1 - 2t^{-3}$

If $A_t = b$ is not an optimal arm a^* , then by design

$$\hat{\mu}_{a^*}(t) + \sqrt{\frac{2 \ln t}{N_{a^*}(t)}} \leq \hat{\mu}_b(t) + \sqrt{\frac{2 \ln t}{N_b(t)}}$$

thus w.h.p. $\mu^* \leq \mu_b + 2\sqrt{\frac{2 \ln t}{N_b(t)}}$

which imposes $N_b(t) \leq \frac{8 \ln T}{(\mu^* - \mu_a)^2}$

Conclude with $R_T = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}[N_a(T)]$

Proof of $\sup_{\nu_1, \dots, \nu_K} R_T \lesssim \sqrt{8KT \ln T}$

We proved $\mathbb{E}[N_b(t)] \lesssim \frac{8 \ln T}{(\mu^* - \mu_a)^2}$

Thus

$$\begin{aligned} R_T &= \sum_{a \in [K]} (\mu^* - \mu_a) \sqrt{\mathbb{E}[N_a(T)]} \sqrt{\mathbb{E}[N_a(T)]} \\ &\leq \sqrt{8 \ln T} \sum_{a \in [K]} \sqrt{\mathbb{E}[N_a(T)]} \\ &\leq \sqrt{8KT \ln T} \end{aligned}$$

Contextual stochastic bandits with K arms

Linear modeling + Logistic modeling

At each round $t = 1, 2, \dots$,

0. A **context** $\mathbf{x}_t \in \mathbb{R}^d$ is **determined** by the environment
1. Statistician picks arm $A_t \in [K]$
2. She gets a reward Y_t with conditional expectation $r(\mathbf{x}_t, A_t)$
3. This is the only feedback she receives

Goal:

Maximize expected rewards \longleftrightarrow Minimize **expected regret**

$$R_T = \sum_{t=1}^T \text{targets?} - \mathbb{E} \left[\sum_{t=1}^T Y_t \right]$$

Structural assumptions handy! E.g., **linearity**:

$$r(\mathbf{x}, a) = \varphi(\mathbf{x}, a)^\top \theta_\star \quad \rightsquigarrow \quad \text{targets} \quad \max_{a \in [K]} \varphi(\mathbf{x}_t, a)^\top \theta_\star$$

Transfer function $\varphi : \mathbb{R}^d \times [K] \rightarrow \mathbb{R}^m$ known,
But parameters $\theta_\star \in \mathbb{R}^d$ unknown

Setting: contexts $\mathbf{x}_t \in \mathbb{R}^d$, pick arms $A_t \in [K]$, get rewards Y_t

Regret
$$R_T = \sum_{t \leq T} \max_{a \in [K]} \varphi(\mathbf{x}_t, a)^\top \theta_\star - \sum_{t \leq T} \mathbb{E}[\varphi(\mathbf{x}_t, A_t)^\top \theta_\star]$$

Key: learn θ_\star (= estimate it while playing)

LinUCB with regularization $\lambda > 0$ for bounded contexts

Abbasi-Yadkori, Pál, Szepesvári [2011]

Based on the idea
$$\sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) Y_s \approx \sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) \varphi(\mathbf{x}_s, A_s)^\top \theta_\star$$

Statement: let
$$M_{t-1} = \lambda \text{Id} + \sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) \varphi(\mathbf{x}_s, A_s)^\top$$

and
$$\hat{\theta}_{t-1} = (M_{t-1})^{-1} \sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) Y_s$$

Setting: bounded contexts $\mathbf{x}_t \in \mathbb{R}^d$, arms $A_t \in [K]$, rewards Y_t

reward function $r(\mathbf{x}, a) = \varphi(\mathbf{x}, a)^\top \theta_*$, with $\mathbb{E}[Y_t | A_t, \mathbf{x}_t] = \varphi(\mathbf{x}_t, A_t)^\top \theta_*$

$$\hat{\theta}_{t-1} = (M_{t-1})^{-1} \sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) Y_s \quad \text{where} \quad M_{t-1} = \lambda \text{Id} + \sum_{s=1}^{t-1} \varphi(\mathbf{x}_s, A_s) \varphi(\mathbf{x}_s, A_s)^\top$$

Confidence region on θ_* :

$$\mathbb{P} \left\{ \left\| \hat{\theta}_{t-1} - \theta_* \right\|_{M_{t-1}} \lesssim \square \sqrt{\ln(t/\delta)} \right\} = 1 - \delta$$

where $\|u\|_M = \sqrt{u^\top M u}$ and provided that λ is well set

Complex proof based on “Laplace’s method of mixtures”

Simultaneous confidence intervals on the $r(\mathbf{x}, a)$: based on

$$\begin{aligned} \left| \varphi(\mathbf{x}, a)^\top \hat{\theta}_{t-1} - \varphi(\mathbf{x}, a)^\top \theta_* \right| &\leq \left\| \hat{\theta}_{t-1} - \theta_* \right\|_{M_{t-1}} \left\| \varphi(\mathbf{x}, a) \right\|_{(M_{t-1})^{-1}} \\ &\leq \underbrace{\square \sqrt{\ln(t/\delta)} \left\| \varphi(\mathbf{x}, a) \right\|_{(M_{t-1})^{-1}}}_{= \varepsilon_{t-1, \delta}(\mathbf{x}, a)} \end{aligned}$$

where $\sum_{t=1}^T \varepsilon_{t-1, \delta}(\mathbf{x}_t, A_t) \lesssim \sqrt{T} \ln(T/\delta)$ by linear algebra

Setting: bounded contexts $\mathbf{x}_t \in \mathbb{R}^d$, arms $A_t \in [K]$, rewards Y_t
reward function $r(\mathbf{x}, a) = \varphi(\mathbf{x}, a)^\top \theta_*$, with $\mathbb{E}[Y_t | A_t, \mathbf{x}_t] = \varphi(\mathbf{x}_t, A_t)^\top \theta_*$

Simultaneous confidence intervals: $|\hat{r}_{t-1}(\mathbf{x}, a) - r(\mathbf{x}, a)| \leq \varepsilon_{t-1, \delta}(\mathbf{x}, a)$
where $\sum_{t \leq T} \varepsilon_{t-1, \delta}(\mathbf{x}_t, A_t) \lesssim \sqrt{T} \ln(T/\delta)$

Optimistic choice: $A_t \in \arg \max_{a \in [K]} \{ \hat{r}_{t-1}(\mathbf{x}_t, a) + \varepsilon_{t-1, \delta}(\mathbf{x}_t, a) \}$

Regret bound: $R_T = \sum_{t=1}^T \max_{a \in [K]} r(\mathbf{x}_t, a) - \sum_{t=1}^T Y_t \leq \tilde{O}(\sqrt{T})$

In high-probability (but algorithm depends on δ)

Or in expectation (set $\delta = t^{-4}$, e.g.)

We could also have obtained high-probability bounds based on the UCB strategy in the non-contextual case

Logistic bandits

Extended from Fauray, Abeille, Calauzènes, Fercoq [2020]

At each round $t = 1, 2, \dots$,

0. A context $\mathbf{x}_t \in \mathbb{R}^d$ is determined by the environment
1. Statistician picks arm $A_t \in [K]$
2. The **outcome** $Y_t \in \{0, 1\}$ is drawn with probability $P(\mathbf{x}_t, A_t)$
3. This is the only feedback Statistician receives
4. Statistician gets the reward $r(\mathbf{x}_t, A_t) Y_t$

Conversion rate P unknown but reward function r known

Structural assumption:

$$P(\mathbf{x}, a) = \eta(\varphi(\mathbf{x}, a)^\top \theta_*) \quad \text{where} \quad \eta(x) = \frac{1}{1 + e^{-x}}$$

Similar results may be achieved as for linear bandits

Estimation based on maximum likelihood

Contextual stochastic bandits with K arms

And now, with budget constraints!

At each round $t = 1, 2, \dots$,

0. A context $\mathbf{x}_t \sim \mathbb{Q}$ is **drawn at random**
1. Statistician picks arm $A_t \in [K]$
2. She gets a reward Y_t with conditional expectation $r(\mathbf{x}_t, A_t)$
3. She also suffers costs \mathbf{Z}_t with conditional expectation $\mathbf{c}(\mathbf{x}_t, A_t)$
4. Her feedback is Y_t and \mathbf{Z}_t

Vector-valued costs: possibly several constraints

Goals:

Maximize $\sum_{t \leq T} Y_t$ while ensuring $\sum_{t \leq T} \mathbf{Z}_t \leq T\mathbf{B}$

Known: budget $T\mathbf{B}$

Unknown: reward function r , cost function \mathbf{c} , distribution \mathbb{Q}

but structural assumptions to be issued on r and \mathbf{c}

Setting called **CBwK** – contextual bandits with knapsacks

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

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

Fairness application

Inspired from Chohlas-Wood, Coots, Zhu, Brunskill, Goel [2021]

Fair budget spending among groups: Z'_t first component of \mathbf{Z}_t

$$\sum_{t=1}^T Z'_t \leq TB_{\text{total}}$$

$$\text{and } \forall g \in \mathcal{G}, \left| \frac{1}{T\gamma_g} \sum_{t=1}^T Z'_t \mathbb{I}_{\{\text{gr}(\mathbf{x}_t)=g\}} - \frac{1}{T} \sum_{t=1}^T Z'_t \right| \leq \tau$$

where $\gamma_g = \mathbb{Q}\{\text{gr}(\cdot) = g\}$

and τ is a tolerance factor, ideally $\sim 1/\sqrt{T}$, i.e., $T\tau$ of order \sqrt{T}

B contains a B_{total} component, as well as components $\pm\gamma_g\tau$

Setting: context $\mathbf{x}_t \sim \mathbb{Q}$, arm $A_t \in [K]$, reward Y_t and costs \mathbf{Z}_t

Conditional expectations: $r(\mathbf{x}_t, A_t)$ and $\mathbf{c}(\mathbf{x}_t, A_t)$

Total budget constraints $T\mathbf{B}$, where some components are as small as \sqrt{T}

Benchmark: static policies $\pi : \mathbf{x} \mapsto (\pi_a(\mathbf{x}))_{a \in [K]} \in \mathcal{P}([K])$

We assume feasibility, and actually do so for $\mathbf{B} - \varepsilon \mathbf{1}$ (OK if a null-cost action exists)

$$\text{opt}(r, \mathbf{c}, \mathbf{B}) = \sup_{\pi} \left\{ \mathbb{E}_{\mathbf{X} \sim \mathbb{Q}} \left[\sum_{a \in [K]} r(\mathbf{X}, a) \pi_a(\mathbf{X}) \right] \right. \\ \left. \text{under } \mathbb{E}_{\mathbf{X} \sim \mathbb{Q}} \left[\sum_{a \in [K]} \mathbf{c}(\mathbf{X}, a) \pi_a(\mathbf{X}) \right] \leq \mathbf{B} \right\}$$

Regret: $R_T = T \text{opt}(r, \mathbf{c}, \mathbf{B}) - \sum_{t \leq T} Y_t$

Hard constraint: $\sum_{t \leq T} \mathbf{Z}_t \leq T\mathbf{B}$

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

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

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

$$= \sup_{\pi} \inf_{\lambda \geq 0} \mathbb{E}_{\mathbf{X} \sim \mathbb{Q}} \left[\sum_{a \in [K]} r(\mathbf{X}, a) \pi_a(\mathbf{X}) - \left\langle \lambda, \sum_{a \in [K]} \mathbf{c}(\mathbf{X}, a) \pi_a(\mathbf{X}) - \mathbf{B} \right\rangle \right]$$

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

→ Suffices to learn r and \mathbf{c} , as well as λ^* \rightsquigarrow parametric problems[†]

Cf. $\mathbf{x}_t \sim \mathbb{Q}$ observed at each round

Learn r and \mathbf{c} : via [†]structural assumptions (linearity or logistic)

Uniform bounds available

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

→ Gradient descent on dual / best response for primal variable(s)

Algorithm with fixed step size γ

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

1. Play $A_t \in \arg \max_{a \in [K]} \left\{ \hat{r}_{t-1}(\mathbf{x}_t, a) - \langle \hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - b\mathbf{1}), \boldsymbol{\lambda}_{t-1} \rangle \right\}$
2. Make gradient step $\boldsymbol{\lambda}_t = \left(\boldsymbol{\lambda}_{t-1} + \gamma (\hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - b\mathbf{1})) \right)_+$
3. Update estimates \hat{r}_t and $\hat{\mathbf{c}}_t$ of functions r and \mathbf{c}

Optimistic estimates: \hat{r}_t upper bounds r and $\hat{\mathbf{c}}_t$ lower bounds \mathbf{c}

Idea already in Agrawal and Devanur [2016]

But the key to handle **smaller budgets** is the **tuning of γ**

From Chzhen, Giraud, Li, Stoltz [2023]

Strategy with fixed γ

1. Play $A_t \in \arg \max_{a \in [K]} \left\{ \hat{r}_{t-1}(\mathbf{x}_t, a) - \langle \hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - \mathbf{b}\mathbf{1}), \boldsymbol{\lambda}_{t-1} \rangle \right\}$
2. Make gradient step $\boldsymbol{\lambda}_t = \left(\boldsymbol{\lambda}_{t-1} + \gamma (\hat{\mathbf{c}}_{t-1}(\mathbf{x}_t, a) - (\mathbf{B} - \mathbf{b}\mathbf{1})) \right)_+$

Analysis for fixed γ : the projected-gradient descent entails

$$\left\| \left(\sum_{t=1}^T \mathbf{z}_t - T(\mathbf{B} - \mathbf{b}\mathbf{1}) \right)_+ \right\| \leq \tilde{\mathcal{O}} \left(\frac{1 + \|\boldsymbol{\lambda}^*\|}{\gamma} \right)$$

Cost margin Tb should be of the same order $(1 + \|\boldsymbol{\lambda}^*\|)/\gamma$ That margin adds a term of order $\|\boldsymbol{\lambda}^*\|(Tb + \sqrt{T})$ to regret→ Oracle choices $b \sim 1/\sqrt{T}$ and $\gamma \sim (1 + \|\boldsymbol{\lambda}^*\|)/\sqrt{T}$ lead to $(1 + \|\boldsymbol{\lambda}^*\|)\sqrt{T}$ regretEstimating $\|\boldsymbol{\lambda}^*\|$ on \sqrt{T} exploration rounds (see, e.g.: Agrawal and Devanur [2016], Han et al. [2022]) imposes $\min \mathbf{B} \geq T^{-1/4}$

We use instead a **careful doubling trick** $\gamma_k = 2^k / \sqrt{T}$

Breaking condition based on budget controls

Theorem:

If $\min \mathbf{B}$ is larger than $1/\sqrt{T}$ up to poly-log terms, then w.h.p.,

$$\sum_{t \leq T} \mathbf{z}_t \leq T\mathbf{B} \quad \text{and} \quad R_T \lesssim \tilde{O}(1 + \|\lambda^*\|)\sqrt{T}$$

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