

Striking a Balance in Fairness for Dynamic Systems Through Reinforcement Learning

Yaowei Hu, Jacob Lear, Lu Zhang

Department of Electrical Engineering and Computer Science

University of Arkansas

Fayetteville, Arkansas, USA

{yaowei.hu, jdlear, lz006}@uark.edu

Abstract—While significant advancements have been made in the field of fair machine learning, the majority of studies focus on scenarios where the decision model operates on a static population. In this paper, we study fairness in dynamic systems where sequential decisions are made. Each decision may shift the underlying distribution of features or user behavior. We model the dynamic system through a Markov Decision Process (MDP). By acknowledging that traditional fairness notions and long-term fairness are distinct requirements that may not necessarily align with one another, we propose an algorithmic framework to integrate various fairness considerations with reinforcement learning using both pre-processing and in-processing approaches. Three case studies show that our method can strike a balance between traditional fairness notions, long-term fairness, and utility.

Index Terms—long-term fairness, sequential decision-making, reinforcement learning, proximal policy optimization

I. INTRODUCTION

Machine learning algorithms are increasingly being widely used in high-stakes decision-making applications, such as college admissions [1], [2], bank loans [3], [4], employment [5], recidivism risk assessment [6], etc. However, our society is rife with systemic inequalities, and real-world data is often influenced by social and historical contexts which may contain historical biases and discrimination based on gender, race, and other factors. As a result, algorithmic fairness in machine learning has received increasing attention. Ensuring that machine learning models trained on biased data do not make discriminatory decisions against vulnerable groups is a key consideration for deploying these algorithmic decision-making systems [7]. In the literature, various notions and metrics have been proposed to measure the concept of fairness in machine learning. Diverse algorithms have been proposed to address the fairness concerns. For surveys on general fair machine learning, please refer to [8]–[10].

Despite the progress in fair machine learning, most of the studies focus on scenarios in which the decision model makes decisions over a static population. These works typically employ fair machine learning algorithms to achieve fairness for a batch of data, but do not consider how these decisions may affect fairness over an extended period of time. Nevertheless, real-world decision-making systems usually operate in a dynamic manner that involves making sequential decisions. Each decision may shift the underlying distribution of features or

user behavior, which in turn affects the subsequent decisions. For example, loan decisions have an impact on an individual's credit score, while the allocation of police forces can influence the crime rate in a particular location. As a result, long-term fairness has been proposed as a fairness notion that concerns how fairness can be achieved and maintained over a time horizon, rather than a single time step, by taking into consideration the system dynamics and feedback loops.

How to achieve long-term fairness has been explored in previous studies through explicit modeling of dynamics and feedback loops. One branch of research is to investigate the impact of current decisions on a target population in application-specific scenarios by leveraging various analytical frameworks, such as the one-step feedback [11] or Pólya urn model [12]. There are analytical results that show that simply enforcing traditional fairness notions at each static decision point may produce adverse effects on disadvantaged groups in the long run [13], [14]. Long-term fairness has also been studied in the context of reinforcement learning (RL) where the system dynamics and feedback loops between decisions and the population are formulated through a Markov Decision Process (MDP) [15]. Following this line of research, [16], [17] establish simulation environments for studying long-term fairness in RL, based on which RL algorithms have been used to learn a decision-making policy that aims to optimize both policy utility and fairness as long-term objectives [18], [19]. However, one limitation of existing methods is that they usually formulate long-term fairness as the difference between the average or instantaneous rewards received by different demographic groups similar to static fairness notions, but do not take into account the inherent difference between traditional fairness notions and long-term fairness requirements.

In this paper, we consider long-term fairness as a requirement on the states rather than the rewards following the definitions in some well-accepted research [13], [14], [20]. We acknowledge that traditional fairness notions and long-term fairness are distinct requirements that may not necessarily align with one another. Traditional fairness considers the equity of the outcomes or performance of the decision model at a single decision point. It is referred to as short-term fairness later in this paper for a clear representation. Long-term fairness, on the other hand, refers to a long-term state in which equity is systematically satisfied. Such a state

may be achieved by gradually reducing the gap between the qualification distributions of different groups. As a result, imposing short-term fairness constraints may not necessarily lead to long-term fairness even if the constraints are incorporated into a long-term objective. For example, suppose a bank uses different thresholds for making loan decisions for the advantaged and disadvantaged groups in order to ensure fair outcomes. However, this approach may not help narrow the gap between the credit score distributions of the two groups.

To address the above issue, we develop an algorithmic framework that promotes both short-term and long-term fairness simultaneously. Similar to prior works, we utilize the MDP framework to leverage its power in optimizing long-term objectives. By recognizing the distinct requirements of short-term and long-term fairness, we incorporate them into the RL algorithm using different approaches. Since the concept of long-term fairness is aligned with the principle of the MDP framework, we employ an in-processing approach to deal with this constraint. We adopt the 1-Wasserstein distance as the metric of the distribution gap and theoretically show that minimizing the distance can lead to a long-term fair state. On the other hand, we adopt a model-agnostic pre-processing approach to deal with short-term fairness to ensure that it is enforced throughout the sequential decision-making process. We extend a classic pre-processing approach called massaging [21] to the RL setting by integrating it with the policy optimization algorithm. Finally, we show the exact implementation of our algorithmic framework using three case studies, where the experimental results demonstrate that our method is capable of striking a desired balance between short-term fairness, long-term fairness, and the utility of the sequential decision system.

We summarize our contributions as below:

- We propose to achieve systematic equity in sequential decision-making by considering short-term and long-term fairness as distinct fairness requirements.
- We develop an efficient and flexible algorithmic framework that integrates short-term and long-term fairness with the MDP framework as distinct constraints.
- Three case studies within simulation environments are used to prove the effectiveness of our method by evaluating the performance of our method and comparing it with the state-of-the-art baselines.

II. RELATED WORK

Extensive research has been conducted on achieving fairness in machine learning within static systems, but the exploration of fairness in dynamic systems has remained relatively limited. It has been first studied in a compound decision-making process called a pipeline [22], [23] where individuals may drop out at any stage and classification in subsequent stages depends on the remaining cohort of individuals. Then, more complex dynamic systems are considered where the sequential decisions have the potential to reshape the distribution of the underlying population. Due to the difficulty of modeling the system dynamics, researchers first focused on specific

applications. For example, the authors in [20] propose a short-term intervention for long-term fairness in the labor market by constructing a dynamic reputational model and adopting a dual labor market composed of a Temporary Labor Market and a Permanent Labor Market. The research in [13] studies the delayed impact of fair machine learning in lending scenarios. It has been realized that, traditional fairness notions cannot guarantee to promote fairness as a delayed impact, even in a one-step feedback model. Researchers in [14] also verify the inadequacy of traditional fairness notions in achieving fairness in the long run. It has also been proposed to model the dynamic system using causal graphs. [24] introduces causal directed acyclic graphs (DAGs) as a framework for studying fairness in dynamical systems. The research in [25], [26] proposes to treat model deployment as soft interventions and infer post-intervention distributions in measuring fairness.

As Markov Decision Process (MDP) is a prevalent framework used to model sequential decision-making, it has also been proposed to study fairness in the context of reinforcement learning (RL). The research in [15] initiates the study of fairness in RL and proposes a fairness constraint that requires an algorithm to never prefer one action over another if the long-term reward of choosing the latter action is higher. However, fairness is defined only based on the reward of each action without considering the demographic information. Along this direction, the authors in [27] study long-term fairness and formulate a Constrained Markov Decision Process (CMDP), but fairness is defined only within the context of recommendation systems. The authors in [19] propose two algorithms for learning policies to satisfy fairness constraints that are defined on the average reward for individuals in different groups. In particular, [16], [17] advocate for the use of simulation studies to understand the long-term behaviors of deployed ML-based decision systems and their potential consequences. The authors explore toy examples of dynamical systems including bank loans, college admissions, allocation of attention, and epidemic control. An extensible open-source software framework for implementing fairness-focused simulation studies is provided. Based on that, [18] proposes to impose fairness requirements in policy optimization by regularizing the advantage evaluation of different actions. The proposed methods make it easy to impose fairness constraints without rewarding engineering or sacrificing training efficiency. However, the authors directly apply traditional fairness notions without considering the different requirements of traditional and long-term fairness. Different from prior works, we consider long-term fairness as a requirement on the states rather than on the rewards and acknowledge traditional and long-term fairness as distinct fairness requirements. We develop an algorithmic framework that promotes both requirements simultaneously.

III. PRELIMINARIES

This section introduces the background and preliminaries of fair machine learning, reinforcement learning (RL), and proximal policy optimization (PPO), a prevalent policy optimization algorithm in RL.

A. Fair Machine Learning

The issue of fairness has become one of the most popular topics in machine learning in recent years. To measure fairness in algorithmic decision-making, a large number of fairness notions have been proposed in the literature. Typical examples include *demographic parity (DP)* and *equal opportunity (EO)*. DP aims to ensure that different demographic groups are represented proportionally in the outcomes of a decision model. EO, on the other hand, refers to the principle of treating individuals or groups fairly by ensuring equal error rates or predictive performance across different demographic subgroups. Then, to address the fairness issues, bias mitigation algorithms are proposed mainly from three perspectives: pre-processing, in-processing, and post-processing. Pre-processing approaches focus on eliminating bias from the training data (e.g., [28]–[30]), in-processing approaches aim to avoid introducing bias in model training by proposing new model structures or loss functions (e.g., [31]–[33]), and post-processing approaches modify predicted outcomes to resolve fairness issues (e.g., [34]–[36]).

B. Reinforcement Learning

Reinforcement learning consists of two interactive components, an agent and an environment, which interact with each other over time. This interaction process is modeled as a Markov decision process (MDP) [37]. An MDP is denoted by a tuple $M = (\mathcal{S}, \mathcal{A}, P, R, \rho_0, \gamma)$, where $\mathcal{S} \in \mathcal{S}$ is a set of states, $\mathcal{A} \in \mathcal{A}$ is a set of actions, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a transition function that represents the probability of next state given the current state and the action, $R : \mathcal{S} \rightarrow \mathbb{R}$ is a reward function, $\rho_0 : \mathcal{S} \rightarrow [0, 1]$ is an initial state distribution, and $\gamma \in [0, 1]$ is a discount factor. At each time step t , the agent observes a state $s_t \in \mathcal{S}$ from the environment and takes an action $a_t \in \mathcal{A}$ following a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ based on the current state. Then the agent observes a new state $s_{t+1} \in \mathcal{S}$ and a reward $r_t \in R$ generated by the environment with the transition probability $P(s_{t+1}|s_t, a_t)$ and the reward function $R(s_t)$ in the next time step $t + 1$. The goal of RL is to learn a optimal policy π_θ which maximizes the expected discounted cumulative rewards, defined as below.

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \right],$$

where $\tau = (s_0, a_0, s_1, a_1, \dots)$ is a trajectory and $\tau \sim \pi_\theta$ means that a trajectory τ is sampled from the policy π_θ following $s_0 \sim \rho_0(s_0)$, $a_t \sim \pi(a_t|s_t)$, $s_{t+1} \sim P(s_{t+1}|s_t, a_t)$.

To address the RL problem, there are several concepts that are often involved in RL algorithms. Let $R(\tau)$ denote the discounted cumulative rewards of a trajectory τ . The state value function V and the state-action value function Q are given by $V(s_t) = \mathbb{E}_{\tau \sim \pi} [R(\tau)|s_t = s]$ and $Q(s_t, a_t) = \mathbb{E}_{\tau \sim \pi} [R(\tau)|s_t = s, a_t = a]$, which evaluate how good a state or a pair of state and action is. The advantage function is the difference between $Q(s_t, a_t)$ and $V(s_t)$, i.e., $A(s_t, a_t) =$

$Q(s_t, a_t) - V(s_t)$, and it can be considered as the advantage of taking a given action over following the policy [38].

C. Proximal Policy Optimization

Policy optimization methods [39] are a type of reinforcement learning algorithms that improve policies directly by estimating policy gradients and optimizing with stochastic gradient ascent. The most commonly used form of gradient estimator is given by

$$\nabla J(\theta) = \mathbb{E}_{(s_t, a_t) \sim \pi_\theta} [A(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t|s_t)] \quad (1)$$

where the expectation is estimated over a batch of samples.

Proximal policy optimization (PPO) [40] is a state-of-the-art policy optimization algorithm stemming from the trust region policy optimization (TRPO) algorithm [41]. It maximizes a clipped surrogate function to prevent the gradient update dramatically, as follows

$$J^{CLIP}(\theta) = \hat{\mathbb{E}} [\min(r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t)] \quad (2)$$

where A_t is short for $A(s_t, a_t)$, $r_t(\theta)$ denotes the probability ratio $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$, and ϵ is a hyperparameter. To compute variance-reduced advantage function, a neural network is used to estimate the state value function $V(s_t)$ with the squared-error loss

$$L^V(\theta) = \mathbb{E}[(V_\theta(s_t) - R(\tau))^2] \quad (3)$$

The clipped surrogate function restricts the magnitude of the gradient update, which not only makes the algorithm more stable, but also allows for multiple updates using a batch of samples, improving the data efficiency.

In this paper, we adopt PPO as the RL algorithm, but our method can be applied to any policy optimization algorithms.

IV. PROBLEM FORMULATION

To develop an RL algorithm for achieving long-term fairness, we start by defining fairness notions in the context of sequential decision-making, presenting the problem formulation, and establishing the setting of the fair RL learning problem.

A. Fairness Definition for Sequential Decision Making

As stated in the last section, traditional fairness notions are usually concerned about whether the machine learning model produces the same outcome or performance across different groups in a static population. In the context of sequential decision-making, we refer to this type of notions as short-term fairness notions, which are defined based on a cohort of individuals who participate in the decision-making system over a specific period.

Definition 1 (Short-term Fairness): In a sequential decision-making system, short-term fairness is defined as the equal outcome or performance of the decision model/policy over a participating cohort.

It is worth emphasizing that, first, short-term fairness notions may be enforced by laws and regulations in some domains, such as the U.S. Equal Employment Opportunity Commission [42] that prohibits employment discrimination.

So, it is essential to enforce short-term fairness throughout the sequential decision-making process including both training and evaluation. Second, different short-term fairness notions may conflict and not be achieved simultaneously if historical and/or systemic biases exist [10]. For example, achieving demographic parity may require preferential treatment to account for historical disadvantages, which could potentially impact equal opportunity. Thus, we adopt a single short-term fairness notion only in our algorithm.

On the other hand, long-term fairness has been proposed to account for fairness and equity of the sequential decision-making system in the long run [13]. The general goal of long-term fairness is to reach a state where the historical disadvantages are rectified and systemic biases are removed. Since the trade-off between short-term and long-term fairness is due to historical disadvantages and systemic biases, we presume that long-term fairness also implies a state in which it becomes easier to simultaneously satisfy different short-term fairness notions.

How to quantify long-term fairness? In the literature, features are regarded as indicators or metrics that assess the qualification or competency levels of individuals. Then, long-term fairness is often formulated by measuring the gap in feature distributions between different groups. For example, in [13], the difference in feature distribution of the disadvantaged group between the starting time step and the ending time step, i.e., $\Delta = \mathbb{E}[x_{t=t^*}|c^-] - \mathbb{E}[x_{t=0}|c^-]$, is defined as the measure of long-term fairness. It is called long-term improvement if $\Delta > 0$, stagnation if $\Delta = 0$, and active harm if $\Delta < 0$. In [20], long-term fairness is defined as the parity in the feature distribution between the advantaged and disadvantaged group when the dynamic system is at equilibrium, i.e., $\Delta = |\mathbb{E}[x_{t=t^*}|c^+] - \mathbb{E}[x_{t=t^*}|c^-]|$. In this paper, we adopt a similar philosophy to [20] and define long-term fairness as follows.

Definition 2 (Long-term Fairness): Long-term fairness is defined as the equal feature distributions of different groups at a long-term state of the sequential decision-making system.

After discussing the two types of fairness we consider, we define our problem formulation as follows.

Problem Formulation 1: Consider a sequential decision-making context. A policy for making the decision is learned through an iterative process of interaction with the environment. Our goal is to learn a fair policy such that: (1) short-term fairness is guaranteed throughout both the training and evaluation processes, and (2) a long-term fair state is reached at the end of the evaluation process.

B. Problem Setting for Fair RL

To formulate the problem of long-term fair sequential decision-making, we consider a finite-horizon RL problem and leverage the MDP framework. Specifically, in our context states $\mathcal{S} = \mathcal{C} \times \mathcal{X}$ where \mathcal{C} is the domain of the sensitive feature and $\mathcal{X} \in \mathbb{R}^m$ represents the domain of the profile features. When \mathcal{C} is a binary domain, we use $\{c^+, c^-\}$ to represent the advantaged and disadvantaged groups respectively. Let \mathcal{A} denote the action space. If \mathcal{A} is a binary domain, we use

$\{1, 0\}$ to represent the positive and negative actions. Our goal is to learn a stochastic policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ which maximizes the agent's cumulative reward while satisfying certain fairness criteria. Note that we generally allow the sensitive attributes to be involved in the policy input and will explicitly adopt constraints to enforce fairness requirements.

We may use the bank loan system as an example to illustrate this problem setting in a real-world scenario. In this example, the bank is treated as the agent, and the population of the applicants is treated as the environment. Assume that the sensitive attribute C is the race of the applicant. As an RL process, at each time step t , an individual from a certain race group is sampled and applies for a loan from the bank, whose state s_t is given by race c_t and personal profile x_t . Then, the bank runs a policy function π_θ to decide whether to approve the loan according to probability distribution $\pi_\theta(a_t|s_t)$. If the loan is approved, depending on whether the individual repays the loan, both the bank's profit and the feature distribution of the population will be affected which determines the new state and the reward. Finally, the goal of the bank is to learn from interactions a lending policy that maximizes its own profit while satisfying fairness.

V. ALGORITHM

In this section, we develop a flexible and effective fair RL algorithm that integrates the pre-processing and in-processing approaches to promote short-term and long-term fairness simultaneously. We assume separate training and evaluation processes where the policy is updated according to the interaction experience with the environment during the training process, and it is further evaluated in a separate environment with the same dynamics after the training completes. However, our method is readily applied to the setting where the evaluation is conducted during the training. In the following, we first provide an overview of the algorithm.

A. Overview

When developing the fair RL algorithm, it is critical to consider the different requirements for long-term and short-term fairness. Long-term fairness is a state that represents the maximization of the equity of the system in the long run. As it is aligned with the principle of the MDP framework that maximizes the expected total reward over time, we adopt an in-processing approach that regularizes the reward to incorporate the long-term fairness objective so that the RL algorithm can be aware of the long-term fairness status in training. Specifically, we regularize the advantage function of a policy optimization algorithm as:

$$A^\lambda(s_t, a_t) = A(s_t, a_t) + \lambda R(s_t), \quad (4)$$

where $R(s_t)$ is the regularization that reflects the long-term fairness requirements and λ is a hyperparameter that controls the degree of regularization.

On the other hand, short-term fairness is an instantaneous constraint that may be enforced by laws or regulations at

every step of the decision-making. Merely incorporating short-term fairness through regularizing the advantage function may not be sufficient to guarantee short-term fairness throughout the entire training process. Thus, we choose to adopt a pre-processing approach to address this issue. Inspired by the classic pre-processing approach named massaging [21], we propose a method called *action massaging* that selectively alters unfair actions produced by the policy network to fair ones. Specifically, after an action a_t is sampled from the policy network π_θ at time step t , we employ a functional mapping $a'_t = m(s_t, a_t)$ where a'_t may or may not be equal to a_t to generate the trajectory. By using this altered trajectory to perform policy optimization, the policy gradient becomes

$$\nabla J(\theta) = \mathbb{E}_{(s_t, a_t) \sim \pi_\theta} [A^\lambda(s_t, a'_t) \nabla_\theta \log \pi_\theta(a'_t | s_t)], \quad (5)$$

which shows that the policy gradient is computed based on the trajectory formed by a'_t . The rationale of the action massaging is to perform fair actions when the policy network generates biased ones and also encourage the policy to generate fair and high-reward actions. Note that this approach differs from the off-policy RL learning algorithm which optimizes the current policy network based on the trajectories generated by a different policy, and hence our approach does not require importance sampling to correct for the bias.

Next, we describe the above pre-processing and in-processing components in detail.

B. Action Massaging for Short-term Fairness

The action massaging altered actions according to a pre-defined short-term fairness criterion that is to meet legal and regulatory requirements for decision-making and prevent discrimination against certain groups. As mentioned earlier, we consider group fairness notions such as DP or EO in our work. At each time step t , short-term fairness relies on the current state and action as well as the past states and actions of the system. To facilitate computation, we adopt a sliding window w so that all the states and actions between time step $t-w$ and t form a participating cohort that will be used to measure short-term fairness. We denote this measure as $\Delta_s(s_t, a_t)$. Then, the action massaging $m(s_t, a_t)$ alters the action a_t to a'_t that minimizes $\Delta_s(s_t, a_t)$ to improve short-term fairness.

When designing the mapping $m(s_t, a_t)$, one principle is that the modifications should not significantly damage the utility of the policy. In our method, we treat $\pi_\theta(a_t | s_t)$ as the confidence level for selecting action a_t , and introduce a constraint that limits the difference between the confidence of the original action a_t and the altered action a'_t . The action massaging only alters the current action to a different action when the above difference is smaller than a predefined threshold. When multiple actions satisfy the constraint, the action massaging chooses the one that leads to best short-term fairness. As a result, the action massaging is formulated as follows:

$$\begin{aligned} m(s_t, a_t) &= \arg \min_{a'_t \in \mathcal{A}} \Delta_s(s_t, a'_t) \\ \text{s.t. } & |\pi_\theta(a_t | s_t) - \pi_\theta(a'_t | s_t)| < \tau. \end{aligned} \quad (6)$$

The constraint in Eq. (6) is an important factor that reflects the trade-off between short-term fairness and utility. It aids in enhancing short-term fairness while minimizing the impact on utility. On one hand, the constraint reduces the number of modifications made by the action massaging, as a large number of modifications will cause instability and deviation of training. On the other hand, the constraint also restricts the modifications to be carried out when current actions have low confidence and hence leads to a smaller reduction in utility. For example, in a special case of binary actions (e.g., the decision of bank loan), the constraint will restrict the modifications to actions with confidence close to 0.5. The exact implementation of the action massaging is task-specific and varies with applications, as will be shown in the case studies in the next section.

C. Advantage Regularization for Long-term Fairness

The action massaging for short-term fairness is not enough to achieve long-term fairness as the objective of short-term fairness may not exactly align with the objective of long-term fairness. As mentioned above, we leverage the MDP's capacity to maximize long-term returns as a means to attain long-term fairness. To quantify long-term fairness, in our work we employ the 1-Wasserstein distance between the distributions of different groups as the long-term fairness metric. The reason is presented in the following proposition.

Proposition 1: Denote by d the 1-Wasserstein distance between the feature distributions of different groups, i.e., $d = W(P(x|c^+), P(x|c^-))$. For any decision model $h : \mathcal{X} \mapsto \mathcal{A}$ that is Lipschitz continuous, its DP is bounded by $l_h \cdot d$ where l_h is the Lipschitz constant of h . If we assume that the true label is given by another decision model $g : \mathcal{X} \mapsto \mathcal{A}$ that is Lipschitz continuous and also satisfies DP, then the EO of h is bounded by $(l_h + l_g)/P(y) \cdot w$ where l_g is the Lipschitz constant of g .

Proof: According to the definition of DP, we have

$$\text{DP}(h) = |\mathbb{E}[h(x)|c^+] - \mathbb{E}[h(x)|c^-]|.$$

Due to the Kantorovich–Rubinstein duality [43], it is straightforward that

$$\begin{aligned} \text{DP}(h) &\leq \sup_{\|h\| \leq l_h} [\mathbb{E}_{x \sim P(x|c^+)}[h(x)] - \mathbb{E}_{x \sim P(x|c^-)}[h(x)]] \\ &= l_h \cdot W(P(x|c^+), P(x|c^-)) = l_h \cdot d. \end{aligned}$$

On the other hand, we have

$$\text{EO}(h) = |\mathbb{E}[h(x)|a = 1, c^+] - \mathbb{E}[h(x)|a = 1, c^-]|.$$

Due to the assumption that the true label is given by g and g satisfies DP, it follows that

$$\begin{aligned} \mathbb{E}[h(x)|a, c] &= \int_{\mathcal{X}} h(x) P(x|a, c) dx \\ &= \int_{\mathcal{X}} h(x) P(x|c) \frac{P(y|x, c)}{P(y|c)} dx = \int_{\mathcal{X}} h(x) P(x|c) \frac{g(x)}{P(y)} dx \\ &= \frac{1}{P(y)} \mathbb{E}_{x \sim P(x|c)} [h(x)g(x)]. \end{aligned}$$

In addition, define $f(x) = h(x)g(x)$ and denote the Lipschitz constant of f as l_f . It is easy to show that $l_f \leq l_h \cdot \sup_x |h(x)| + l_g \cdot \sup_x |g(x)|$. Since $h(x) \leq 1$ and $g(x) \leq 1$, we have $l_f \leq l_h + l_g$. As a result, we have

$$EO(h) \leq \frac{l_h + l_g}{P(y)} W(P(x|c^+), P(x|c^-)) = \frac{l_h + l_g}{P(y)} \cdot d.$$

The above proposition shows that, by approaching a long-term state where the 1-Wasserstein distance between the feature distributions of different groups is minimized, we can mitigate at that state both the DP and EO of any decision model that is Lipschitz continuous. This implies that a long-term fair state has been reached.

Denote the long-term fairness measure computed at time step t as $\Delta_l(s_t)$. Similar to short-term fairness, we adopt a sliding window to form a participating cohort for estimating the feature distributions. Then, we incorporate $\Delta_l(s_t)$ into the advantage function as the regularization. However, rather than directly adding $\Delta_l(s_t)$ to the advantage function, we further consider the trade-off between short-term fairness and long-term fairness when the two objectives are not aligned. Specifically, we promote the advantage when both short-term and long-term fairness can be improved while demoting the advantage when both short-term and long-term fairness is damaged. When there is a conflict between short-term and long-term fairness, we keep the current advantage unchanged. The regularization term is defined as follows:

$$R(s_t) = \begin{cases} \min(0, \Delta_l(s_t) - \Delta_l(s_{t+1})) & \Delta_s(s_t, a_t) > \delta \\ \max(0, \Delta_l(s_t) - \Delta_l(s_{t+1})) & \Delta_s(s_t, a_t) \leq \delta \end{cases} \quad (7)$$

As can be seen, the first term is active when short-term fairness $\Delta_s(s_t, a_t)$ is larger than the threshold δ . Then, this term will penalize the original advantage when the long-term fairness measure does not decrease at the next time step $t+1$ compared to the current time step t . The second term is active when short-term fairness $\Delta_s(s_t, a_t)$ is less or equal to δ . In this case, we reward the advantage function when the long-term fairness measure reduces. This approach allows for a continuous improvement of long-term fairness throughout the entire sequence, rather than a sudden change at a specific point, while it remains simple and effective.

Combining the above two methods for fairness, we present the pseudocode of our algorithm in Algorithm 1, which is referred to as F-PPO.

VI. EXPERIMENTS

For demonstrating the performance of our proposed method, we make use of the simulation environments [16], [17] that implement toy examples of dynamic systems for supporting studies of long-term consequences of ML-based decision systems. We conduct three case studies in the context of bank loans, allocation of attention, and epidemic control. The proposed method is evaluated with respect to utility, short-term fairness, and long-term fairness. As mentioned earlier, the policies will be first trained by interacting with the environment and then

Algorithm 1: Fair Proximal Policy Optimization (F-PPO)

```

1 Initialize policy network  $\pi_\theta$  and value function network  $v_\phi$ ;
2 for  $k = 0, 1, 2, \dots$  do
3   Collect trajectories  $\mathcal{D}_k$  from policy  $\pi_\theta$  where actions  $a_t$  are sampled from  $\pi_\theta(a_t|s_t)$ ;
4   Compute  $\Delta_s(s_t, a_t)$  and apply action massaging according to Eq. (6) to obtain altered trajectories;
5   Compute  $\Delta_l(s_t)$  and penalized advantage  $A^\lambda(s_t, a_t)$  according to Eqs. (4) and (7);
6   Update the policy by maximizing the clipped surrogate function  $J^{CLIP}$  according to Eq. (2);
7   Update the value function network by minimizing the squared-error loss  $L^V$  according to Eq. (3);
8 end

```

tested separately in the environment. All the code and data are available online at <https://github.com/yaoweihu/Fairness-in-RL>.

Baselines. We consider two different categories of baseline agents in our experiments. The first category is human-designed policy agents used in [16], [18], including the EO agent for bank loans, the CPO agent for attention allocation, and the Max agent for epidemic control. The second category consists of learning-based policy agents. We consider the original PPO algorithm that only maximizes the cumulative reward, and the A-PPO algorithm that is the state-of-the-art fair RL algorithm proposed in [18] for achieving fairness through advantage regularization. As an ablation study, we also consider a variant of our method named F-PPO-L that only consists of the long-term fairness component but with the short-term fairness component removed.

A. Case Study: Bank Loans

In this case study, the bank lending scenario is simulated where an agent plays the role of a bank to make decisions about whether to grant loans to a stream of applicants. The qualification of applicants is described by a discrete credit score, which changes with the loan decisions.

Environment. In this environment, at each time step the bank observes a loan applicant s_t which is sampled with replacement from the pool of applicants. Each applicant consists of a credit score (qualification feature) and a group membership (sensitive feature). The group membership c_t is uniformly sampled from $\{c^+, c^-\}$. The credit score x_t , on the other hand, is drawn from a group-dependent discrete distribution over $X \in \{1, 2, \dots, X_{max}\}$ with $\mathbb{E}_{c^-}[X] < \mathbb{E}_{c^+}[X]$. The bank employs a policy to make binary decisions of whether to deny or approve loan applications. If a loan applicant receives a loan and defaults, his/her credit score C_i drops, which is simulated in the distribution by moving a small portion of mass from $P_c(X_i)$ to $P_c(X_{i-1})$ if $i \neq 1$. Similarly, if the loan applicant receives a loan and repays, a small portion of mass will be moved from $P_c(X_i)$ to $P_c(X_{i+1})$ if $i \neq X_{max}$. There is no change in the distribution if the applicant does not receive the

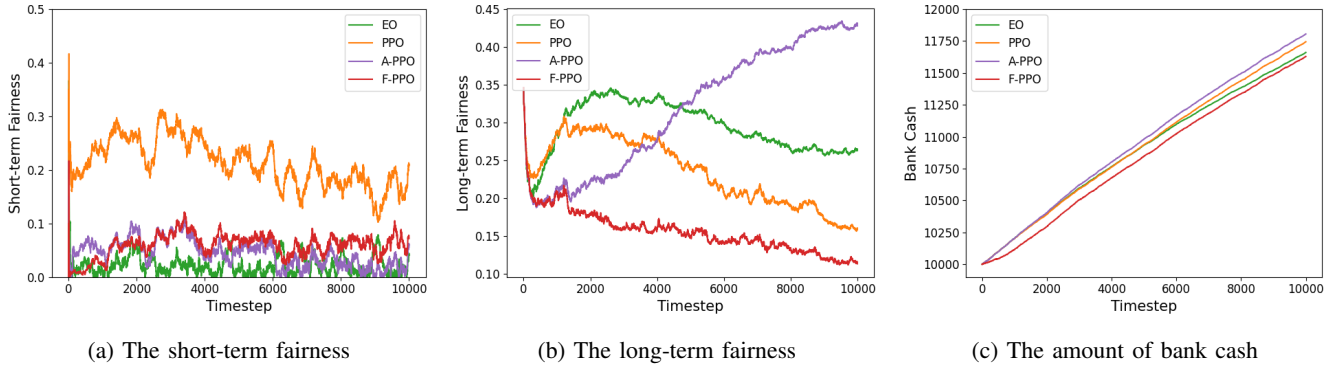


Fig. 1: Experimental results for bank loans. The recorded values are averages over 10 evaluation runs.

loan. The bank's profit increases by the amount of the loan plus the interest on successfully repaid loans, and decreases by the loan amount on defaults. The probability of default is given by a deterministic function of the credit score. The reward of the bank at each time step is defined as the change in its profit at the next time step.

Implementation of F-PPO. In this case study, we implement a policy network $\pi_\theta(a_t|s_t)$ to make binary decisions $a_t \in \{1, 0\}$. We adopt EO as the short-term fairness notion, which is defined as follows:

$$\Delta_s(s_t, a_t) = \left| \frac{\sum_{t-w}^t \text{successful_loan}_{tc^+}}{\sum_{t-w}^t \text{will_repay}_{tc^+}} - \frac{\sum_{t-w}^t \text{successful_loan}_{tc^-}}{\sum_{t-w}^t \text{will_repay}_{tc^-}} \right|, \quad (8)$$

where $w = 300$ is the sliding window size. For the action massaging, the action will be flipped if the alternative action is fairer in terms of short-term fairness and the confidence of the action is lower than the threshold. The threshold in Eq. (6) is dynamically adjusted according to the number of training iterations. The idea is to perform a cold start in action massaging so that the actions are not altered at the beginning of training. The threshold is initially zero and increases after a certain number of iterations. Specifically, the threshold τ at the i th iteration is defined as $\tau = 1 - 2\tau(i)$ where

$$\tau(i) = \begin{cases} \tau_s \cdot \gamma^{i-i_s} & i \geq i_s \\ 0.5 & \text{otherwise} \end{cases}$$

In our experiments, we set $\tau_s = 0.5$, $i_s = 17$, and $\gamma = 0.985$. Finally, long-term fairness is computed as

$$\Delta_l(s_t) = W(P_{t-w:t}(X|c^+), P_{t-w:t}(X|c^-))$$

where $P_{t-w:t}(X|c)$ is the distribution mass of the credit score of group c measured within the sliding window. For other hyperparameters, we set $\lambda = 1$ and $\delta = 0.05$.

Agents. We include PPO, A-PPO, and EO as baselines to compare with our F-PPO, where EO is the agent that maximizes the bank profits subject to constraints of equal opportunity at every time step.

Results. The short-term fairness, long-term fairness, and reward obtained by different agents during the test are shown in Fig. 1. As can be seen, despite similar performance in terms of reward (Fig. 1(c)), the fairness performance of different methods is diverse. For short-term fairness (Fig. 1(a)), EO, A-PPO and F-PPO are all able to keep the bias values below 0.1, while the original PPO produces much larger bias values. For long-term fairness (Fig. 1(b)), our F-PPO approach exhibits superior performance among all the methods considered, as it consistently achieves the smallest bias values, and these bias values continue to decrease over time. On the other hand, A-PPO produces the worst performance where the bias values increase over time. This result shows that simply adding traditional fairness constraints into a long-term objective does not necessarily achieve long-term fairness. By combining the three results, we see that our F-PPO algorithm strikes a desirable balance between short-term fairness, long-term fairness, and utility of the policy.

For the ablation study, Fig. 2(a) shows the mean and standard deviation of short-term fairness achieved by F-PPO and F-PPO-L in each iteration of the training process over 350 iterations. The results show the effectiveness of the action massaging in maintaining short-term fairness throughout both the training and evaluation processes, while the long-term fairness component alone cannot guarantee short-term fairness.

B. Case Study: Attention Allocation

This scenario aims to simulate incident monitoring and mitigation. In the simulation, the agent's role is to assign attention units to a set of locations. Each attention unit can prevent, or catch, one incident at the location to which it is assigned. The incident rates at each location vary over time in accordance with the number of incident occurrences as well as the agent's decisions on how to assign attention units.

Environment. In the attention allocation environment, let N represent the number of attention units, and K be the number of locations. At each time step t , the agent assigns all N units over the K locations and $a_{k,t}$ denotes the number of units assigned to location k . The number of incidents that occur at each location is sampled from a Poisson distribution

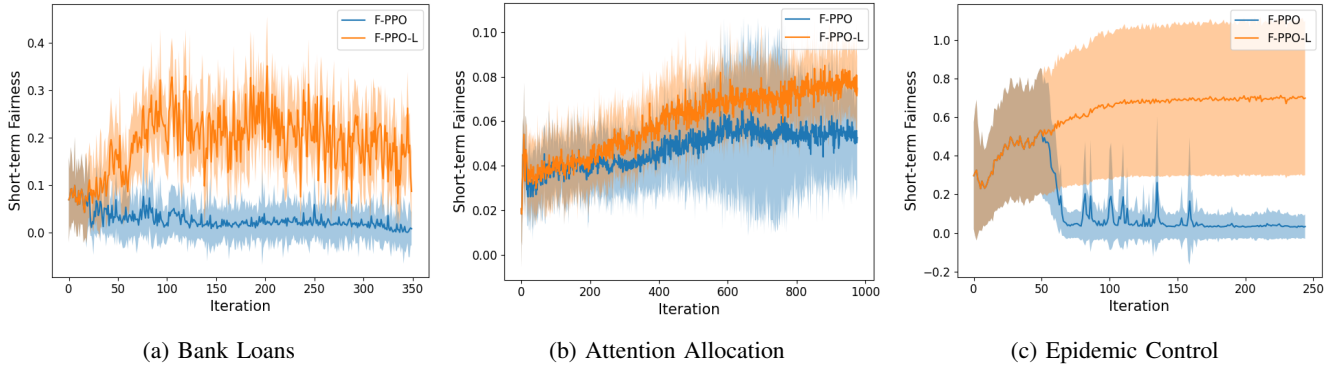


Fig. 2: Ablation study: mean and standard deviation of short-term fairness in each iteration measured during training.

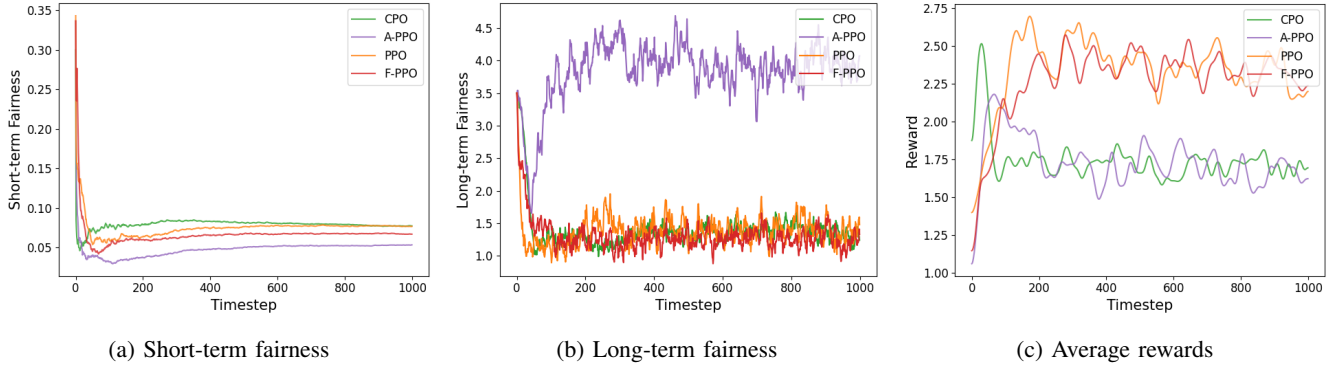


Fig. 3: Results for the Attention Allocation environment. The recorded values are the averages over 10 evaluation episodes.

as $y_{k,t} \sim \text{Poisson}(\Lambda_{k,t})$, where $\Lambda_{k,t}$ is a dynamic parameter which changes according to

$$\Lambda_{k,t+1} = \begin{cases} \Lambda_{k,t} + \lambda_k^I & \text{if } a_{k,t} = 0 \\ \Lambda_{k,t} - \lambda_D \cdot a_{k,t} & \text{otherwise.} \end{cases}$$

Here, λ_k^I is the increase rate, which may vary between locations k , and λ_D is the decrease rate which is the same across locations. The number of incidents discovered at a location is given by $\hat{y}_{k,t} = \min(a_{k,t}, y_{k,t})$. The reward is defined as $r(s_t) = \zeta_0 \sum_{k=1}^K \hat{y}_{k,t} - \zeta_1 \sum_{k=1}^K (y_{k,t} - \hat{y}_{k,t})$ which is determined by the fraction of incidents discovered. The parameters ζ_0 and ζ_1 weight the reward function in terms of incidents discovered versus incidents missed. The state s_t is an observation history of length H and each observation is a tuple of vectors $(\hat{y}_t, y_t, a_t, \hat{y}_t \odot y_t)$, where \odot denotes the Hadamard division operation.

Implementation of F-PPO. The policy network for attention allocation produces a K dimensional vector of logits which are converted into a probability distribution $P(k)$ using the softmax function. The action is constructed by iteratively assigning attention units to the locations until all have been assigned. In each iteration, one attention unit is assigned to the location with the highest probability, from which the amount of $\frac{1}{N}$ is removed before the next iteration. For short-term fairness,

we adopt DP as the metric defined as follows:

$$\Delta_s(s_t, a_t) = \max_k \left| \frac{\sum_{t'=t-w}^t a_{k,t'}}{N \cdot w} - \frac{1}{K} \right|, \quad (9)$$

which requires that the number of units assigned should be equal across different locations. The action massaging checks for each pair of locations k_1, k_2 where at least one unit is assigned to k_1 , if reallocating one unit from k_1 to k_2 would improve short-term fairness. To minimize the impact of the action massaging on the utility, the algorithm also checks if the difference between $P(k_1)$ and $P(k_2)$ is less than the threshold. When both conditions are met by multiple pairs, the algorithm selects the one that leads to the best short-term fairness performance and performs the reallocation. For simplicity, we use a static threshold of 0.08. Finally, we measure long-term fairness according to the incident distribution over all locations. In training, we estimate the incident distribution based on the number of incident occurrences, but in evaluation, we use $\Lambda_{k,t}$ as the ground truth of the incident distribution.

Agents. For baselines used to evaluate our F-PPO agent, we consider PPO, A-PPO, as well as the CPO agent that aims to discover the most incidents.

Results. The experimental results are shown in Fig. 3. Our F-PPO achieves the best long-term fairness performance while maintaining relatively low short-term bias values and high rewards. As a comparison, although A-PPO produces the

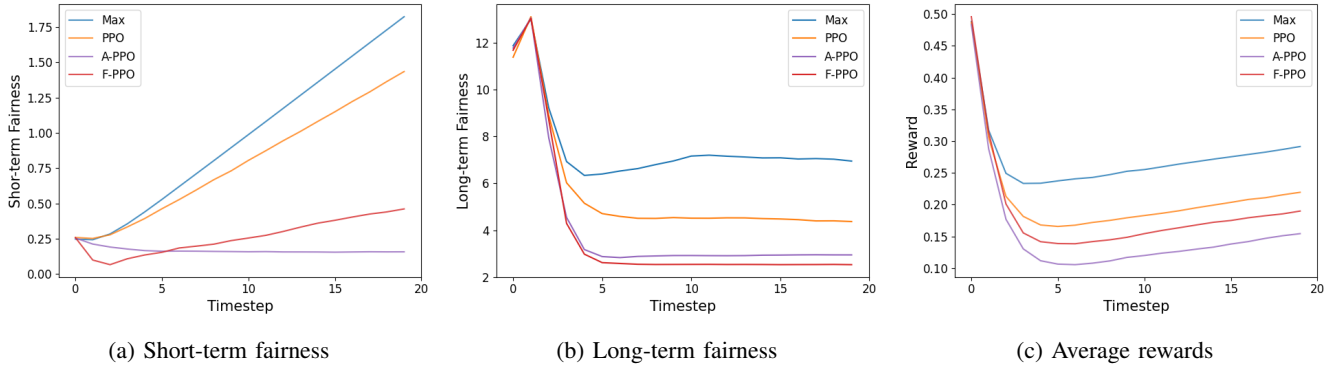


Fig. 4: Experimental results for epidemic control. The recorded values are averages over 200 evaluation episodes.

best short-term fairness performance, its long-term fairness performance and utility are among the worst. For the ablation study, Fig. 2(b) shows the effectiveness of the short-term fairness component.

C. Case Study: Epidemic Control

The third case study simulates an infectious disease scenario where vaccines are allocated within a social network in a step-by-step manner. At each step, one individual is selected by the policy to receive the vaccine. Meanwhile, healthy individuals have the chance to get infected, and sick individuals have the chance to recover. The task is to optimally allocate vaccines in order to mitigate the spread of disease.

Environment. In this environment, we have a social network $G = (V, E)$ where V is a set of individuals and E is a set of edges representing social connections. The state of the environment is a vector of the health state of all individuals. For each individual, the health state is represented by a four-dimensional one-hot encoding $H = \{S, I, R\}$ that represents the three possible states that an individual can be in, including susceptible (healthy), infected, and recovered. At the initial time step, a random set of individuals V_0 are infected. Then, at each time step t , the probability of a susceptible individual i transitioning to the state of infected is given by $P_I(v_{i,t}) = 1 - (1 - \tau)^{|N_I(v_{i,t})|}$, where $N_I(v_{i,t})$ represents the number of infected individuals in the neighbor of individual i and $\tau \in [0, 1]$ is an infectious factor. Meanwhile, the probability of an infected individual recovering is given by $P_R(v) = \rho$ where $\rho \in [0, 1]$ is the recovering factor. If a susceptible individual receives the vaccine, his/her state directly transitions from S to R. The reward is defined as the proportion of the population who are not infected. To study fairness, the Girvan-Newman algorithm [44] is used to partition the network G into two communities corresponding to groups c^+ and c^- .

Implementation of F-PPO. The policy network is a multiclass classifier that outputs the probabilities of $|V| + 1$ actions representing either not vaccinating or vaccinating any of the $|V|$ individuals. EO is still adopted as the short-term fairness metric, which measures the vaccination ratio among newly

infected individuals in different groups as follows

$$\Delta_s(s_t, a_t) = \left| \frac{\sum_{t-w}^t \text{vaccine_given}_{tc^+}}{\sum_{t-w}^t \text{new_infected}_{tc^+} + 1} - \frac{\sum_{t-w}^t \text{vaccine_given}_{tc^-}}{\sum_{t-w}^t \text{new_infected}_{tc^-} + 1} \right|. \quad (10)$$

For the action massaging, at each time step it checks if providing the vaccine to an individual from the other community would result in a fairer allocation. If this condition is met, the algorithm proceeds to check if there exists an individual from the other community whose predicted probability is sufficiently close to that of the current individual. If such an individual is found, the algorithm modifies the action accordingly. A dynamic threshold is again adopted. For the i th iteration, the threshold is defined as:

$$\tau(i) = \begin{cases} \min(\tau_e, \tau_s \cdot \gamma^{i-i_s}) & i \geq i_s \\ 0 & \text{otherwise} \end{cases}$$

Specifically, we set $\tau_s = 0.01$, $\tau_e = 0.35$, $i_s = 50$, $\gamma = 1.2$. Finally, long-term fairness is measured as the distance between the health states of the two communities. For other hyperparameters, we set $\lambda = 0.25$ and $\delta = 0.05$.

Agents. We consider the Max agent in addition to PPO, A-PPO, and F-PPO. The Max agent vaccinates the most susceptible individual each time, which is considered as the individual with the most number of infected neighbors.

Results. The experimental results are shown in Fig. 4. As can be seen, F-PPO achieves the best performance in terms of long-term fairness and significantly improves short-term fairness compared with the Max and PPO agents. A-PPO achieves the best performance in terms of short-term fairness, but produces the worst performance in gaining rewards. As expected, the Max agent achieves the highest utility performance, but it also demonstrates the poorest fairness performance. The combination of the results also demonstrates the capability of the F-PPO. The ablation study in Fig. 2(c) shows similar results to the other two case studies.

VII. CONCLUSIONS

In this paper, we studied the problem of achieving long-term fairness in sequential decision-making systems. We modeled

the system as a Markov Decision Process (MDP) and tackled the problem by developing a fair reinforcement learning (RL) algorithm. By acknowledging that short-term fairness and long-term fairness are distinct requirements that may not necessarily align with one another, we developed an algorithmic framework that incorporates both requirements using different bias mitigation approaches, including pre-processing and in-processing approaches. We conducted three simulation case studies. The results show that our method can strike a balance between short-term fairness, long-term fairness, and utility.

ACKNOWLEDGMENTS

This work was supported in part by NSF 2142725.

REFERENCES

- [1] K. D. Johnson, D. P. Foster, and R. A. Stine, "Impartial predictive modeling: Ensuring fairness in arbitrary models," *Statistical Science*, p. 1, 2016.
- [2] R. S. Baker and A. Hawn, "Algorithmic bias in education," *International Journal of Artificial Intelligence in Education*, pp. 1–41, 2021.
- [3] L. Zhang, Y. Wu, and X. Wu, "Achieving non-discrimination in data release," *Proceedings of ACM SIGKDD'17*, 2017.
- [4] M. S. A. Lee and L. Floridi, "Algorithmic fairness in mortgage lending: from absolute conditions to relational trade-offs," *Minds and Machines*, vol. 31, no. 1, pp. 165–191, 2021.
- [5] C. Schumann, J. Foster, N. Mattei, and J. Dickerson, "We need fairness and explainability in algorithmic hiring," in *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 2020.
- [6] R. Berk, H. Heidari, S. Jabbari, M. Kearns, and A. Roth, "Fairness in criminal justice risk assessments: The state of the art," *Sociological Methods & Research*, vol. 50, no. 1, pp. 3–44, 2021.
- [7] Y. Li, H. Chen, S. Xu, Y. Ge, J. Tan, S. Liu, and Y. Zhang, "Fairness in recommendation: A survey," *arXiv preprint arXiv:2205.13619*, 2022.
- [8] Z. Tang, J. Zhang, and K. Zhang, "What-is and how-to for fairness in machine learning: A survey, reflection, and perspective," *arXiv preprint arXiv:2206.04101*, 2022.
- [9] D. Pessach and E. Shmueli, "A review on fairness in machine learning," *ACM Computing Surveys (CSUR)*, vol. 55, no. 3, pp. 1–44, 2022.
- [10] G. Alves, F. Bernier, M. Couceiro, K. Makhoul, C. Palamidessi, and S. Zhioua, "Survey on fairness notions and related tensions," *EURO Journal on Decision Processes*, p. 100033, 2023.
- [11] L. T. Liu, S. Dean, E. Rolf, M. Simchowitz, and M. Hardt, "Delayed impact of fair machine learning," in *Proceedings of ICML'18*, ser. Proceedings of Machine Learning Research, J. G. Dy and A. Krause, Eds., vol. 80. PMLR, 2018, pp. 3156–3164.
- [12] D. Ensign, S. A. Friedler, S. Neville, C. Scheidegger, and S. Venkatasubramanian, "Runaway feedback loops in predictive policing," in *Conference on Fairness, Accountability and Transparency*. PMLR, 2018, pp. 160–171.
- [13] L. T. Liu, S. Dean, E. Rolf, M. Simchowitz, and M. Hardt, "Delayed impact of fair machine learning," in *International Conference on Machine Learning*. PMLR, 2018, pp. 3150–3158.
- [14] X. Zhang, R. Tu, Y. Liu, M. Liu, H. Kjellstrom, K. Zhang, and C. Zhang, "How do fair decisions fare in long-term qualification?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 18457–18469, 2020.
- [15] S. Jabbari, M. Joseph, M. Kearns, J. Morgenstern, and A. Roth, "Fairness in reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1617–1626.
- [16] A. D'Amour, H. Srinivasan, J. Atwood, P. Baljekar, D. Sculley, and Y. Halpern, "Fairness is not static: deeper understanding of long term fairness via simulation studies," in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020, pp. 525–534.
- [17] J. Atwood, H. Srinivasan, Y. Halpern, and D. Sculley, "Fair treatment allocations in social networks," *arXiv preprint arXiv:1911.05489*, 2019.
- [18] E. Y. Yu, Z. Qin, M. K. Lee, and S. Gao, "Policy optimization with advantage regularization for long-term fairness in decision systems," *NeurIPS*, 2022.
- [19] M. Wen, O. Bastani, and U. Topcu, "Algorithms for fairness in sequential decision making," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 1144–1152.
- [20] L. Hu and Y. Chen, "A short-term intervention for long-term fairness in the labor market," in *Proceedings of the 2018 World Wide Web Conference*, 2018, pp. 1389–1398.
- [21] F. Kamiran and T. Calders, "Classifying without discriminating," in *2009 2nd international conference on computer, control and communication*. IEEE, 2009, pp. 1–6.
- [22] A. Bower, S. N. Kitchen, L. Niss, M. J. Strauss, A. Vargas, and S. Venkatasubramanian, "Fair pipelines," *arXiv preprint arXiv:1707.00391*, 2017.
- [23] C. Dwork and C. Ilvento, "Fairness under composition," *arXiv preprint arXiv:1806.06122*, 2018.
- [24] E. Creager, D. Madras, T. Pitassi, and R. Zemel, "Causal modeling for fairness in dynamical systems," in *International Conference on Machine Learning*. PMLR, 2020, pp. 2185–2195.
- [25] Y. Hu, Y. Wu, L. Zhang, and X. Wu, "Fair multiple decision making through soft interventions," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17965–17975, 2020.
- [26] Y. Hu and L. Zhang, "Achieving long-term fairness in sequential decision making," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 2022, pp. 9549–9557.
- [27] Y. Ge, S. Liu, R. Gao, Y. Xian, Y. Li, X. Zhao, C. Pei, F. Sun, J. Ge, W. Ou *et al.*, "Towards long-term fairness in recommendation," in *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, 2021, pp. 445–453.
- [28] F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," *Knowledge and Information Systems*, vol. 33, no. 1, pp. 1–33, 2012.
- [29] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, "Certifying and removing disparate impact," in *proceedings of ACM SIGKDD'15*, 2015, pp. 259–268.
- [30] D. Xu, S. Yuan, L. Zhang, and X. Wu, "Fairgan: Fairness-aware generative adversarial networks," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 570–575.
- [31] N. Quadrianto and V. Sharmanska, "Recycling privileged learning and distribution matching for fairness," *Advances in neural information processing systems*, vol. 30, 2017.
- [32] M. P. Kim, A. Korolova, G. N. Rothblum, and G. Yona, "Preference-informed fairness," *arXiv preprint arXiv:1904.01793*, 2019.
- [33] M. B. Zafar, I. Valera, M. Rodriguez, K. Gummadi, and A. Weller, "From parity to preference-based notions of fairness in classification," *Advances in neural information processing systems*, vol. 30, 2017.
- [34] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in *NIPS*, 2016.
- [35] S. Corbett-Davies, E. Pierson, A. Feller, S. Goel, and A. Huq, "Algorithmic decision making and the cost of fairness," in *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*, 2017, pp. 797–806.
- [36] A. K. Menon and R. C. Williamson, "The cost of fairness in binary classification," in *Conference on Fairness, accountability and transparency*. PMLR, 2018, pp. 107–118.
- [37] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [38] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *arXiv preprint arXiv:1506.02438*, 2015.
- [39] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," *Advances in neural information processing systems*, vol. 12, 1999.
- [40] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [41] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.
- [42] S. Barocas, M. Hardt, and A. Narayanan, "Fairness in machine learning," *Nips tutorial*, vol. 1, p. 2017, 2017.
- [43] C. Villani, *Topics in optimal transportation*. American Mathematical Soc., 2021, vol. 58.
- [44] M. Girvan and M. E. Newman, "Community structure in social and biological networks," *Proceedings of the national academy of sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.