

Meritocratically Fair Algorithms for Infinite and Contextual Bandits

M. Joseph, M. Kearns, J. Morgenstern, S. Neel, A. Roth

High-Level Motivation

- ▶ Machine learning can be *unfair* in many ways: biased data; different populations with different properties; less data about minorities, etc.
- ▶ How do we define *fair learning*? What is the performance cost of fairness?

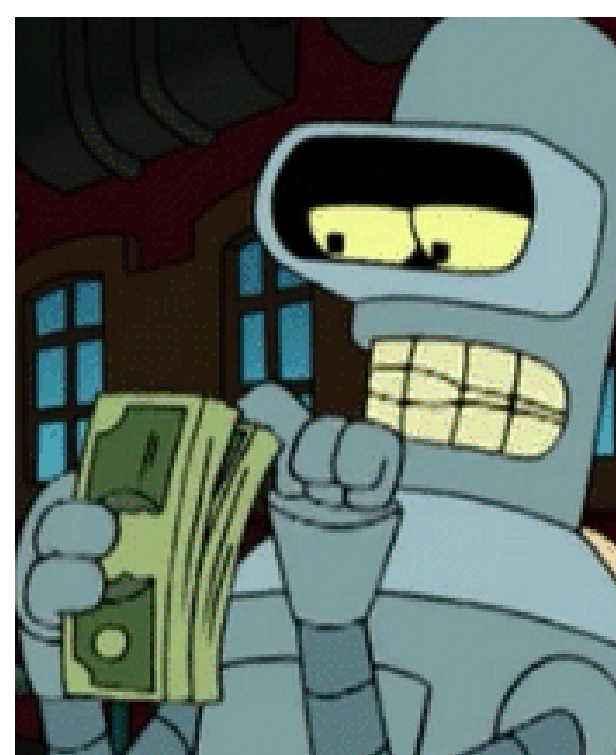
Previous Work

- ▶ JKMR16 [1] studied fairness for finite contextual bandits
- ▶ **Problem 1:** unrealistic assumptions (one individual per group per day; choose exactly one individual per day – artificial inter-group competition?)
- ▶ **Problem 2:** results do not scale well when number of arms is large

Finite Setting

- ▶ Goal: address Problem 1 above
- ▶ In each round \mathbf{t} we see set \mathbf{C}_t of at most k contexts in \mathbb{R}^d , choose a subset $\mathbf{P}_t \subset \mathbf{C}_t$ of exactly m contexts, and observe noisy linear reward $\mathbf{r}_i^t = \langle \beta, \mathbf{x}_i^t \rangle + \epsilon_i^t$ for each $i \in \mathbf{P}_t$
- ▶ Addresses problem 1: can see multiple individuals per population per round, can choose multiple individuals per round
- ▶ Group membership can be encoded in context in \mathbb{R}^d or not
- ▶ Goal: maximize $\sum_t \sum_{i \in \mathbf{P}_t} \mathbb{E}[r_i^t]$, measure performance by regret $\mathbf{R}(\mathbf{T}) = \sum_t [\mathbb{E}[\sum_{i \in \mathbf{P}_t^*} r_i^t] - \mathbb{E}[\sum_{j \in \mathbf{P}_t} r_j^t]]$ (loss from choosing subset \mathbf{P}_t instead of best expected subset \mathbf{P}_t^* across rounds)

- ▶ Models a program that learns to grant loans by granting m loans daily



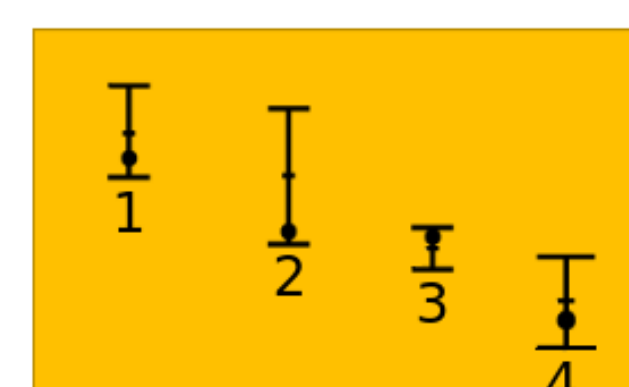
General Fairness Definition

- ▶ Algorithm \mathcal{A} is **fair** if (whp) for all $\mathbf{t} \in \mathbf{T}$ and for all $i, j \in \mathbf{C}_t$ $\mathbb{E}[r_i^t] \geq \mathbb{E}[r_j^t] \Rightarrow \pi_i^t \geq \pi_j^t$ where $\pi_i^t = \mathbb{P}[\text{choose } i \text{ in round } \mathbf{t}]$ (omitting histories for simple notation)

Interpreting Fairness Definition

- ▶ “Whp, never more likely to choose a worse arm than a better arm”
- ▶ Optimal policies always play the expected best arm and therefore are fair. Challenge: how to *learn* the optimal policy fairly?

A Fair Algorithm: RidgeFair_m



• μ
• $\hat{\mu}$
• CI bounds

RidgeFair₅ would select arms 1, 2, 3, and 4, and also select one of 5, 6, and 7 at random.

- ▶ Uses confidence intervals around estimated means to reason about relative quality; fairness forces *chaining*

- ▶ In round \mathbf{t} : Choose all arms in highest connected component of confidence intervals, then choose the last arms by randomizing when you reach a connected component in which you cannot choose all arms

FairUCB [1] vs RidgeFair_m

- ▶ Encoded in our setting (using contexts in \mathbb{R}^{dk}) FairUCB achieves regret $\mathbf{R}(\mathbf{T}) = \mathbf{O}(\max[\mathbf{T}^{4/5} k^{6/5} d^{3/5}, k^3])$
- ▶ RidgeFair_m achieves regret $\mathbf{R}(\mathbf{T}) = \mathbf{O}(dk^2 \sqrt{\mathbf{T}})$
- ▶ Improvement via better (and more technical) confidence intervals for $\hat{\beta}$
 - ▶ Uses martingale matrix concentration results from APS11 [2]

Infinite Setting

- ▶ Goal: address Problem 2 above
- ▶ In each round \mathbf{t} we see a convex set \mathbf{C}_t of choices contained in a ball of radius r , select exactly one, and observe (single) noisy reward $\mathbf{r}_t = \langle \beta, \mathbf{x} \rangle + \epsilon_t$
- ▶ Goal: maximize $\sum_t \mathbb{E}[\mathbf{r}_t]$, measure performance by regret $\mathbf{R}(\mathbf{T}) = \sum_t \mathbb{E}[\mathbf{r}_t^* - \mathbf{r}_t]$ where \mathbf{r}_t^* is an optimal choice in round \mathbf{t} and \mathbf{r}_t is the actual choice

A Fair Algorithm: FairGap

- ▶ Uses convexity of each \mathbf{C}_t : optimal point must be an *extremal* point
- ▶ Plays randomly until confidence interval around $\hat{\beta}$ shrinks enough to separate optimal extremal point from suboptimal extremal points
- ▶ Performance thus depends on Δ_{gap} – the “gap” in expected reward between an optimal and next sub-optimal extremal point
- ▶ Instance-dependent regret bound:

$$\mathbf{R}(\mathbf{T}) = \tilde{\mathbf{O}} \left(\frac{r^6 \mathbf{R}^2}{\kappa^2 \lambda^2 \Delta_{\text{gap}}^2} \right)$$
 where $\kappa = 1 - r \sqrt{\frac{2}{\mathbf{T} \lambda}}$ and $\lambda = \min_{1 \leq t \leq \mathbf{T}} [\lambda_{\min}(\mathbb{E}_{\mathbf{x}_t \sim \mathbf{C}_t} [\mathbf{x}_t^T \mathbf{x}_t])]$
- ▶ Regret independent of k

Instance-dependent Lower Bound

- ▶ **Thm:** Let $\mathbf{C}_t = [-1, 1]^d$ for each \mathbf{t} and choose some $\beta \in [-1, 1]^d$. Then for every Δ_{gap} there exists an instance distribution for which any fair algorithm whp experiences $\tilde{\Omega}(1/\Delta_{\text{gap}})$ regret.
 - ▶ Adapts Bayesian lower-bound argument from JKMR16 [1]
- ▶ (Some) instance-dependence is therefore necessary for any fair algorithm in the infinite setting
 - ▶ FairGap’s $\mathbf{O}(1/\Delta_{\text{gap}}^2)$ regret is almost tight

Instance-independent Lower Bound

- ▶ **Thm:** Let \mathbf{C}_t be \mathbf{S}^1 (the unit circle) for each \mathbf{t} . Then for any $\beta \in \mathbf{S}^1$, no fair algorithm achieves nontrivial regret.
- ▶ Consequence of $\Delta_{\text{gap}} = 0$ – continuity means you can never actually identify an optimal point

References

- [1] Matthew Joseph, Michael Kearns, Jamie H Morgenstern, and Aaron Roth. Fairness in learning: Classic and contextual bandits. In *NIPS 2016*.
- [2] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *NIPS 2011*.