

RUTGERS



# Towards Long-term Fairness in Recommendation

Presented by Yingqiang Ge

Rutgers University

3/3/2021



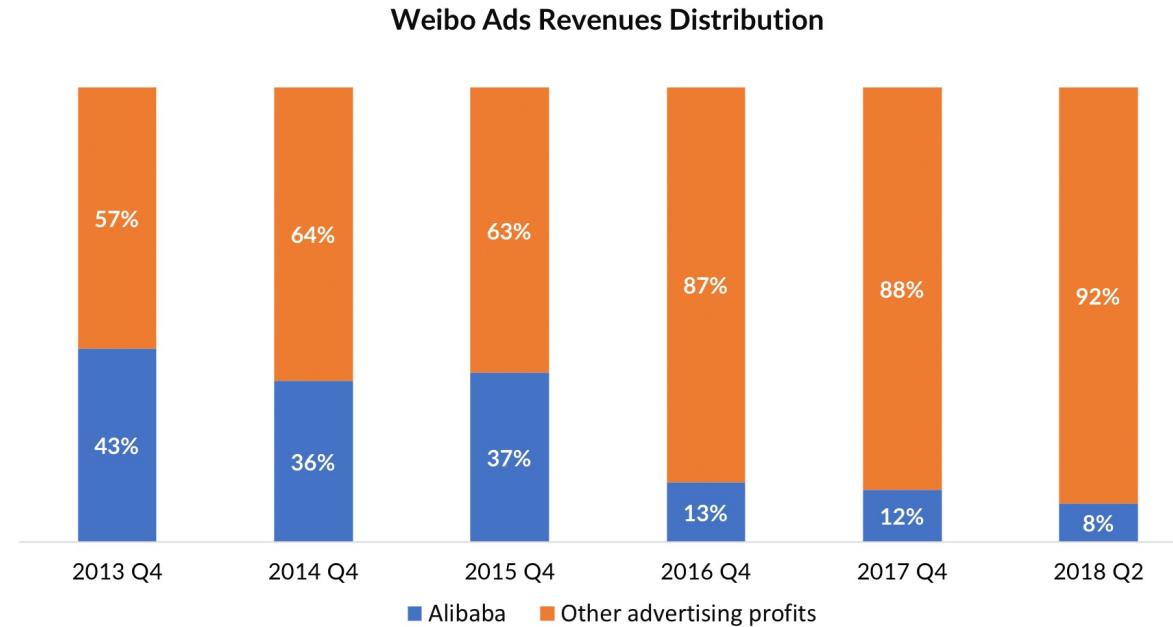
# Agenda

- ❖ Motivation
- ❖ Problem Formulation
- ❖ Model Specification
- ❖ Experiments
- ❖ Contributions
- ❖ Q&A



# Motivation

## ➤ Why Fairness in E-commerce



Source: Tech.sina.com, Walkthechat Analysis

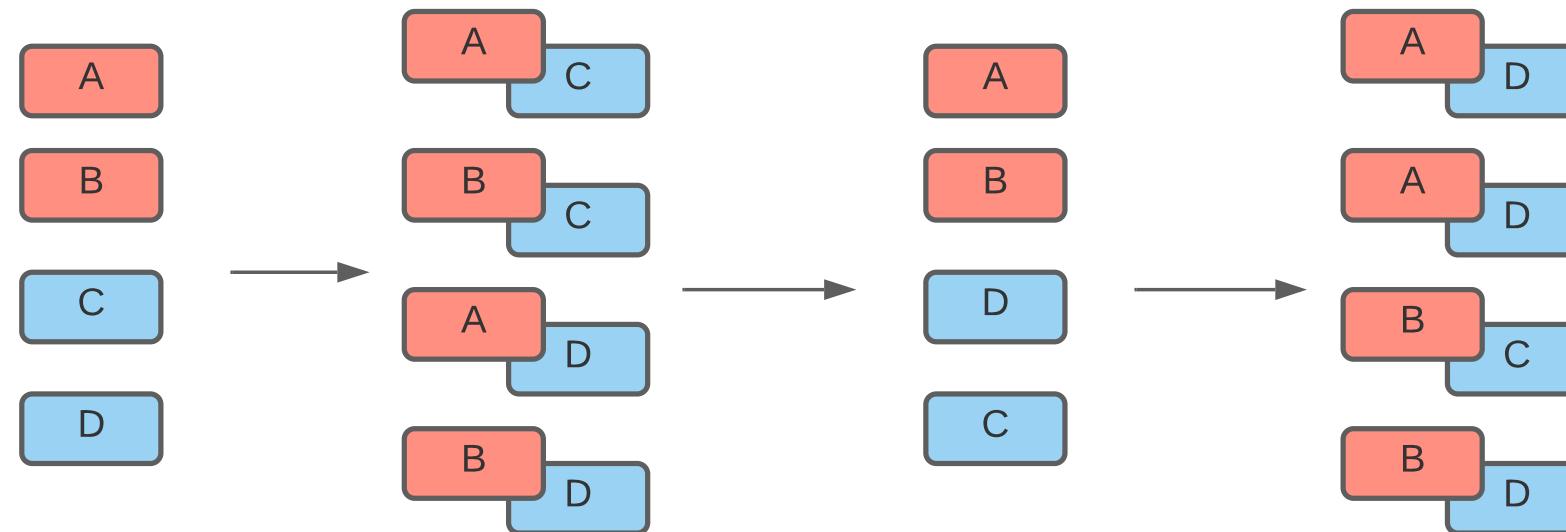
Once the Superstar left the platform(i.e., move to your competitor), the whole ecosystem of the e-commerce platform would crush immediately.

# Motivation

- Existing approaches fail to consider the dynamic nature of the recommender systems, where attributes such as item popularity may change over time
- Imposing seemingly fair decisions through static criteria can lead to unexpected unfairness in the long run.
- In essence, fairness cannot be defined in static or one-shot setting without considering the long-term impact, and long-term fairness cannot be achieved without understanding the underlying dynamics.

# Motivation

➤ Consider a simple example of ``Matthew Effect'' in a RS



# Motivation

- **Static fairness** does not consider the changes in the recommendation environment.
- **Dynamic fairness** learns a strategy that accommodates the dynamics.
- **Long-term fairness** views the recommendation as a long-term process and aims to maintain fairness in the long run by achieving dynamic fairness over time.

Need a fairness-aware system that

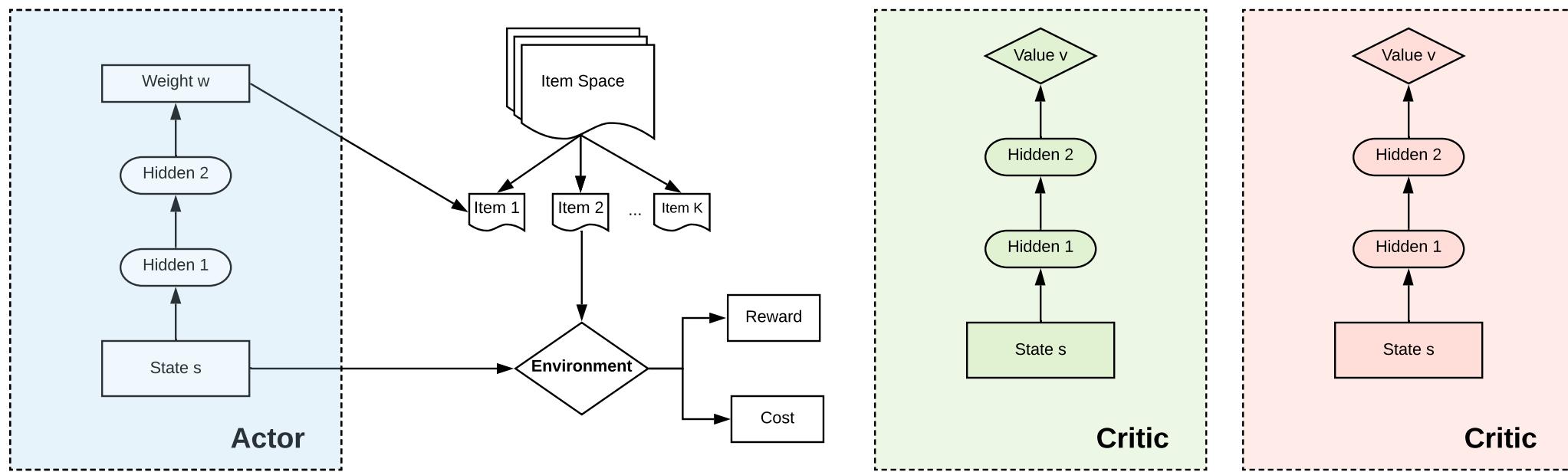
- is dynamic > static.
- considers long-term impact > gives one-time solution.

# Problem Formulation

In order to capture the interactive nature of recommendation scenarios, reinforcement learning (RL) based solutions have become an important topic.

- **State:** state  $s_t$  of a user
  - $H_t$  - user's most recent positive interaction history
  - Demographic information (if exists).
- **Action:** a recommendation list  $a_t = \{a_t^1, \dots, a_t^K\}$  with current state  $s_t$ .
- **Reward:** the immediate feedback  $R(st, at)$  given the action  $a_t$  and the user state  $s_t$ 
  - Typical user feedback includes click, skip, or purchase, etc.
- **Cost:** a cost value  $C(s_t, a_t)$  given by the problem-specific cost function
  - i.e., #items that come from the sensitive group
- **Discount rate:**  $\gamma_r$  and  $\gamma_c$ :
  - $\gamma_r \in [0,1]$  is for long-term rewards
  - $\gamma_c \in [0,1]$  is for long-term costs.

# Problem Formulation



# Problem Formulation

$$\pi_{k+1} = \arg \max_{\pi \in \Pi_\theta} \mathbb{E}_{\substack{s \sim d^{\pi_k} \\ a \sim \pi}} [A^{\pi_k}(s, a)],$$

$$\text{s.t. } J_{C_i}(\pi_k) + \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d^{\pi_k} \\ a \sim \pi}} \left[ A_{C_i}^{\pi_k}(s, a) \right] \leq \mathbf{d}_i, \forall i$$

$$\bar{D}_{KL}(\pi \| \pi_k) \leq \delta$$

Particularly, for problems with only one linear constraint, there is an analytical solution, which is also given by Achiam et al.

$$\theta_{k+1} = \arg \max_{\theta} g^\top (\theta - \theta_k)$$

$$\text{s.t. } c + b^\top (\theta - \theta_k) \leq 0$$

$$\frac{1}{2} (\theta - \theta_k)^\top H (\theta - \theta_k) \leq \delta$$

# Fairness Constraints

$$R(s_t, a_t, s_{t+1}) = \sum_{l=1}^K \mathbb{1}(a_t^l \text{ gets positive feedback})$$

$$C(s_t, a_t, s_{t+1}) = \sum_{l=1}^K \mathbb{1}(a_t^l \text{ is in sensitive group})$$

## Exact-K Fairness Constraint

$$\frac{\text{Exposure}_t(G_0)}{\text{Exposure}_t(G_1)} \leq \alpha$$

$$\text{Exposure}_t(G_0) \leq \alpha \text{Exposure}_t(G_1)$$

$$(1 + \alpha) \text{Exposure}_t(G_0) \leq \alpha \text{Exposure}_t(G_0) + \alpha \text{Exposure}_t(G_1)$$

$$(1 + \alpha) \text{Exposure}_t(G_0) \leq \alpha K$$

$$C(s_t, a_t, s_{t+1}) \leq \frac{\alpha}{1 + \alpha} K = \alpha' K$$

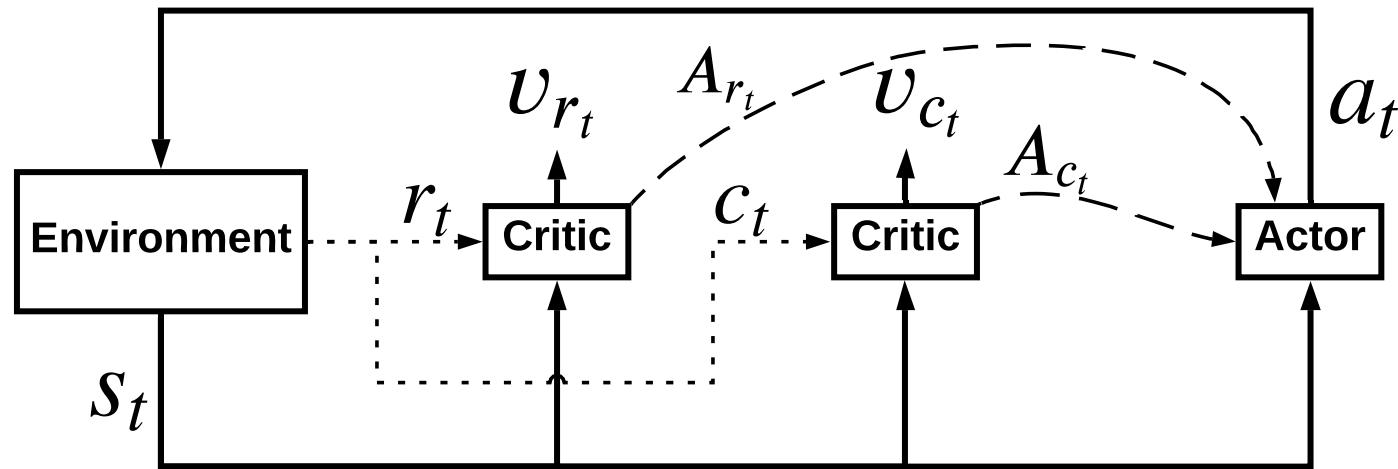
# Fairness Constraints

Goal: use CPO to learn a policy that maximize the cumulative reward while the cost is within the bound.

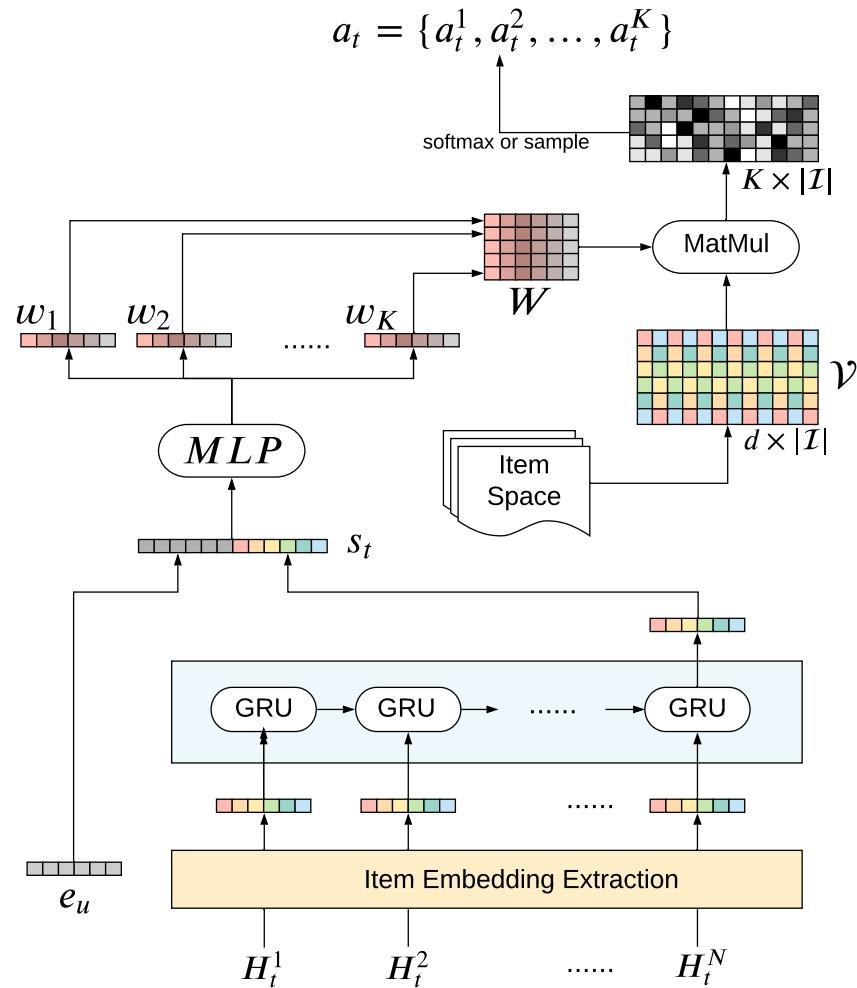
$$J_C(\pi) = \underset{\tau \sim \pi}{\mathbb{E}} \left[ \sum_{t=0}^T \gamma_c^t C(s_t, a_t, s_{t+1}) \right] \leq \sum_{t=0}^T \gamma_c^t \boxed{\alpha' K}$$

Ratio bound between two groups as in exact-K fairness constraints

# Model Specification



# Model Specification



$$s_t = [\mathbf{e}_u; \mathbf{h}_u], \quad \mathbf{h}_u = \text{GRU}(H_t)$$

$$H_{t+1} = \begin{cases} \{H_t^2, \dots, H_t^N, a_t^l\} & r_t^l > 0 \\ H_t & \text{Otherwise} \end{cases}$$

$$P_k = \text{softmax}(W_k \mathcal{V}^\top), \quad k = 1, \dots, K,$$

$$a_t^k = \arg \max_{i \in \{1, \dots, |\mathcal{I}|\}} P_{k,i}, \quad \forall k = 1, \dots, K,$$

# Model Specification

---

**Algorithm 1:** Parameters Training for FCPO

---

```
1 Input: step size  $\delta$ , cost limit value  $\mathbf{d}$ , and line search ratio  $\beta$ 
2 Output: parameters  $\theta$ ,  $\omega$  and  $\phi$  of actor network, value function,
   cost function
3 Randomly initialize  $\theta$ ,  $\omega$  and  $\phi$ .
4 Initialize replay buffer  $D$ ;
5 for  $Round = 1 \dots M$  do
6   Initialize user state  $s_0$  from log data;
7   for  $t = 1 \dots T$  do
8     Observe current state  $s_t$  based on Eq. (12);
9     Select an action  $a_t = \{a_t^1, \dots, a_t^K\} \in \mathcal{I}^K$  based on Eq.
      (14) and Eq. (15)
10    Calculate reward  $r_t$  and cost  $c_t$  according to environment
        feedback based on Eq. (8) and Eq. (9);
11    Update  $s_{t+1}$  based on Eq. (12);
12    Store transition  $(s_t, a_t, r_t, c_t, s_{t+1})$  in  $D$  in its
        corresponding trajectory.
13  end
14  Sample minibatch of  $N$  trajectories  $\mathcal{T}$  from  $D$ ;
15  Calculate advantage value  $A$ , advantage cost value  $A_c$ ;
16  Obtain gradient direction  $d_\theta$  by solving Eq. (4) with  $A$  and  $A_c$ ;
17  repeat
18     $\theta' \leftarrow \theta + d_\theta$ 
19     $d_\theta \leftarrow \beta d_\theta$ 
20  until  $\pi_{\theta'}(s)$  in trust region & loss improves & cost  $\leq \mathbf{d}$ ;
21  (Policy update)  $\theta \leftarrow \theta'$ ;
22  (Value update) Optimize  $\omega$  based on Eq.(16);
23  (Cost update) Optimize  $\phi$  based on Eq.(17);
24 end
```

---

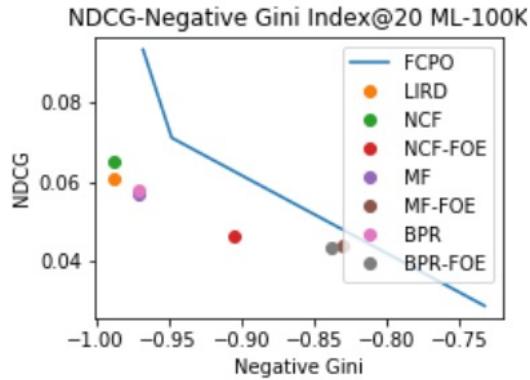
# Experiment

Methods	Recall (%) ↑			F1 (%) ↑			NDCG (%) ↑			Gini Index (%) ↓			Popularity Rate (%) ↑		
	K=5	K=10	K=20	K=5	K=10	K=20	K=5	K=10	K=20	K=5	K=10	K=20	K=5	K=10	K=20
Movielens-100K															
MF	1.847	3.785	7.443	2.457	3.780	5.074	3.591	4.240	5.684	98.99	98.37	97.03	99.98	99.96	99.92
BPR-MF	1.304	3.539	8.093	1.824	3.592	5.409	3.025	3.946	5.787	98.74	98.17	97.01	99.87	99.87	99.78
NCF	<u>1.995</u>	3.831	6.983	<u>2.846</u>	<u>4.267</u>	<u>5.383</u>	<u>5.319</u>	<u>5.660</u>	<u>6.510</u>	99.70	99.39	98.80	100.0	100.0	100.0
LIRD	1.769	<u>5.467</u>	<u>8.999</u>	2.199	4.259	4.934	3.025	3.946	5.787	99.70	99.41	98.81	100.0	100.0	100.0
MF-FOE	1.164	2.247	4.179	1.739	2.730	3.794	3.520	3.796	4.367	<u>86.29</u>	<u>84.05</u>	<u>82.98</u>	92.90	91.89	90.98
BPR-FOE	0.974	2.053	4.404	1.496	2.568	3.933	3.127	3.514	4.332	86.50	84.38	83.78	<u>92.17</u>	<u>91.36</u>	<u>90.70</u>
NCF-FOE	1.193	1.987	4.251	1.759	2.398	3.698	4.033	3.897	4.633	96.92	94.53	90.44	100.0	100.0	100.0
FCPO-1	<b>4.740</b>	<b>8.607</b>	<b>14.48</b>	<b>4.547</b>	<b>5.499</b>	<b>5.855</b>	<b>6.031</b>	<b>7.329</b>	<b>9.323</b>	98.73	98.07	96.75	92.60	90.42	85.85
FCPO-2	3.085	5.811	10.41	3.270	4.164	4.953	4.296	5.203	7.104	97.95	96.88	94.78	70.07	68.28	65.55
FCPO-3	0.920	1.668	3.329	1.272	1.807	2.535	2.255	2.369	2.871	<u>75.23</u>	<b>74.06</b>	<b>73.23</b>	<b>36.52</b>	<b>36.66</b>	<b>36.94</b>
Movielens-1M															
MF	1.152	2.352	4.650	1.701	2.814	4.103	3.240	3.686	4.574	99.44	99.18	98.74	99.92	99.90	99.86
BPR-MF	1.240	2.627	5.143	1.773	2.943	4.197	3.078	3.593	4.632	98.93	98.44	97.61	99.40	99.23	98.96
NCF	1.178	2.313	4.589	1.832	2.976	<u>4.382</u>	<u>4.114</u>	<u>4.380</u>	<u>5.080</u>	99.85	99.71	99.42	100.0	100.0	100.0
LIRD	<u>1.961</u>	<u>3.656</u>	<u>5.643</u>	<u>2.673</u>	<u>3.758</u>	4.065	3.078	3.593	4.632	99.87	99.73	99.46	100.0	100.0	95.00
MF-FOE	0.768	1.534	3.220	1.246	2.107	3.345	3.321	3.487	4.021	92.50	91.06	91.32	98.89	98.78	98.68
BPR-FOE	0.860	1.637	3.387	1.374	2.233	3.501	3.389	3.594	4.158	<u>90.48</u>	<u>88.92</u>	<u>89.01</u>	<u>96.56</u>	<u>96.12</u>	<u>95.78</u>
NCF-FOE	0.748	1.403	2.954	1.230	1.980	3.175	3.567	3.589	4.011	97.73	96.57	<u>95.04</u>	100.0	100.0	100.0
FCPO-1	<b>2.033</b>	<b>4.498</b>	<b>8.027</b>	<b>2.668</b>	<b>4.261</b>	<b>5.201</b>	<b>4.398</b>	<b>5.274</b>	<b>6.432</b>	99.81	99.67	99.34	99.28	96.93	91.70
FCPO-2	1.520	3.218	6.417	2.015	3.057	4.145	3.483	3.920	5.133	99.47	99.10	97.41	72.66	68.27	71.35
FCPO-3	0.998	1.925	3.716	1.449	2.185	2.948	2.795	2.987	3.515	<u>88.97</u>	<b>88.34</b>	<b>87.70</b>	<b>63.43</b>	<b>62.73</b>	<b>61.45</b>

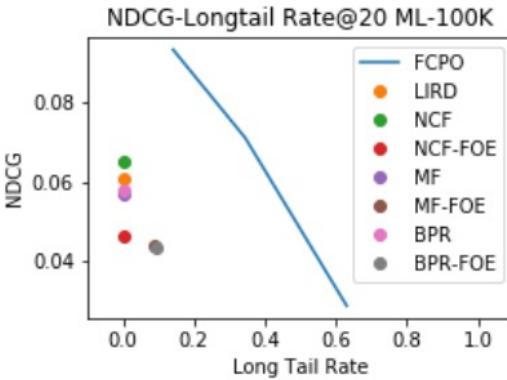
# Experiment

- Our FCPO approach achieves the best top-K recommendation performance against all baselines on both datasets, implying that the proposed method does have the ability to capture dynamic user-item interactions,
- We can easily see that there exists a trade-off between the recommendation performance and the fairness performance both in FCPO and FOE.

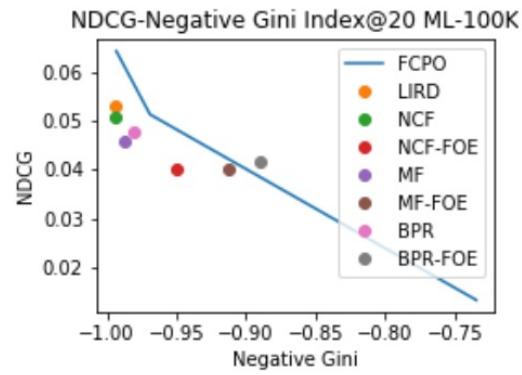
# Experiment



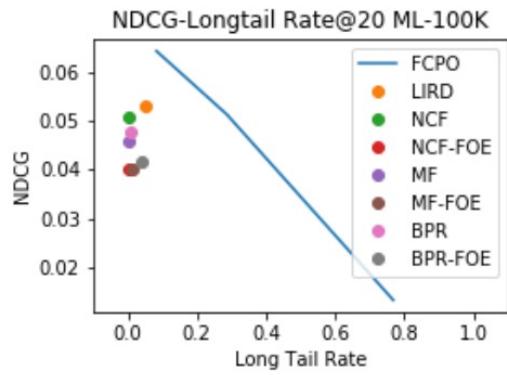
(a) NDCG vs Negative Gini on ML100K



(b) NDCG vs Long-tail Rate on ML100K



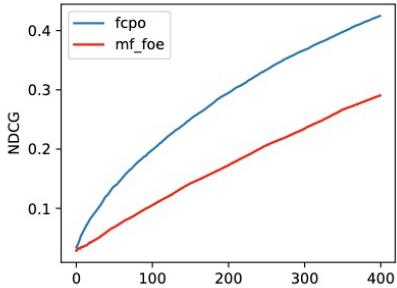
(c) NDCG vs Negative Gini on ML1M



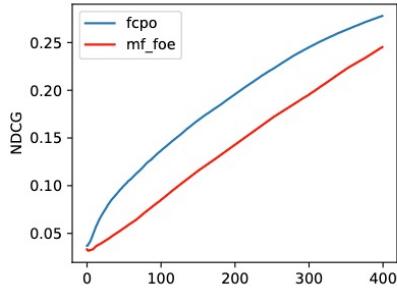
(d) NDCG vs Long-tail Rate on ML1M

- For the same Gini Index, our method achieves much better NDCG; meanwhile, under the same NDCG scores, our method achieves better fairness.
- FCPO can achieve much better trade-off than FOE in both individual fairness and group fairness.

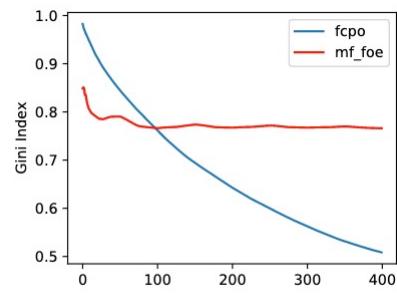
# Experiment



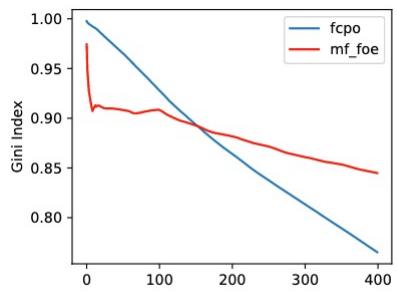
(a) NDCG on Movielens100K



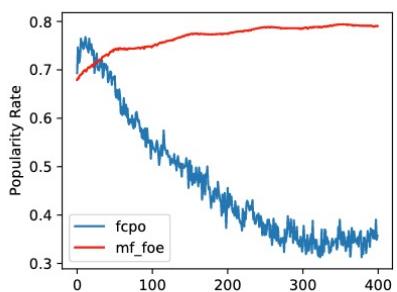
(b) NDCG on Movielens1M



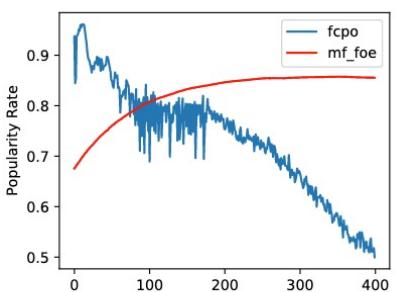
(c) Gini on Movielens100K



(d) Gini on Movielens1M



(e) Popularity Rate on Movielens100K



(f) Popularity Rate on Movielens1M

- We compared FCPO with a static short-term fairness solution (i.e., MF-FOE) for 400 steps of recommendation.
- Since FCPO makes adjustment of its policy according to the fairness feedback, it can successfully and continuously suppress the fairness metric to a much lower value during testing.
- The overall ranking performance of MF-FOE is consistently outperformed by FCPO, which indicates that MF-FOE sacrifices the recommendation performance more than FCPO in order to control fairness.

# Contributions

1. We propose to model the sequential interactions between consumers and recommender systems under the fairness constraint of item exposure as a **CMDP**.
2. We leverage the **Constrained Policy Optimization** (CPO) with adapted neural network architecture to automatically learn the optimal policy under different fairness constraints.
3. To the best of our knowledge, this is the first attempt to model the dynamic nature of fairness with respect to **changing group labels**, and to show its effectiveness in the long term.

# Q&A

A large, colorful word cloud centered around the words "thank you" in various languages. The word "thank" is in blue, "you" is in yellow, and "you" is in red. The background is white with a light gray grid. The word cloud includes many other words related to gratitude and thanks in different languages, such as "danke" (German), "gracias" (Spanish), "merci" (French), "obrigado" (Portuguese), "спасибо" (Russian), "谢谢" (Chinese), "감사합니다" (Korean), "谢谢" (Mandarin), "merhaba" (Turkish), "mānana" (Hawaiian), and "mānana" (Maori). The words are in different colors and sizes, creating a dense and visually appealing pattern.