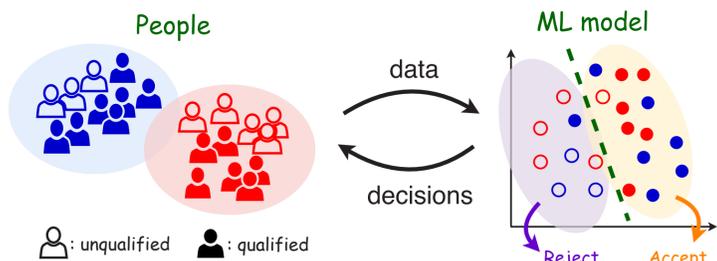# OBJECTIVES

- **Setting:** a decision-maker aims to select people from applicants that are qualified for tasks.
- Impose fairness constraint to make fair decisions (e.g., same acceptance rates across groups)
- Interplay between ML models and people
  - ML decisions affect people's behaviors
  - People generate data for training ML models



People — ML model

data / decisions

👤: unqualified  👤: qualified

Reject  Accept

**Goal:** study the **long-term** impact of the fairness constraints on **qualifications** of different groups.
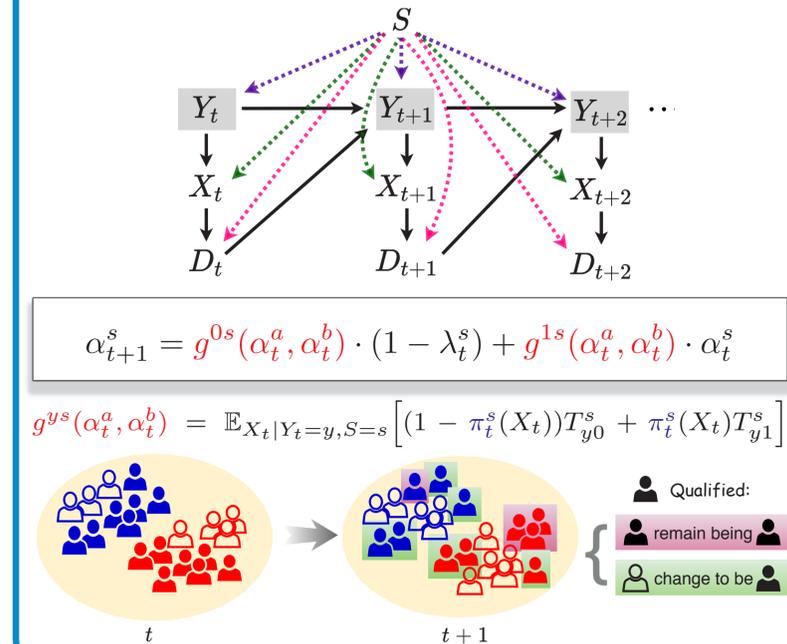
# MODEL

**Two demographic groups $\mathcal{G}_a$, $\mathcal{G}_b$**

- Sensitive attribute $S \in \{a, b\}$
- **Time-varying** feature $X_t \in \mathbb{R}^d$ and qualification state $Y_t \in \{0, 1\}$
  - ***Feature generation process:*** time-invariant $P_{X|Y,S}(x|y,s) = \mathbb{P}(X_t = x|Y_t = y, S = s)$
  - ***Transitions of qualification state:*** time-invariant $T_{yd}^s = \mathbb{P}(Y_{t+1} = 1|Y_t = y, D_t = d, S = s)$
- Qualification rate $\alpha_t^s = P_{Y_t|S}(1|s)$
- Inequality measure: disparity between $\alpha_t^a$ and $\alpha_t^b$

**Myopic decision-maker's optimal fair policies $\pi_t^a, \pi_t^b$**

$\max\limits_{\pi^a, \pi^b} \quad U_t(\pi^a, \pi^b) = \mathbb{E}[R(D_t, Y_t)]$

s.t. $\quad \mathbb{E}_{X_t \sim \mathcal{P}_\mathcal{C}^a}[\pi^a(X_t)] = \mathbb{E}_{X_t \sim \mathcal{P}_\mathcal{C}^b}[\pi^b(X_t)]$

- Unconstrained (UN)
- Demographic Parity (DP): $\mathcal{P}_{DP}^s(x) = P_{X|S}(x|s)$
- Equal Opportunity (EqOpt): $\mathcal{P}_{EqOpt}^s(x) = P_{X|Y,S}(x|1, s)$

- Decision $D_t \in \{0, 1\}$ is based on $\pi_t^s(x) = \mathbb{P}(D_t = 1|X_t = x, S = s)$
- Utility function $R(1, 1) = u_+, R(1, 0) = -u_-, R(0, 1) = R(0, 0) = 0$

# DYNAMICS



$\alpha_{t+1}^s = g^{0s}(\alpha_t^a, \alpha_t^b) \cdot (1 - \lambda_t^s) + g^{1s}(\alpha_t^a, \alpha_t^b) \cdot \alpha_t^s$

$g^{ys}(\alpha_t^a, \alpha_t^b) = \mathbb{E}_{X_t|Y_t=y, S=s}\left[(1 - \pi_t^s(X_t))T_{y0}^s + \pi_t^s(X_t)T_{y1}^s\right]$

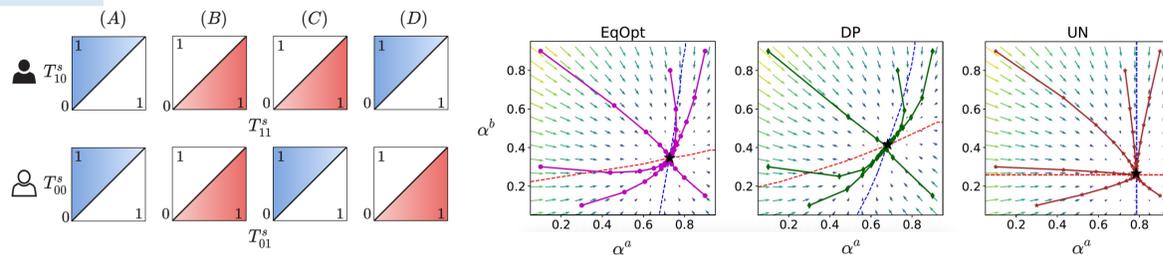Qualified:
👤 remain being 👤
👤 change to be 👤

$t$ → $t+1$

# EQUILIBRIUM ANALYSIS

- **Optimal (fair) policies:** threshold policies are optimal.
- **Existence of equilibrium:** $\forall T_{dy}^s \in (0, 1)$, the dynamics have at least one equilibrium $(\hat{\alpha}^a, \hat{\alpha}^b)$.
- **Uniqueness of equilibrium:** sufficient conditions for the uniqueness of equilibrium under (A)(B).
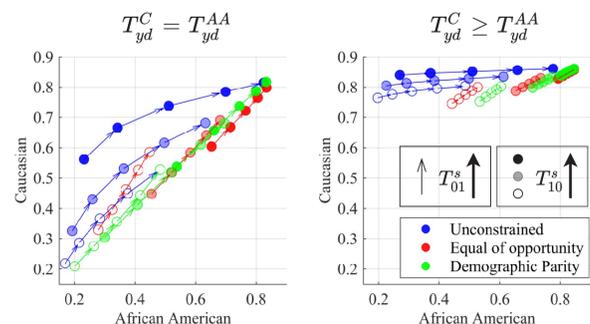
Two effects on people
- "Lack of motivation"
  $T_{y1}^s \leq T_{y0}^s$
- "Leg-up"
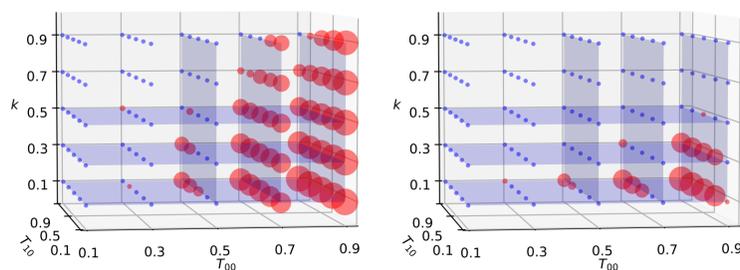  $T_{y1}^s \geq T_{y0}^s$



# NUMERICAL RESULTS

- **FICO score dataset**
  - Effect of transition intervention



$T_{yd}^C = T_{yd}^{AA}$

$T_{yd}^C \geq T_{yd}^{AA}$

↑ $T_{01}^s$ ↑   ○ $T_{10}^s$ ↑

- Unconstrained
- Equal of opportunity
- Demographic Parity

- **COMPAS dataset**
  - Oscillation may happen in the long-run



# LONG-TERM IMPACT OF FAIRNESS CONSTRAINTS

- **Natural equality:** $\forall P_{X|Y,S}$ and $\forall \alpha \in (0, 1)$, $\exists$ transitions $T_{yd}^s$ under (A) or (B) s.t. $\hat{\alpha}_{UN}^a = \hat{\alpha}_{UN}^b = \alpha$.
  - If $P_{X|Y,S=a} = P_{X|Y,S=b}$, then fairness $\mathcal{C} = $ DP or EqOpt **maintains** equality: $\hat{\alpha}_\mathcal{C}^a = \hat{\alpha}_\mathcal{C}^b$
  - If $P_{X|Y,S=a} \neq P_{X|Y,S=b}$, then fairness $\mathcal{C} = $ DP or EqOpt **violates** equality: $\hat{\alpha}_\mathcal{C}^a \neq \hat{\alpha}_\mathcal{C}^b$
- **Natural inequality ($\hat{\alpha}_{UN}^a \neq \hat{\alpha}_{UN}^b$):**
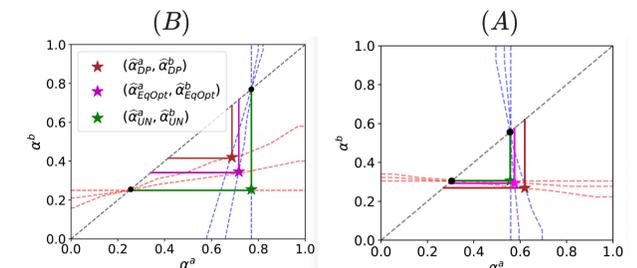
  **Case 1:** due to different transitions
  - Under (A), DP and EqOpt **exacerbate** inequality
  - Under (B), DP and EqOpt **mitigate** inequality
  - Disadvantaged group **remains** being disadvantaged

  **Case 2:** due to different feature that generated
  Under some conditions on $P_{X|Y,S}$, $u_+$, $u_-$ and $T_{yd}^s$ satisfying (B):
  - EqOpt **mitigates** inequality and disadvantaged group **remains** being disadvantaged
  - DP either **mitigates** inequality, or **flips** disadvantaged group



# EFFECTIVE INTERVENTION

- **Policy Intervention:**
  - **Sub-optimal** fair policies can improve $(\hat{\alpha}^a, \hat{\alpha}^b)$.
  - $\exists$ threshold policies s.t. $\hat{\alpha}^a = \hat{\alpha}^b$ as long as $T_{yd}^a$ and $T_{yd}^b$ are not different significantly.
- **Transition Intervention:**
  - Increasing any $T_{yd}^s$ increases qualification rate $\hat{\alpha}^s$.

# CONCLUSIONS

- Construct a POMDP framework for sequential decision-making and analyze its equilibrium.
- Imposing fairness constraints may or may **not** help in promoting long-term equality.
- Importance of understanding real-world dynamics in decision-making systems.