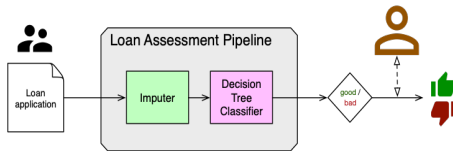


Mitigating Biases in Decision-Making Systems: a Control Systems Perspective

Giulia De Pasquale

ETH Zürich

December 11, 2024



Applications

employment

health

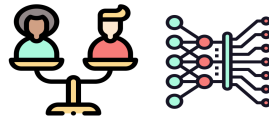
education

law

...

✓ High scalability

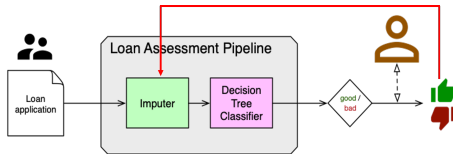
✗ Exacerbate existing biases
and even introduce new ones



Algorithmic fairness

💡 Enforce group fairness
metrics to mitigate biases

✗ solutions are designed for
stationary systems



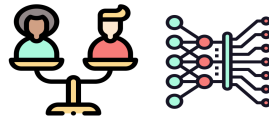
Applications

employment
health
education
law

...

✓ High scalability

✗ Exhacerbate existing biases
and even introduce new ones

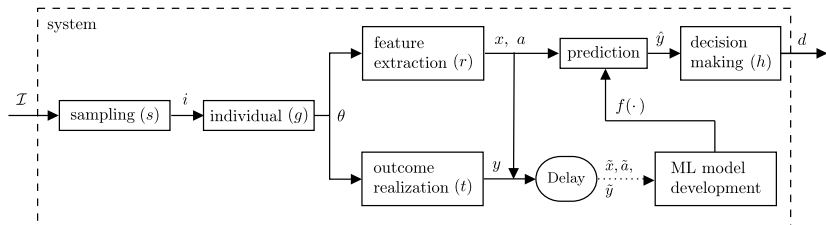


Algorithmic fairness

💡 Enforce group fairness
metrics to mitigate biases

✗ solutions are designed for
stationary systems

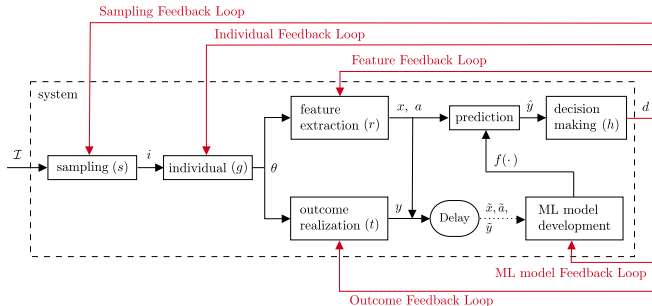
1



The ML-based decision making pipeline as an open loop system

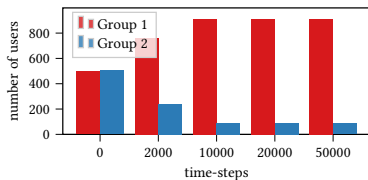
¹"A classification of feedback loops and their relation to biases in automated decision-making systems", J. Baumann, N. Pagan, E. Elokda, GDP, S. Bolognani, A. Hannak, Conference on Equity and Access in Algorithms, Mechanisms, and Optimization

1



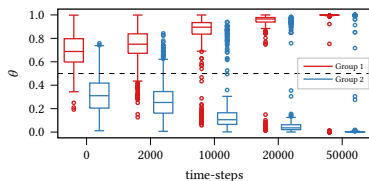
The ML-based decision making pipeline as a closed loop system

¹"A classification of feedback loops and their relation to biases in automated decision-making systems", J. Baumann, N. Pagan, E. Elokda, GDP, S. Bolognani, A. Hannak, Conference on Equity and Access in Algorithms, Mechanisms, and Optimization



Sampling FL: Representation bias

The available data is not representative of the population:
the ML model does not generalize well for the disadvantaged group, e.g. Amazon's Alexa.



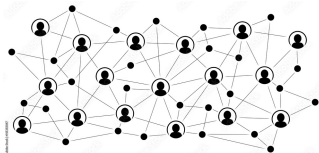
Individual FL: Historical bias

Users with high initial interests get recommended the item: θ increases over time. **Decisions change individual properties**, leads to polarization of interests.

A Solution to Representation Bias²

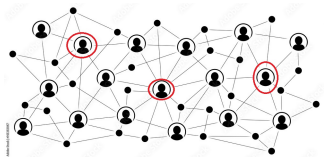
²"Fairness in Social Influence Maximization via Optimal Transport", S. Chowshary, GDP*, N. Lanzetti*, A. Stoica, F. Dörfler, NeurIPS 2024

Suppose you want to sell a product, or make an information spread as much as possible in a social network:



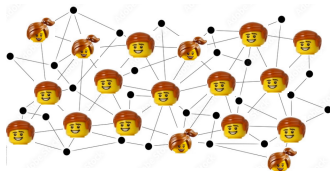
Social Influence Maximization (SIM) is the problem of how to strategically select seeds that spread information throughout a network in order to **maximize the outreach**.

Suppose you want to sell a product, or make an information spread as much as possible in a social network:



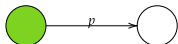
Social Influence Maximization (SIM) is the problem of how to strategically select seeds that spread information throughout a network in order to **maximize the outreach**.

Suppose you want to spread the news about an open position as Assistant Professor in Control Engineering:

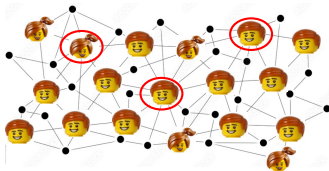


Fairness in SIM: solve SIM by ensuring **balanced outreach** among different communities, e.g. demographic groups.

Spreading mechanism: Independent cascade model

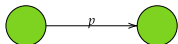


Suppose you want to spread the news about an open position as Assistant Professor in Control Engineering:



Fairness in SIM: solve SIM by ensuring **balanced outreach** among different communities, e.g. demographic groups.

Spreading mechanism: Independent cascade model



Given the groups C_1, \dots, C_m , a configuration is said to be

Equal, if the SIM algorithm chooses a seed set S such that

$$\frac{\mathbb{E}[|v \in S|v \in C_i|]}{|C_i|} = \frac{\mathbb{E}[|v \in S|v \in C_j|]}{|C_j|} \quad \forall i, j.$$

Equitable, if the SIM algorithm chooses a seed set S such that

$$\frac{\mathbb{E}[|v \text{ reached}|v \in C_i|]}{|C_i|} = \frac{\mathbb{E}[|v \text{ reached}|v \in C_j|]}{|C_j|} \quad \forall i, j.$$

Max-Min Fair, if the SIM algorithm chooses a seed set S such that

$$\min_{i \in [m]} \frac{\mathbb{E}[|v \text{ reached}|v \in C_i|]}{|C_i|}$$

is maximized.

Given the groups C_1, \dots, C_m , a configuration is said to be

Equal, if the SIM algorithm chooses a seed set S such that

$$\frac{\mathbb{E}[|v \in S|v \in C_i|]}{|C_i|} = \frac{\mathbb{E}[|v \in S|v \in C_j|]}{|C_j|} \quad \forall i, j.$$

Equitable, if the S

$$\mathbb{E}[|v \text{ rea}$$

Max-Min Fair, if

is maximized

ch that

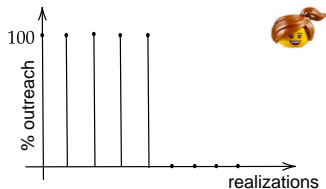
$$\mathbb{E}[|v \in S|v \in C_i|] \geq \mathbb{E}[|v \in S|v \in C_j|] \quad \forall i, j.$$

S such that



What's wrong with the Expectation?

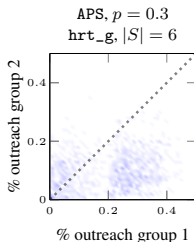
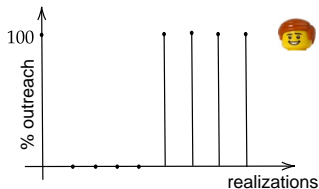
Consider the outcome: "In 50% of the cases, no one in group 1 gets the information and everyone in group 2 does, and in the other 50 % it is the opposite."



$$\frac{\mathbb{E}[\% \text{ reached} | v \in C_1]}{|C_1|} = \frac{\mathbb{E}[\% \text{ reached} | v \in C_2]}{|C_2|}$$

The outcome is classified as **equitable**, however it is **highly unfair**.

Note: this also happens in experimental settings!



We want to answer questions such as as:

- i) When group 1 receives the information, will group 2 also receive it?
- ii) Even if the two groups have the same marginal outreach probability distributions, will the final configurations always be **fair**?

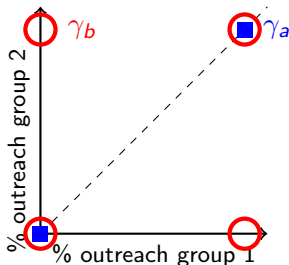


Figure: Illustration of the (γ_a, γ_b) example.


Marginals: $\mu_i = \frac{1}{2}\delta_0 + \frac{1}{2}\delta_1, i \in \{1, 2\}$

Distributions:

$$\gamma_a = 0.5 \cdot \delta_{(0,0)} + 0.5 \cdot \delta_{(1,1)}, \quad \gamma_b = 0.25 \cdot \delta_{(0,0)} + 0.25 \cdot \delta_{(1,1)} + 0.25 \cdot \delta_{(0,1)} + 0.25 \cdot \delta_{(1,0)}$$



Use the **joint** outreach probability distribution to capture the correlation between the two groups!

 Quantify fairness by computing the distance of the probability distribution γ from an ideal reference distribution γ^* along the diagonal.

Optimal Transport Problem: quantifies the minimum transportation cost to morph γ into γ^* when transporting a unit of mass from (x_1, x_2) to (y_1, y_2) costs $c((x_1, x_2), (y_1, y_2))$.

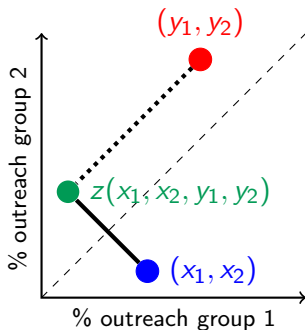
$$W_c(\gamma, \gamma^*) = \min_{\pi \in \Pi(\gamma, \gamma^*)} \mathbb{E}_{(x_1, x_2), (y_1, y_2) \sim \pi} [c((x_1, x_2), (y_1, y_2))]$$

Ingredients:

- i) transportation cost;
- ii) reference distribution.

Transportation Cost:

- moving mass **along** the diagonal costs 0, as it does not affect fairness
- moving mass **orthogonally** towards the diagonal comes at a price. We quantify the price as the **Euclidean distance**.



$$c((x_1, x_2), (y_1, y_2)) = \|z(x_1, x_2, y_1, y_2) - (x_1, x_2)\| = \frac{\sqrt{2}}{2} |(x_2 - x_1) - (y_2 - y_1)|,$$

Definition (Mutual Fairness)

Given a network with communities $(C_i)_{i \in [2]}$, a SIM algorithm is said to be *mutually fair* if the algorithm propagation is such that it maximizes

$$\text{FAIRNESS}(\gamma) := 1 - \sqrt{2}W_c(\gamma, \gamma^*),$$

$$W_c(\gamma, \gamma^*) = \min_{\pi \in \Pi(\gamma, \gamma^*)} \mathbb{E}_{(x_1, x_2), (y_1, y_2) \sim \gamma, [c((x_1, x_2), (y_1, y_2))]} \text{ and } \gamma^* = \delta_{(1,1)}.$$

Observations:

- $\min \text{FAIRNESS}(\gamma) = 0$; $\text{argmin} = \gamma = \delta_{(0,1)}$;
- $\max \text{FAIRNESS}(\gamma) = 1$; $\text{argmax} = \gamma^*$.
- since γ^* is a delta distribution, we can solve the OT problem in closed form and $\text{FAIRNESS}(\gamma) = 1 - \frac{1}{N} \sum_{i=1}^N |x_{1,i} - x_{2,i}|$

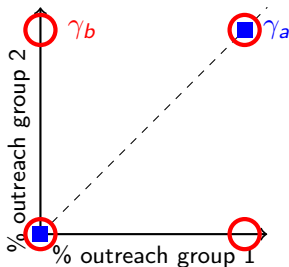


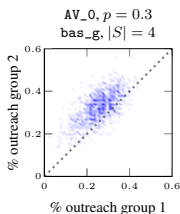
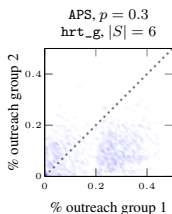
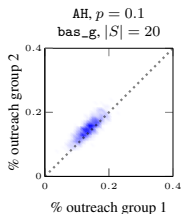
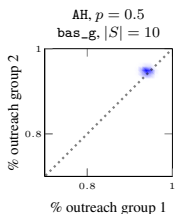
Figure: Illustration of the (γ_a, γ_b) example.

$$\text{FAIRNESS}(\gamma_a) = 1$$

$$\text{FAIRNESS}(\gamma_b) = 0.5.$$

Joint outreach probability distribution for different real datasets, each with a chosen demographic partitioning the population in two groups.

Four qualitative outcomes:



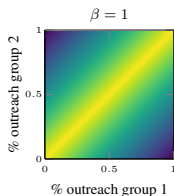
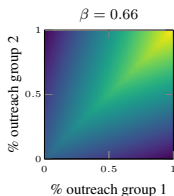
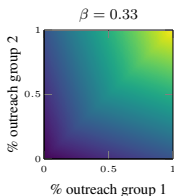
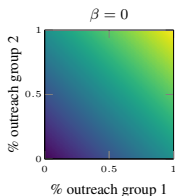
For both $\gamma = \delta_{(0,0)}$ and $\gamma^* = \delta_{(1,1)}$ the fairness score is maximal:

We need a fairness-efficiency trade-off!

We can define the transportation cost as a weighted sum:

$$\begin{aligned} c_\beta((x_1, x_2), (y_1, y_2)) &= \\ \beta \|z(x_1, x_2, y_1, y_2) - (x_1, x_2)\| &+ (1 - \beta) \|z(x_1, x_2, y_1, y_2) - (y_1, y_2)\| = \\ \beta \frac{\sqrt{2}}{2} |(x_2 - x_1) - (y_2 - y_1)| &+ (1 - \beta) \frac{\sqrt{2}}{2} |(x_1 + x_2) - (y_1 + y_2)|. \end{aligned}$$

Heatmap of c_β :



Definition (β -Fairness)

Consider a network with groups C_1, C_2 , a SIM algorithm is said to be β -fair if the algorithm propagation is such that it maximizes

$$\beta - \text{FAIRNESS}(\gamma) := 1 - \frac{\sqrt{2}}{\max\{1, 2 - 2\beta\}} W_{c_\beta}(\gamma, \gamma^*),$$

The OT problem can be solved in closed form

$$\beta - \text{FAIRNESS}(\gamma) = \mathbb{E}_{(x_1, x_2) \sim \gamma} \left[1 - \frac{\beta|x_1 - x_2| + (1 - \beta)|x_1 + x_2 - 2|}{\max\{1, 2 - 2\beta\}} \right]$$

In particular, for $\beta = 1$, we recover the mutual fairness $\text{FAIRNESS}(\gamma)$ and for $\beta = 0$ we obtain the efficiency metric $\mathbb{E}_{(x_1, x_2) \sim \gamma} \left[1 - \frac{x_1 + x_2 - 2}{2} \right]$.

Algorithm 1 Stochastic Seedset Selection Descent

Input: Social Graph $G(V_G, E_G)$, initial seed set S_0 , β fairness weight, ϵ -tolerance

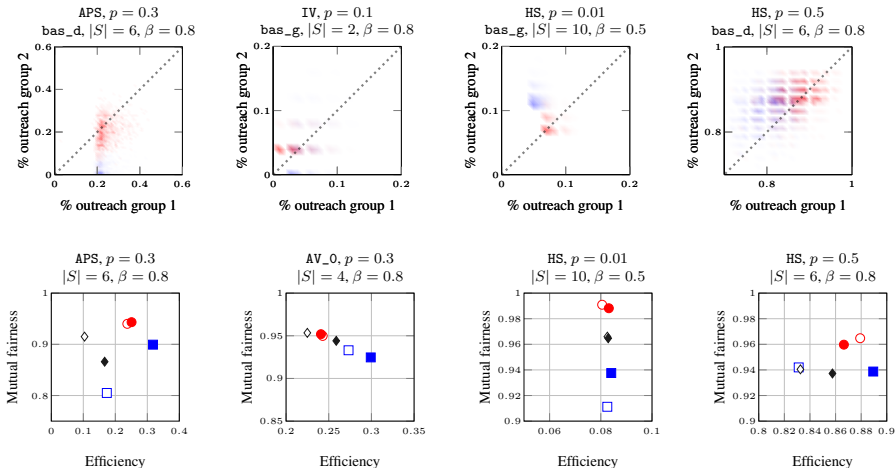
Output: Optimal seedset S^*

```

1:  $\mathcal{S} \leftarrow \{\}$ ,  $S \leftarrow S_0$                                 ▷ initial collection of candidates, running seedset
2: for  $k$  iterations do                                       ▷ configurable  $k$ 
3:    $V_S \leftarrow$  nodes reachable from  $S$  via cascade, using SEEDSET_REACH routine
4:    $S' \leftarrow \{\}$ 
5:   for  $|S|$  iterations do                                     ▷ searching nearby states,  $V_{S'}$ , to get  $S'$  (Appendix E.3)
6:      $S' \leftarrow S' \cup \{v\} \mid v \sim V_S$ 
7:      $V_{S'} \leftarrow$  nodes reachable from  $S'$  in a fixed horizon, using SEEDSET_REACH
8:      $V_S \leftarrow V_S \setminus V_{S'}$ 
9:      $E_S \leftarrow -\text{BETA\_FAIRNESS}(S, \beta)$                    ▷ expected potential energy defined on  $\beta$ -fairness
10:     $E_{S'} \leftarrow -\text{BETA\_FAIRNESS}(S', \beta)$ 
11:     $p_{\text{accept}} \leftarrow \min\{1, e^{E_S - E_{S'}}\}$              ▷  $S'$  acceptance on energy minimization
12:    if  $x \sim \mathcal{B}(p_{\text{accept}})$  then                             ▷ Metropolis sampling
13:       $S^+ \leftarrow S'$                                        ▷ get a better seedset
14:    else
15:      if  $x \sim \mathcal{B}(\epsilon)$  then                                 ▷ for some small constant  $\epsilon$ 
16:         $S^+ \leftarrow \{v_i\}_{i=1}^{|S|} \overset{|S|}{\sim} V_G$              ▷ random seedset
17:      else
18:         $S^+ \leftarrow S$                                        ▷ retain existing choice
19:       $\mathcal{S} \leftarrow \mathcal{S} \cup \{S^+\}$ 
20:       $S \leftarrow S^+$                                        ▷ for next iteration
21:  $S^* \leftarrow S \in \mathcal{S} \mid \text{BETA\_FAIRNESS}(S, \beta)$  is maximum  ▷ via S3D_ITERATE
22: return  $S^*$ 

```

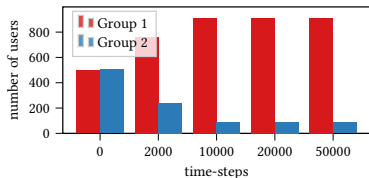

Are the outcomes more fair?



Greedy-based algorithms: ■ = bas_g , ● = S3D_g , and ◆ = hrt_g .
Degree-based algorithms: □ = bas_d , ○ = S3D_d , and ◇ = hrt_d .

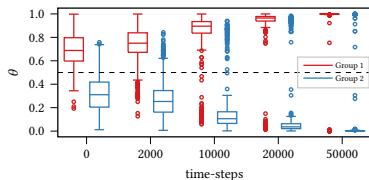
- New fairness metric for SIM that captures new fairness-related aspects;
- We leverage β -fairness to design a new seed selection strategy that tradeoffs fairness and efficiency;
- We show superior fairness performance with minor decrease in efficiency.

Note: Mutual fairness is applicable whenever you have empirical distributions associated with groups.



Sampling FL: Representation bias

The available data is not representative of the population:
the ML model does not generalize well for the disadvantaged group, e.g. Amazon's Alexa



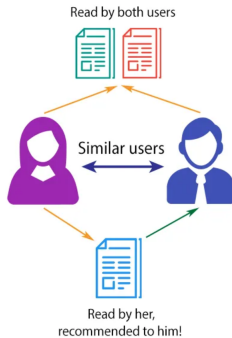
Individual FL: Historical bias

Users with high initial interests get recommended the item: θ increases over time. **Decisions change individual properties**, leads to polarization of interests.

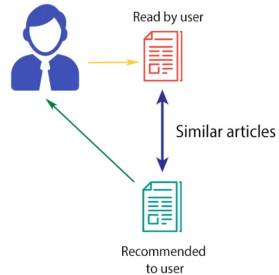
A Solution to Historical Bias³

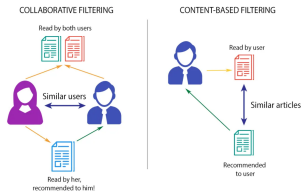
³S. Chandrasekaran, GDP, G. Belgioioso, F. Dörfler, "Mitigating Polarization in Recommender Systems via Network-aware Feedback Optimization", submitted.

COLLABORATIVE FILTERING

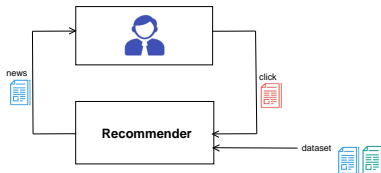


CONTENT-BASED FILTERING



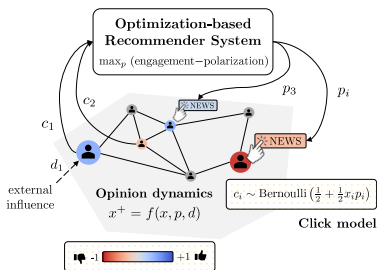


Make the feedback loop explicit, to understand



- i) the impact of recommendation on users opinions;
- ii) how recommender systems should depart from engagement maximization to mitigate polarization.

We leverage on **online feedback optimization** to design a RS as a dynamic feedback controller that mitigates polarization by providing user personalized content, using only **implicit feedback**.



Assumption: one single topic of discussion

Assumption: The dynamics is exponentially stable and admits a unique steady-state map

$$h(p, d) = f(h(p, d), p, d)$$

with $h(p, d)$ continuously-differentiable and L -lipschitz wrt p .

$$\min_{p,x} \varphi^{\text{clk}}(p,x) + \gamma \varphi^{\text{pol}}(x)$$

$$\text{s.t. } x = h(p,d)$$

$$p \in [-1, 1]^n$$

$$\varphi^{\text{clk}} = - \sum_{i \in [n]} \mathbb{E}_{c_i \sim \mathcal{B}(g_i(x_i, p_i))} [c_i]$$

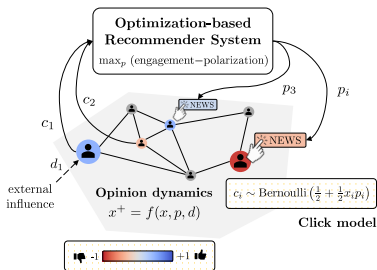
$$\varphi^{\text{pol}}(x) = \|x\|^2$$

Challenges:

- only clicks are available: opinions, opinion dynamics, network topology, clicking behaviour, external influence unknown \rightarrow the problem must be solved online
- non-convex problem

The recommender system only relies on clicks:

$$\frac{\#\text{clk}}{\#\text{news}} \approx \mathbb{E}[\mathcal{B}(g(p,x))] = g(p,x).$$



The recommender system dynamically generates recommendation via projected gradient descent

$$p^+ = \text{proj}_{[-1,1]}[p - \eta \underbrace{(\nabla_p \varphi(p, x) + \nabla_p h(p, d)^\top \nabla_x \varphi(p, x))}_{\nabla \varphi}]$$

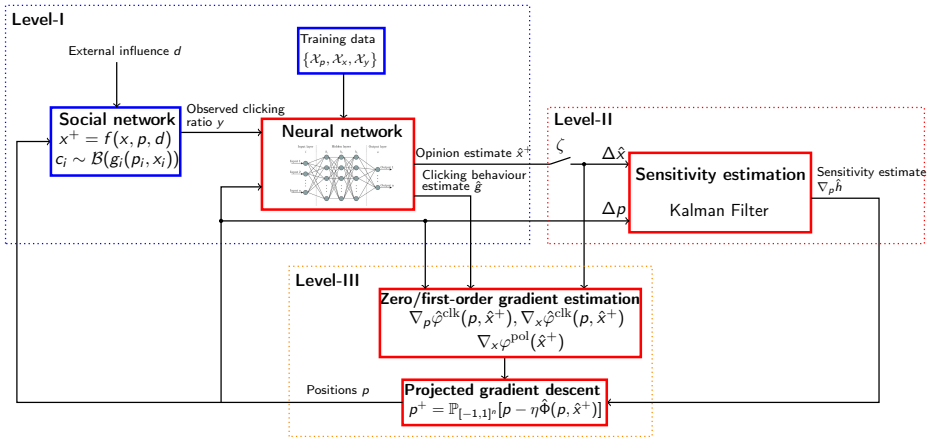
$$\varphi = \varphi^{\text{clk}} + \varphi^{\text{pol}}.$$

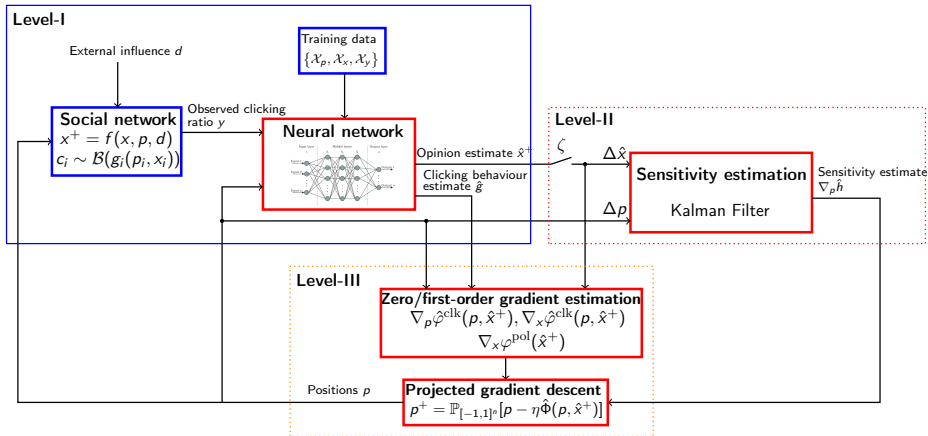
Challenges

Evaluating $\nabla \varphi$ requires access to:

- i) Online opinions x
- ii) Sensitivity mapping $\nabla_p h(p, d)$
- iii) Gradients $\nabla_p \varphi(p, x)$, $\nabla_x \varphi(p, x)$

None of these information is available online!



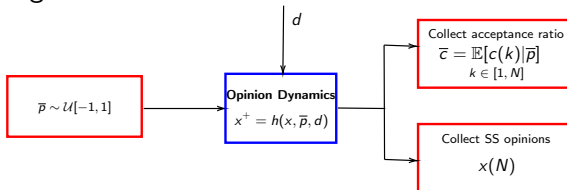


The recommender system dynamically generates recommendation via projected gradient descent

$$p^+ = \text{proj}_{[-1,1]}[p - \eta \underbrace{(\nabla_p \varphi(p, x) + \nabla_p h(p, d)^\top \nabla_x \varphi(p, x))}_{\nabla \varphi}]$$

Training data collection

Repeat #training times:



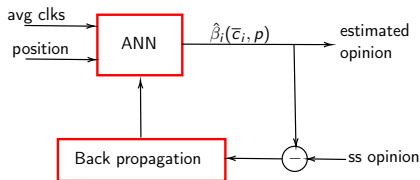
Assumption:

- i) There exists a continuous mapping $\beta(\bar{c}, p) = x + \theta(x)$, $\|\theta(x)\| \leq \theta$
- ii) $g(x, p)$ is Lipschitz and globally smooth.

There exists α s.t.

$$g(p, \beta(\bar{c}, p)) = \bar{c} + \nabla_x g(p, x)^\top \theta(x) + \alpha(\bar{c}),$$

$$\|\alpha(\bar{c})\| \leq \alpha$$



$$\begin{cases} \hat{x}_i^+ = \hat{\beta}_i(\bar{c}_i, p) \\ p \text{ via OFO} \end{cases}$$

Opinion estimation error

$$\| \overbrace{h(p, d)}^{\epsilon_x} - \hat{\beta} \| \leq \sqrt{n} (\sup_{\bar{c}, p} \|\beta - \hat{\beta}\|_\infty + \text{ANN bias})$$

Training is carried out distributedly

³Tabuada, Charesifard, "Universal approximation power of deep residual neural networks through the lens of control", TAC, 2023

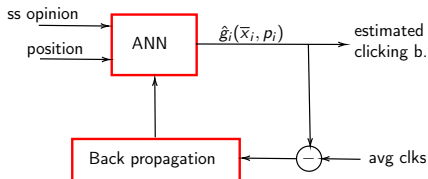
Assumption:

- i) There exists a continuous mapping $\beta(\bar{c}, p) = x + \theta(x)$, $\|\theta(x)\| \leq \theta$
- ii) $g(x, p)$ is Lipschitz and globally smooth.

There exists α s.t.

$$g(p, \beta(\bar{c}, p)) = \bar{c} + \nabla_x g(p, x)^\top \theta(x) + \alpha(\bar{c}),$$

$$\|\alpha(\bar{c})\| \leq \alpha$$

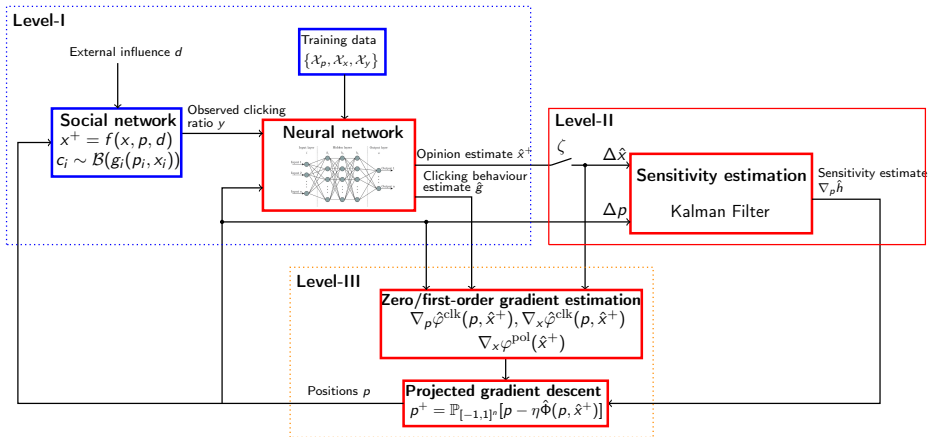


$$\begin{cases} \hat{c}_i^+ = \hat{g}_i(\hat{x}_i^+, p_i) \\ p \text{ via OFO} \end{cases}$$

clicking behaviour estimation error

$$\underbrace{\|\hat{g}(p, \hat{x}) - g(p, h(p, d))\|}_{\epsilon_g} \leq \sqrt{n}(\sup_{p,x} \|g(p, x) - \hat{g}(p, x)\|_\infty + \text{ANN bias} + f(\theta, \alpha))$$

³Tabuada, Charesifard, "Universal approximation power of deep residual neural networks through the lens of control", TAC, 2023



The recommender system dynamically generates recommendation via projected gradient descent

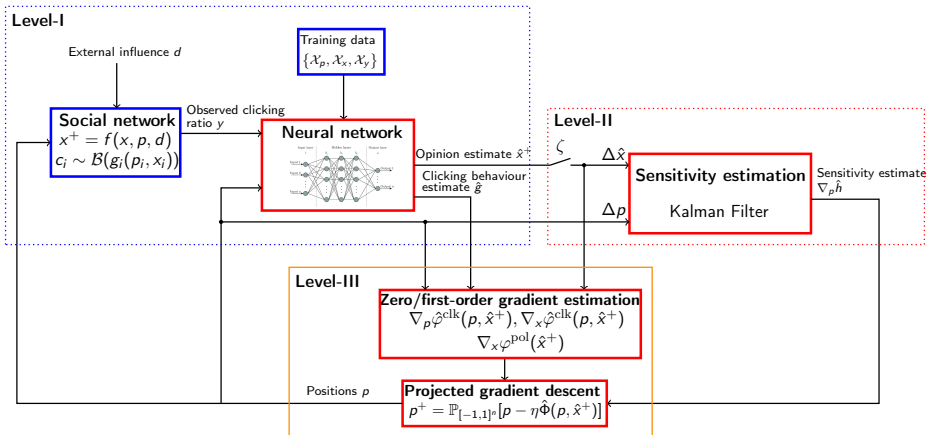
$$p^+ = \text{proj}_{[-1,1]}[p - \eta \underbrace{(\nabla_p \varphi(p, x) + \nabla_p h(p, d)^\top \nabla_x \varphi(p, x))}_{\nabla \varphi}]$$

To estimate the sensitivity online we rely on **Kalman filter**.

Note: $\nabla_p h_{ij}(p, d) \neq 0 \rightarrow j$ and i are connected

To ensure the sensitivity estimate is accurate:

Assumption: The inputs p are persistently exciting.



The recommender system dynamically generates recommendation via projected gradient descent

$$p^+ = \text{proj}_{[-1,1]}[p - \underbrace{\eta(\nabla_p \varphi(p, x) + \nabla_p h(p, d)^\top \nabla_x \varphi(p, x))}_{\nabla \varphi}]$$

$\varphi = \varphi^{\text{clk}} + \varphi^{\text{pol}}$. Estimation via forward difference method

$$\nabla_x \hat{\varphi}_i^{\text{clk}}(p, x) = \frac{\hat{\varphi}^{\text{clk}}(p, x + \mu e_i) - \hat{\varphi}^{\text{clk}}(p, x)}{\mu}$$

$$\nabla_p \hat{\varphi}_i^{\text{clk}}(p, x) = \frac{\hat{\varphi}^{\text{clk}}(p + \mu e_i, x) - \hat{\varphi}^{\text{clk}}(p, x)}{\mu},$$

Gradient estimation error

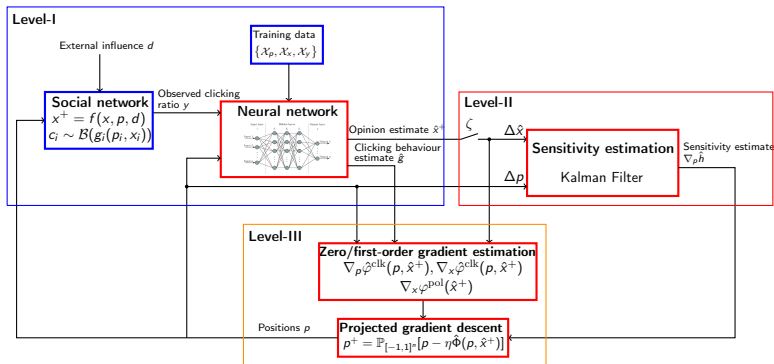
Under the previous regularity assumptions on β, g

$$\|\nabla \hat{\varphi}^{\text{clk}} - \nabla \varphi^{\text{clk}}\| \leq \frac{1}{2} L_x \mu + 2 \frac{\sqrt{n} \epsilon_g}{\mu}; \quad \mu^* = 2n^{1/4} \sqrt{\frac{\epsilon_g}{L}}$$

Smoothing parameter μ , requires fine tuning: small, but not too much!

We now collected all the ingredients to run gradient descent for the recommender system algorithm:

$$p^{k+1} = \text{proj} \left[p^k - \eta \zeta^k \left(\nabla_p \hat{\varphi}^{\text{clk}}(p^k, \hat{x}^k) + \nabla_p \hat{h}(p^k, d)^\top \nabla_x \hat{\varphi}^{\text{clk}}(p^k, \hat{x}^k) \right) \right]$$



Initialization

Collect data during training

Build opinion and clicking behaviour estimators $(\hat{\beta}, \hat{g})$

Optimization phase

for $k \geq 0$ do

Collect clicks $c_i^k \sim \mathcal{B}(g_i(p_i^k, x_i^k))$ from users

CTR $y^k \leftarrow \frac{\sum_{t=\tau_i}^k c^t}{k-\tau_i+1}$, $\tau_i = (i-1)T < k$

Estimate opinions $\hat{x}_i^{k+1} \leftarrow \hat{\beta}_i(y_i^k, p^k)$

if $\zeta^k = 1$ then

$\mathcal{T} \leftarrow \text{append}[k]$

Estimate sensitivity \hat{H}^k

Estimate gradient

Update positions p^{k+1}

else

$\hat{H}^k \leftarrow \hat{H}^{k-1}$; $p^{k+1} \leftarrow p^k$

end if

end for

We ensure convergence by using the **gradient mapping**

$$\mathcal{G}(p) := \frac{1}{\eta} \left(p - \text{proj}_{[-1,1]}[p - \eta(\nabla\varphi)] \right)$$

a common metric to quantify convergence in non convex-regimes.

OFO Convergence

Under all the previous assumptions, for $\eta \in (0, \frac{1}{2(L')})$, $\mu = \mu^*$, the position sequence generated by the projected gradient descent algorithm satisfies

$$\frac{1}{|\mathcal{T}|} \sum_{\substack{l \in \mathcal{T} \\ l \leq k}} \mathbb{E} \left[\|\mathcal{G}(p^l)\|^2 \right] \leq K_1, \quad \forall k \geq T$$

$K_1 \propto \varphi(p^0, h(p^0, d)) - \varphi^*, \epsilon_x^2, \epsilon_g^2, L'^2, \frac{1}{\eta^2}$, gradient est. error

Opinion Dynamics and Clicking Behaviour

Extended FJ model

$$x^+ = (I - \Gamma_p - \Gamma_d)Ax + \Gamma_p p + \Gamma_d d$$

Users follow two clicking behaviours

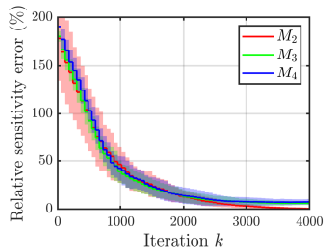
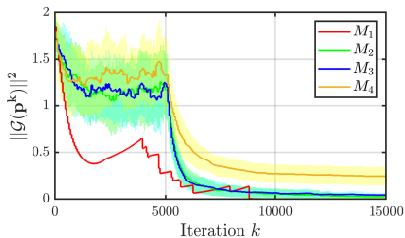
$$c_i \sim \mathcal{B}\left(\underbrace{\frac{1}{2} + \frac{1}{2}x_i p_i}_{C_a}\right), \quad c_i \sim \mathcal{B}\left(\underbrace{\frac{1}{2} + \frac{1}{2}e^{-c(x_i - p_i)^2}}_{C_b}\right)$$

we perform our algorithm over a network of 15 users, with C_a and C_b randomly distributed. Initial opinion $\sim \mathcal{U}[-1, 1]$, A substochastic, $d^k = x^0 + \text{noise}$, $\Gamma_p \sim \mathcal{U}[10^{-2}, 0.5]$

Training We train the NN for opinion and clicking behaviour with horizon $N = 100$ and collect 75 data points, with trigger period $T = 60$, with the clicks being recorded in the interval $[N - T, N]$. We take $m = 375$ training and 125 testing points.

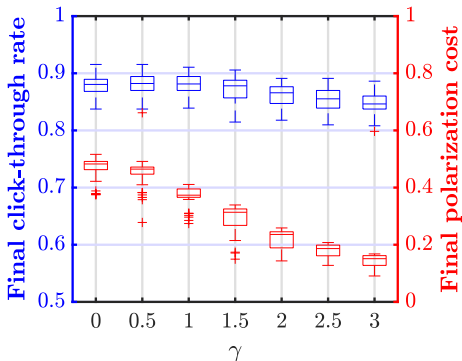
Online We set $p^0 = 0$ (neutral recommendations). All simulations are conducted for $N = 10^3$ over 50 Monte-Carlo trials.

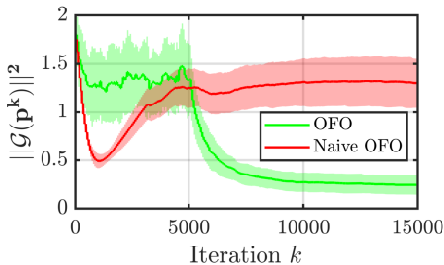
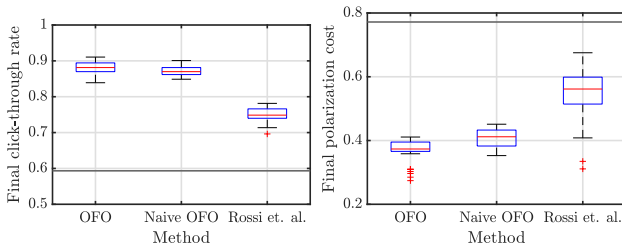
Method	Sensitivity	Opinions	Clicking behaviour
M_1 (Oracle)	✓	✓	✓
M_2	✗	✓	✓
M_3	✗	✗	✓
M_4 (Alg. 1)	✗	✗	✗



$$\min_{p,x} \varphi^{\text{clk}}(p, x) + \gamma \varphi^{\text{pol}}(x)$$

s.t. $x = h(p, d)$
 $p \in [-1, 1]^n$



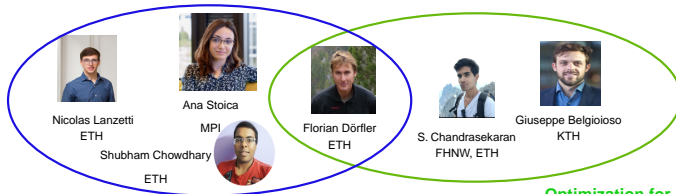


Conclusions

- A Model-free recommender system algorithm that balances engagement maximization and polarization mitigation;
- Theoretical guarantees for CL stability;
- Validation on synthetic data

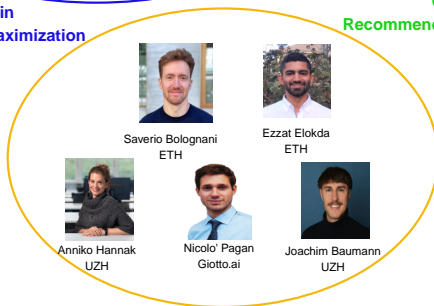
Future Directions

- Relax smoothness hypothesis on clicking behaviour;
- Consider other interests drivers than confirmation bias, e.g. repulsion.



**Fairness in
Social Influence Maximization**

**Optimization for
Recommendation Systems**



**Biases in
Automated Decision Making**

Thanks for your attention

Appendix

Sensitivity dynamics as a random process¹:

$$\text{vec}(\nabla_p h(p, d))^+ = \text{vec}(\nabla_p h(p, d)) + w \quad \text{Process model}$$

$$\Delta x_{\text{SS}}^+ = \Delta \tilde{p} * \text{vec}(\nabla_p h(p, d)) + v \quad \text{Measurement model}$$

where

- $\Delta x_{\text{SS}}^+ = h(p^k, d) - h(p^{k-1}, d)$
- $w^k \sim \mathcal{N}(0, Q^k)$
- $v^k \sim \mathcal{N}(0, R^k)$, accounts for the external influence
- $\Delta \tilde{p} = (p^k - p^{k-1})^\top \otimes I_n$

Sensitivity and covariance updates:

$$\text{vec}(\nabla_p h)^k = \text{vec}(\nabla_p h)^{k-1} + \zeta^k (K^{k-1} \Delta \hat{x}^{k+1} - \Delta \tilde{p}^k \text{vec}(\nabla_p h)^{k-1})$$

$$\Sigma^k = \Sigma^{k-1} + \zeta^k (Q^k - K^{k-1} \Delta \tilde{p}^k \Sigma^{k-1})$$

Trigger mechanism: Enforces time-scale separation and ensures that a sufficient number of clicks is collected (clicking ratio accuracy).

¹Picallo, Ortman, Bolognani, Dörfler, *Adaptive real time grid operation via online feedback optimization with sensitivity estimation* Electric Power Systems Research, 2022

Note: The CTR is recorded over a time horizon with constant p . The dynamics is exponentially stable: the opinion estimate is close to the steady state opinion $h(p, d) \rightarrow$ we can treat the opinion dynamics as a static map.

CL Convergence

Under all the previous assumptions, the sensitivity estimation error $e^k := \text{vec}(h^k) - \text{vec}(\hat{h}^k)$ has bias and variance bounded in norm, with

$$\|\mathbb{E}[e^k]\| \leq J_1 \quad \mathbb{E}[\|e^k\|^2] \leq J_2$$

with $J_1, J_2 \propto \epsilon_x, \frac{1}{T}$ and $J_2 \propto \sigma_r^2$