

Fairness in Learning: Classic and Contextual Bandits

M. Joseph, M. Kearns, J. Morgenstern, A. Roth
 {majos, mkearns, jamiemor, aaroth}@cis.upenn.edu

High-Level Motivation

- ▶ Machine learning can be unfair in many ways: data that encodes existing biases; data collection feedback loops; different populations having different properties; less data about minority populations . . .
- ▶ How do we define “fair learning”?
- ▶ What is the performance cost of being fair?

General Problem Setting

- ▶ We study the *bandits* setting: k arms, on day $t \in T$ choose arm i^t and observe noisy reward $r_{i^t}^t$
- ▶ Goal: maximize $\sum_t \mathbb{E}[r_{i^t}^t]$, measure performance by regret $R(T) = \sum_t [\max_{i \in [k]} \mathbb{E}[r_i^t] - r_{i^t}^t]$

- ▶ Models a program that learns to grant loans to k different groups by granting loans to one member of one group each day



General Fairness Definition

- ▶ Algorithm \mathcal{A} is **fair** if with probability $\geq 1 - \delta$, for all days $t \in T$ and for all $i, j \in [k]$

$$\mathbb{E}[r_i^t] \geq \mathbb{E}[r_j^t] \Rightarrow \pi_{i|h_1, \dots, h_{t-1}}^t \geq \pi_{j|h_1, \dots, h_{t-1}}^t$$

where $\pi_{i|h_1, \dots, h_{t-1}}^t =$

$\mathbb{P}[\text{choose } i \text{ in round } t \text{ after observing } h_1, \dots, h_{t-1}]$.

- ▶ “With high probability, never more likely to choose a worse arm than a better arm”

Why is Fairness Hard?

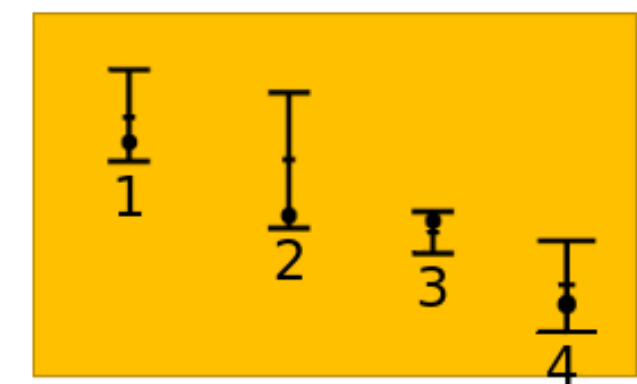
- ▶ Optimal policies always play the expected best arm and therefore are fair. Challenge: how to *learn* the optimal policy fairly?

Classic Bandits Setting

- ▶ μ_i for each arm i such that for all i and t $\mathbb{E}[r_i^t] = \mu_i$
- ▶ Fair: $\mu_i \geq \mu_j \Rightarrow \pi_{i|h_1, \dots, h_{t-1}}^t \geq \pi_{j|h_1, \dots, h_{t-1}}^t$
- ▶ “With high probability, never more likely to choose an arm with lower μ than an arm with higher μ ”

A Fair Classic Bandit Algorithm: FairBandits

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



• μ
 - $\hat{\mu}$
 - CI bounds

FairBandits plays randomly from chain (Arms 1 to 4)

- ▶ In round t : pick uniformly at random from “chain” of top arms (top connected component of overlapping confidence intervals)

Cost of Fairness in Classic Bandits

- ▶ FairBandits regret upper bound $R(T) = \tilde{O}(\sqrt{k^3 T})$
- ▶ Regret lower bound (any fair algorithm) $R(T) = \Omega(k^3)$, while $R(T) = \tilde{O}(\sqrt{k T})$ (unfair)

Contextual Bandits Setting

- ▶ Function $f_i \in \mathcal{C}$ for $i \in [k]$; $x_t^t \in \mathbb{R}^d$ for $t \in T$, $i \in [k]$ such that $\mathbb{E}[r_i^t] = f_i(x_t^t)$
- ▶ Fair: $f_i(x_t^t) \geq f_j(x_t^t) \Rightarrow \pi_{i|h_1, \dots, h_{t-1}}^t \geq \pi_{j|h_1, \dots, h_{t-1}}^t$
- ▶ “With high probability, never more likely to choose an arm with lower $f(x^t)$ than an arm with higher $f(x^t)$ ”

Fair Contextual Bandits and KWIK Learning

- ▶ \mathcal{C} is KWIK-learnable [1] with poly KWIK bound $\Leftrightarrow \mathcal{C}$ can be learned fairly with poly regret
- ▶ For d -dimensional *linear functions*, KWIK bounds [2] imply fair learning with $R(T) = \tilde{O}(\max(T^{4/5} k^{6/5} d^{3/5}, k^3))$
- ▶ For d -dimensional *conjunctions*, KWIK bounds [3] imply that no fair learning algorithm has a worst-case regret bound better than $R(T) = \Omega(2^d)$

References

- [1] Lihong Li, Michael L Littman, Thomas J Walsh, and Alexander L Strehl. Knows what it knows: a framework for self-aware learning. *Machine learning*, 82(3):399–443, 2011.
- [2] Alexander L Strehl and Michael L Littman. Online linear regression and its application to model-based reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 1417–1424, 2008.
- [3] Lihong Li. *A unifying framework for computational reinforcement learning theory*. PhD thesis, Rutgers, The State University of New Jersey, 2009.