
Remembering to Be Fair: Non-Markovian Fairness in Sequential Decision Making

Parand A. Alamdari^{*123} Toryn Q. Klassen^{*123} Elliot Creager⁴²³ Sheila A. McIlraith¹²³

Abstract

Fair decision making has largely been studied with respect to a single decision. Here we investigate the notion of fairness in the context of sequential decision making where multiple stakeholders can be affected by the outcomes of decisions. We observe that fairness often depends on the history of the sequential decision-making process, and in this sense that it is inherently non-Markovian. We further observe that fairness often needs to be assessed at time points *within* the process, not just at the end of the process. To advance our understanding of this class of fairness problems, we explore the notion of non-Markovian fairness in the context of sequential decision making. We identify properties of non-Markovian fairness, including notions of long-term, anytime, periodic, and bounded fairness. We explore the interplay between non-Markovian fairness and memory and how memory can support construction of fair policies. Finally, we introduce the FairQCM algorithm, which can automatically augment its training data to improve sample efficiency in the synthesis of fair policies via reinforcement learning.

1. Introduction

In many real-world decision-making settings, decisions are of consequence to a diversity of entities—the *stakeholders* of the decision-making process. Furthermore, decision making is often sequential, involving multiple decisions executed over time in support of realizing near- and/or longer-term objectives. Here we study fairness in the context of sequential decision making, where multiple stakeholders

^{*}Equal contribution ¹University of Toronto, Toronto, Canada ²Vector Institute, Toronto, Canada ³Schwartz Reisman Institute for Technology and Society, Toronto, Canada ⁴University of Waterloo, Waterloo, Canada. Correspondence to: Parand A. Alamdari <parand@cs.toronto.edu>.

Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

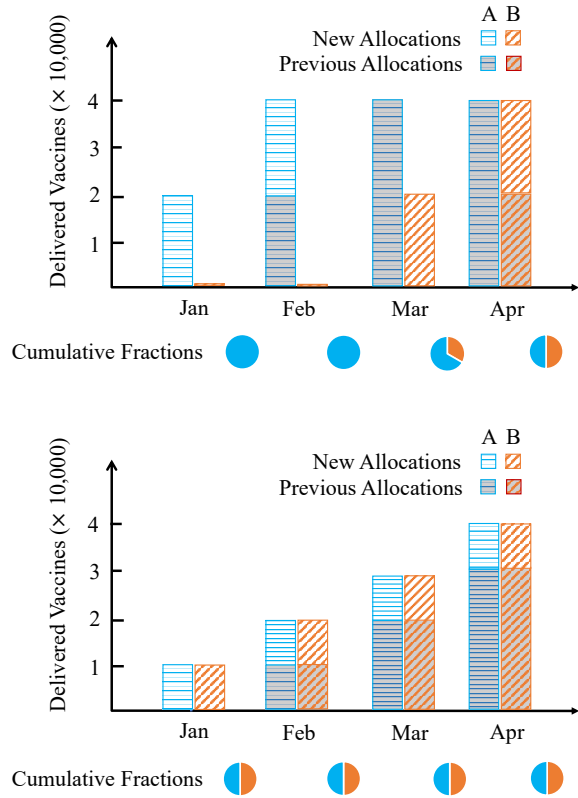


Figure 1. Two processes for distributing vaccines to countries A and B. Both result in an equal distribution of vaccine at the end. Monthly evaluation shows that the first process favors A for a time.

can be affected by the outcomes of decisions.

To ground this discussion, consider the problem of distributing vaccines to countries around the world, and our aspiration that the allocation to different countries be “fair” in some manner. Much of the decision making required to distribute vaccines is in service of getting vaccine safely from origin to destination. Delivery can be complex and costly, restricting the set of feasible plans. Indeed many such decisions will not immediately affect the fairness of the vaccine distribution. Does it even make sense to ask whether these intermediate decisions are “fair?” Perhaps what we need to aspire to is *long-term fairness* of the entire

decision-making process, in the limit when all the vaccines have been delivered. Unfortunately, this has its problems.

To illustrate, consider two greatly simplified processes to distribute 80,000 vaccines to countries A and B, as depicted in Figure 1; the first process reflects an economical plan that delivers 40,000 vaccines to country A in two consecutive monthly shipments of 20,000 each, followed by two consecutive monthly shipments of the same size to country B. In the second process, A and B each receive 10,000 vaccines each month for four months. Each process results in the delivery of 40,000 vaccines to each of A and B. However the latter process may seem “more fair,” particularly when we treat early vaccine access as a proxy for reduced negative health outcomes (e.g., Du et al., 2024).

This simple example suggests the need to measure and enforce fairness at significant timepoints throughout the process—a notion of *periodic* fairness, perhaps monthly, quarterly, or at year’s end. In the extreme we might aspire to a notion of *anytime fairness*, or perhaps a notion of *bounded fairness* where the fairness of the allocation of vaccines is judged at a time point that is determined by some property of the system, such as after delivering each million vaccines.

This example also exposes an important property of fairness in sequential decision making, that it can be inherently non-Markovian. That is, the assessment of fairness of a sequential decision-making process does not just rely on the current state, or more generally (s, a, s') , the state s' resulting from deciding to perform action a in state s . Rather, it is a function of the history of the decision-making process, the history of state-action pairs. This has implications not only for how we define various notions of fairness but also how we compute fair policies, and the role of memory.

To the best of our knowledge, this is the first work to explore these important concepts. Our contributions include:

1. We introduce the notion that multi-stakeholder fairness can be non-Markovian in sequential decision making and define formative concepts relating to the assessment of fairness at varying time intervals including long-term, periodic, anytime, and bounded fairness.
2. We study the role of memory in converting non-Markovian problems into Markovian problems so that the generation of fair policies, which are inherently non-Markovian, can be addressed using standard Markovian as well as non-Markovian methods.
3. We propose FairQCM, an algorithm that encourages sample-efficient fair policy learning by generating counterfactual experiences during the learning process. We demonstrate that FairQCM promotes sample-efficient fair learning compared with baseline methods that learn explicitly from the entire history, or use a neural memory module.

These contributions are relevant both to the synthesis of fair sequential decision making, and to the assessment of fairness with respect to historical traces of human- and/or machine-generated decision making.

There are many practical instances of sequential decision making problems where the assessment of fairness is potentially non-Markovian in nature and where some form of intermittent fairness assessment and enforcement may be necessary. There is a rich literature on vaccine allocation (Erdoğan et al., 2024). Interestingly, the COVID vaccine allocation program COVAX had the explicitly temporal objective of participating countries progressing *at the same rate*, and they developed an allocation algorithm (World Health Organization, 2021) whose inputs included the historical allocation of vaccines from the program.¹ Another practical example of sequential decision making, where the assessment and optimization of fairness can require consideration of history, is in establishing fair *waiting times* in healthcare, such as for hospital admissions or surgeries (e.g., Qi, 2017; Ala et al., 2021).

We mention these applications to reinforce the practical import and potential impact of the framing and foundational ideas presented in this paper. To simplify our discussion going forward, we will migrate from our vaccine example to an even simpler example involving the distribution of indivisible goods—in this case doughnuts.

Running Example – Doughnut Allocation: Given n people, and m doughnuts for distribution, define a policy to distribute these doughnuts in a fair manner over time. This problem presents at least two challenges: (i) defining what constitutes a fair allocation of these goods, and (ii) prescribing, or *learning* a policy to realize a fair distribution. We will use variants of this problem to illustrate our framework.

2. Related Work

We are not the first to study fairness over time. Indeed, a flurry of research activity has followed from the observation that intervening to promote fairness in the short-term can lead to distinct and sometimes unexpected results in the long run (Liu et al., 2018; Hashimoto et al., 2018; Hu & Chen, 2018; D’Amour et al., 2020).

Zhang & Shah (2014) introduced a method to maximize returns for the worst-off agent in a setting where agents have local interests. Jabbari et al. (2017) explored a more conservative notion of fairness, where a reinforcement learning (RL) agent was tasked with considering long-term discounted rewards when comparing two actions, reminiscent of a dynamic take on the individual fairness principle of

¹We note that COVAX has been subject to various criticisms (see, e.g., Usher, 2021).

“treating likes alike” (Dwork et al., 2012). Group notions of fairness have also been explored and theoretically analyzed within RL. Wen et al. (2021) offer practical approaches to encourage parity in rewards between a majority and minority group, as well as characterizations of the sample complexity involved. Deng et al. (2022) consider group fairness constraints at each step and provide theoretical guarantees for fairness violations.

Defining what constitutes fairness is an open-ended design choice with normative and political implications (Narayanan, 2018; Binns, 2018; Xinying Chen & Hooker, 2023). Given the importance of the scalar reward signal in RL, one key question is how to aggregate rewards across stakeholders and time to measure the overall fairness of a policy.

A (cardinal) social welfare function aggregates individual utilities together into one number. Suppose we have a sequence of utilities $\vec{u} = u_1, \dots, u_n$. Some well-known social welfare functions are the utilitarian welfare $Util(\vec{u}) = \sum_i u_i$, the Rawlsian welfare $Rawls(\vec{u}) = \min(u_1, \dots, u_n)$, and the Nash welfare $Nash(\vec{u}) = \prod_i u_i$ (see, e.g., Caragiannis et al., 2019). Some of these can be viewed, in part, as measuring fairness. Recent work has proposed optimizing Nash welfare, Gini social welfare, and utilitarian objectives (Mandal & Gan, 2022; Siddique et al., 2020; Fan et al., 2023). Our contributions are distinct in our focus on measuring fairness over a history of decisions—as opposed to immediate choices or long-term aggregated discounted rewards—and the special role of memory that this non-Markovian perspective on fairness entails.

Our motivating sequential doughnut allocation example is related to an established literature on allocation problems (Ibaraki & Katoh, 1988), where fairness concerns have also been explored (Kash et al., 2014; Bouveret & Lemaître, 2016; Caragiannis et al., 2019). Our focus on formal frameworks for sequential decision making differentiates us in this regard; although temporally-extended notions of fairness have also been suggested in the context of computational social choice (Boehmer & Niedermeier, 2021), they are underexplored, even in this context.

3. Preliminaries

We consider an environment that is fully observable with dynamics that are governed by a Markov Decision Process (MDP) (Puterman, 2014). We utilize an MDP variant that makes explicit the stakeholders that are affected by the sequential decision-making process and their individual reward functions. Our definition is functionally equivalent to that of a multi-objective MDP (Rojers et al., 2013).

Definition 3.1 (Multi-stakeholder Markov Decision Process). A multi-stakeholder MDP is a tuple $\langle S, s_{\text{init}},$

$A, P, R_1, \dots, R_n, \gamma \rangle$ where S is a finite set of states, $s_{\text{init}} \in S$ is the initial state, A is a finite set of actions, and $P(s_{t+1} | s_t, a_t)$ is the transition probability distribution, giving the probability of transitioning to state s_{t+1} by taking action a_t in s_t . For $i = 1, \dots, n$, $R_i : S \times A \times S \rightarrow \mathbb{R}$ is the reward function of the i th stakeholder, and $\gamma \in (0, 1]$ is the discount factor.

In such an environment, if executing action a_t in state s_t results in state s_{t+1} , each stakeholder i receives the reward $R_i(s_t, a_t, s_{t+1})$. The formalism is agnostic with respect to who—stakeholders or others—actually executes the actions.

A (Markovian) policy $\pi(a | s)$ is a probability distribution over the actions $a \in A$, given a state $s \in S$. We can also define a non-Markovian policy as a mapping from histories to distributions over actions: $\pi(a_t | s_1, a_1, \dots, a_{t-1}, s_t)$.

A trace τ of a multi-stakeholder MDP is a (possibly infinite) sequence of alternating states and actions: $s_1, a_1, s_2, a_2, s_3, a_3, \dots$ (where $s_1 = s_{\text{init}}$). A finite trace always ends with a state. Following in the spirit of Sutton & Barto (2018), given a trace τ we can define the discounted return for stakeholder i as the discounted sum of that stakeholder’s rewards accumulated over the trace: $G_i(\tau) \stackrel{\text{def}}{=} \sum_t \gamma^{t-1} R_i(s_t, a_t, s_{t+1})$ (t ranges over the time-points in the trace, whether there are finitely or infinitely many.) A policy may generate various traces (i.e., when actions are selected by it) and could be evaluated in terms of expected discounted returns.

For an MDP, which can be thought of as a multi-stakeholder MDP where $n = 1$, it is standard to try to find a policy that maximizes the expected discounted return. It is well-known that (for infinite traces) there always exists a Markovian policy that does so (Watkins & Dayan, 1992).

Which policies are to be preferred when there are more stakeholders? Each stakeholder might prefer a policy that gives itself higher expected discounted return. We consider a number of ways to assess a policy’s fairness in Section 4.

Finally, here is some further notation we use:

- For a finite trace $\tau_T = s_1, a_1, s_2, \dots, s_T, a_T, s_{T+1}$ we may use a subscript T to indicate the length (number of actions) as just shown.
- Given a trace $\tau = s_1, a_1, \dots, s_{t+1}, a_{t+1}, s_{t+2}, \dots$ of length $> t$, we can write τ_t for the prefix s_1, a_1, \dots, s_{t+1} . Recall that a finite trace ends with a state.
- For $n \in \mathbb{N}$, we define the set $[n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$.

4. Fairness over Time

In Section 1 we observed that the assessment of fairness depends on some function of the history of states and actions and consequently is non-Markovian. In this section we argue

for the importance of *non-Markovian* fairness and propose (Section 4.1) a formalism that allows for expressing many different forms of fairness over time.

Example: Returning to our simple doughnut example, consider the case where we have 3 children (stakeholders) A, B, C, and 24 doughnuts to distribute. There are three actions *giveA*, *giveB*, *giveC* that, in the deterministic case, make the propositional variables *gotA*, *gotB*, and *gotC* true, respectively. Not surprisingly received doughnuts disappear near instantaneously! This is a simple allocation problem and there are many ways of characterizing fairness, but for ease of explication, a simple strategy to assess whether, at the end of the process, the 24 doughnuts were allocated fairly is to inspect the history, count the number of times each child received a doughnut and if the sums for each child agree, we can proclaim fairness of the process. ■

To formalize this intuition, suppose that we have a function $U : (S \times A)^* \times S \rightarrow \mathbb{R}^n$ that tells us how well each stakeholder is doing at the endpoint of the given trace. The idea is that $U(\tau_T)$'s i th entry, $U(\tau_T)_i$, is some measure of how good τ_T has been to stakeholder i . We will call U the *stakeholder status function*, and the value it returns the *stakeholder status vector*.

In our doughnut example, the stakeholder status function would be a function over the entire history, resulting in a vector of size $n = 3$ denoting the number of doughnuts distributed to each of A, B, and C. For other applications, possible stakeholder status functions could be stakeholder returns or some measure of how much *envy* stakeholder i feels on τ (see Shams et al., 2021).

Suppose also for now that we have some *aggregation* function $W : \mathbb{R}^n \rightarrow \mathbb{R}$ that we can use to aggregate stakeholder scores to determine how fair things are. W could for instance be a standard social welfare function like Nash welfare (Section 2), or a boolean-valued function which outputs 1 just in case all stakeholder scores are equal. This is indeed what is required for our doughnut example.

Given U and W , when can we say that a trace τ was fair? It might seem that we have an obvious answer: to equate the fairness of the trace with fairness at the end.

Definition 4.1 (Long-term fairness). Given a status function U and an aggregation function W , the *long-term fairness* of a finite trace τ_T is $W(U(\tau_T))$.

As argued in Section 1, in many cases, we wish to assess the fairness of a process more frequently, perhaps weekly or monthly, rather than just at the end. To do so, we define a notion of periodic fairness. There is a variety of options for aggregating periodic assessments. One simple option would be to take the status vector at the end of each week, apply the aggregation function, and then sum the results. In such a case, if comparing traces of equal length, the greater the

sum, the more fair the trace. For now let us just say that we have some function $W_{\text{ex}} : (\mathbb{R}^n)^* \rightarrow \mathbb{R}$ that lets us aggregate a sequence of stakeholder status vectors. We will call W_{ex} the *extended aggregation* function (higher scores are better).

Definition 4.2 (Periodic fairness). Given a status function U and extended aggregation function W_{ex} , the *periodic fairness* (with period p) of a finite trace τ_T is $W_{\text{ex}}(U(\tau_p), U(\tau_{2p}), \dots, U(\tau_{\lfloor T/p \rfloor p}))$.

A special case is to assess the trace at every time point:

Definition 4.3 (Anytime fairness). Given a status function U and extended aggregation function W_{ex} , the *anytime fairness* of a trace τ_T is its periodic fairness with period 1.

Long-term fairness could be thought of as a special case of anytime fairness – the case in which the extended aggregation function W_{ex} ignores all but the last of its inputs.

A generalization of periodic fairness is *bounded fairness*, where we have some “filter” function $B(\tau_t)$ to indicate at which points fairness should be assessed (for example, after a dozen doughnuts has been distributed):

Definition 4.4 (Bounded fairness). Given a status function U , an extended aggregation function W_{ex} , and a binary-valued function $B(\tau_t)$, the *bounded fairness* of a finite trace τ_T is $W_{\text{ex}}(U(\tau_{t_1}), U(\tau_{t_2}), \dots, U(\tau_{t_k}))$ where (t_1, t_2, \dots, t_k) is the subsequence of $(1, 2, \dots, T)$ for which $B(\tau_{t_i}) = 1$ for each t_i .

In Section 4.1, we bring these notions together into a formalism that can express many conceptions of fairness for traces and also for policies.

4.1. Fairness Schemes

To compare traces (and later policies) in a general way, we introduce *fairness schemes*, which include as elements the status, extended aggregation, and filter functions.

Definition 4.5 (Fairness scheme). Given a multi-stakeholder MDP with state space S and action space A , a *fairness scheme* is a tuple $\mathfrak{F} = \langle U, W_{\text{ex}}, B \rangle$ where

- $U : (S \times A)^* \times S \rightarrow \mathbb{R}^n$ is the *stakeholder status function*.
- $W_{\text{ex}} : (\mathbb{R}^n)^* \rightarrow \mathbb{R}$ is the *extended aggregation function*.
- $B : (S \times A)^* \times S \rightarrow \{0, 1\}$ is the *filter function*.

In some cases we may not specify a filter function and just write $\mathfrak{F} = \langle U, W_{\text{ex}} \rangle$, which means that we assume that B is the constant function $B(\tau) = 1$ (indicating that fairness should be assessed at all time points.)

We typically pair a fairness scheme with a multi-stakeholder MDP, so for convenience we have the following definition:

Definition 4.6 (Non-Markovian Fair Decision Process (NMFDP)). A *Non-Markovian Fair Decision Process* is a tuple $\langle \mathcal{M}, \mathfrak{F} \rangle$ where \mathcal{M} is a multi-stakeholder MDP and \mathfrak{F} is a fairness scheme defined w.r.t. \mathcal{M} .

Now, we can define the fairness score of a trace.

Definition 4.7 (Fairness score of a trace). Given an NMFDP $\langle \mathcal{M}, \mathfrak{F} \rangle$ where $\mathfrak{F} = \langle U, W_{\text{ex}}, B \rangle$, the *fairness score* of a finite trace τ_T , which with a slight abuse of notation we will write as $\mathfrak{F}(\tau_T)$, is

$$\mathfrak{F}(\tau_T) \stackrel{\text{def}}{=} W_{\text{ex}}(U(\tau_{t_1}), U(\tau_{t_2}), \dots, U(\tau_{t_k}))$$

where (t_1, t_2, \dots, t_k) is the subsequence of $(1, 2, \dots, T)$ for which $B(\tau_{t_i}) = 1$ for each t_i .

That is, the fairness score is the extended aggregation of the status vectors that pass the filter (note that this definition coincides with bounded fairness).

In the doughnut example, t_i could, for example, denote each time at which another dozen doughnuts have been distributed and $U(\tau_{t_i})$ could be a vector corresponding to the number of doughnuts distributed to each of A, B, and C at that time. W_{ex} then corresponds to a function over that sequence of vectors that measures their fairness.

In the case where we are considering policies rather than historical traces, we can similarly use a fairness scheme \mathfrak{F} to score a policy by taking an expectation:

Definition 4.8 (Fairness score of a policy). Given an NMFDP $\langle \mathcal{M}, \mathfrak{F} \rangle$ and a time horizon T , the *fairness score of a policy* π according to fairness scheme $\langle U, W_{\text{ex}}, B \rangle$ is

$$\mathfrak{F}^T(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{\tau_T \sim \pi, \mathcal{M}} [\mathfrak{F}(\tau_T)]$$

where $\mathfrak{F}(\tau_T)$ is the fairness score of the trace per Definition 4.7. To evaluate fairness over infinite traces, we can consider the limit (if it exists) as $T \rightarrow \infty$; i.e., $\mathfrak{F}^\infty(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{\tau \sim \pi, \mathcal{M}} [\lim_{T \rightarrow \infty} \mathfrak{F}(\tau_T)]$.

For a given choice of fairness scheme, an **optimal policy** would be the one that maximizes the fairness score, either for some given time horizon or in the limit. Finding such a policy can be considered in different settings with or without access to the dynamics model (if we do have access, then a form of planning might be used; otherwise, we might use RL). Furthermore, we might or might not have access to the fairness scheme itself (perhaps fairness scores just come as feedback from the environment). We will consider computing policies in Section 5.

4.2. More on the extended aggregation function

For the extended aggregation function W_{ex} , one option would be to define it to first aggregate stakeholders statuses at each time point, and then aggregate across time.

Definition 4.9 (Timepoint-first). Given a fairness scheme $\langle U, W_{\text{ex}}, B \rangle$, we will say the extended aggregation function is *timepoint-first* if it can be written as

$$W_{\text{ex}}(u_1, \dots, u_k) = W_{\text{temporal}}(W(u_1), \dots, W(u_k))$$

in terms of two aggregation functions: an aggregation function $W : \mathbb{R}^n \rightarrow \mathbb{R}$ that is applied to the stakeholder status vector at each time point, and a temporal aggregation function $W_{\text{temporal}} : \mathbb{R}^* \rightarrow \mathbb{R}$ that combines the results.

Here are some example temporal aggregation functions:

- Long-term fairness: $W_{\text{temporal}}(w_1, \dots, w_k) = w_k$
- Average: $W_{\text{temporal}}(w_1, \dots, w_k) = \text{mean}(w_1, \dots, w_k)$
- The discounted sum of the aggregated values for each time point (treating them like rewards):

$$W_{\text{temporal}}(w_1, \dots, w_k) = \sum_{t=1}^k \gamma^{t-1} w_t$$

Furthermore, W_{temporal} could be a standard social welfare function, where we treat the aggregations from different time points as the individual utilities to be aggregated. For example, we could be “Rawlsian with respect to time” and only care about the worst-off timepoint: $W_{\text{temporal}}(w_1, \dots, w_k) = \min(w_1, \dots, w_k)$.

However, not all desirable extended aggregation functions are naturally expressible in a timepoint-first way. What if a trace is unfair *to* a particular stakeholder at one time, and unfair *in their favor* at another? We might want to say that those events balance out. However, from a single real number $W(U(\tau_t))$, it is awkward to extract who is being (un)fairly treated at t . That information would be easier to define in terms of $U(\tau_t)$ itself. We consider this next.

4.2.1. UNFAIRNESS FOR INDIVIDUAL STAKEHOLDERS

We can use the status function U to define when a specific stakeholder is being unfairly treated. In the doughnut example where 24 doughnuts were being distributed, a distribution of 8 doughnuts to each of A, B, and C would constitute a fair distribution, whereas a distribution of 6, 8, and 10, respectively, would correspond to a distribution that was *unfair to A in favor of C*.

Definition 4.10 (Unfair to / unfair in favor of). Given an NMFDP $\langle \mathcal{M}, \mathfrak{F} \rangle$ with status function U , a trace τ is *unfair to* stakeholder i at time t on the trace if $U(\tau_t)_i < \text{mean}(U(\tau_t))$, and is *unfair in favor of* stakeholder i at time t on the trace if $U(\tau_t)_i > \text{mean}(U(\tau_t))$.

We can also make a quantified version of unfairness:

Definition 4.11 (Unfairness to / overall unfairness to). Given an NMFDP $\langle \mathcal{M}, \mathfrak{F} \rangle$ with status function U , the *unfairness* of trace τ to stakeholder i at time t on the trace is $\text{unfairness}_i(\tau, t) \stackrel{\text{def}}{=} U(\tau_t)_i - \text{mean}(U(\tau_t))$. The *overall unfairness* of a finite trace τ_T to i is $\text{unfairness}_i(\tau_T) \stackrel{\text{def}}{=} \sum_{t \in [T]} \text{unfairness}_i(\tau_T, t)$.

An extended aggregation function could take unfairness to individuals into account in various ways; for example, we could set W_{ex} to be the negative of the sum of squared overall unfairness values: $-\sum_{i \in [n]} (\text{unfairness}_i(\tau_T))^2$.

5. Computing Fair Policies

In this section we consider approaches to computing policies for NMFDPs in certain cases—for fairness schemes that can be made to behave like reward functions.

5.1. The Role of Memory

How should actions be selected in a NMFDP? In general, an optimal policy may need to consider the whole trace τ that has occurred so far, even just to determine what the current stakeholder status vector is. In this section, we will focus on how to transform the NMFDP to make the stakeholder status function Markovian, in the following sense:

Definition 5.1 (Markovian function). Given an NMFDP $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}}, B \rangle\rangle$, we will say a function $f : (S \times A)^* \times S \rightarrow \text{codomain}(f)$ is Markovian if

$$f(s_1, a_1, \dots, s_t, a_{t+1}, s_{t+1}) = f_M(s_{t+1})$$

for some function f_M ; i.e., f only varies with the last state.

Having a Markovian status function U does not guarantee the existence of an optimal policy that is Markovian—that still depends on the rest of the fairness scheme $\mathfrak{F} = \langle U, W_{\text{ex}}, B \rangle$. (Also, in the finite horizon case the policy might have to be non-stationary.) However, there exist choices of the extended aggregation function W_{ex} and filter function B that ensure there will be optimal Markovian policies.

Suppose $B(\tau) = 1$ (which we will be assuming until Section 5.3). Suppose also that W_{ex} is timepoint-first (Definition 4.9) and so we have that $W_{\text{ex}}(U(\tau_1), \dots, U(\tau_T)) = W_{\text{temporal}}(W(U(\tau_1)), \dots, W(U(\tau_T)))$, and further that W_{temporal} treats its inputs like rewards (e.g., by returning their discounted sum). Then if U is Markovian we can apply standard MDP techniques – treating the composite function $W \circ U$ as a reward function – to find a Markovian policy. Indeed, to do so we would only need the weaker condition that $W \circ U$ is Markovian.

Therefore, it is of interest to consider how to make the status function Markovian. In some cases this can be done by augmenting the state space S to remember past information.²

Definition 5.2 (Memory augmentation/augmented). Given an NMFDP $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}} \rangle\rangle$, a *memory augmentation* is a tuple $\langle M, m_{\text{init}}, \mu \rangle$ where M is the finite set of memory states, $m_{\text{init}} \in M$ is the initial memory state, and $\mu : M \times A \times S \rightarrow M$ is the memory update function. The resulting *memory-augmented* NMFDP is an NMFDP $\langle\langle S \times M, \langle s_{\text{init}}, m_{\text{init}} \rangle, A, P', R'_1, \dots, R'_n, \gamma \rangle, \langle U', W_{\text{ex}} \rangle\rangle$ where

- $P'(\langle s_{t+1}, m_{t+1} \rangle \mid \langle s_t, m_t \rangle, a_t) = P(s_{t+1} \mid s_t, a_t)$ if $m_{t+1} = \mu(m_t, a_t, s_{t+1})$ and otherwise 0.
- $R'_i(\langle s_t, m_t \rangle, a_t, \langle s_{t+1}, m_{t+1} \rangle) = R_i(s_t, a_t, s_{t+1})$ for each i
- $U'(\langle s_1, m_1 \rangle, a_1, \dots, \langle s_t, m_t \rangle, a_t, \langle s_{t+1}, m_{t+1} \rangle) = U(s_1, a_1, \dots, s_t, a_{t+1}, s_{t+1})$

The idea is that the memory starts storing the value m_{init} , and given a trace $s_1, a_1, \dots, s_t, a_t, s_{t+1}, \dots$ the value of the memory at time $t + 1$ would be $m_{t+1} = \mu(m_t, a_t, s_{t+1})$ where m_t is the value of the memory at time t . The original environment’s dynamics are preserved, meanwhile.

While we wrote U' as not depending on the memory, the cases we are interested in are those in which it is Markovian in the memory-augmented NMFDP, and so is equivalent to a function $U'_M(\langle s_{t+1}, m_{t+1} \rangle)$. There is a whole class of status functions, which we will call *value-regular*, for which we will show that there is a memory augmentation that makes them Markovian.

Definition 5.3 (value-regular). Given an NMFDP $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}} \rangle\rangle$, the stakeholder status function U is *value-regular* if there is some finite set $V \subseteq \mathbb{R}^n$ such that for any finite trace τ we have that $U(\tau) \in V$, and for each $v \in V$ the set of traces that U maps to v is *regular* in the following sense: the set of strings $\{\sigma \in (A \times S)^* : U(s_{\text{init}}, \sigma) = v\}$ is a regular language (e.g., could be described with a regular expression) (see, e.g., Sipser, 1997).

Theorem 5.4. Let $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}} \rangle\rangle$ be an NMFDP where U is value-regular. Then there exists a memory augmentation $\langle M, m_{\text{init}}, \mu \rangle$ so that, in the resulting memory-augmented NMFDP, U' is Markovian.

For the proof, see Appendix A.1. It turns out that value-regular stakeholder status functions are the *only* ones which we can make Markovian through memory augmentation, per the following theorem (again, the proof is in Appendix A.1).

Theorem 5.5. Let $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}} \rangle\rangle$ be an NMFDP and $\langle M, m_{\text{init}}, \mu \rangle$ be a memory augmentation such that, in the resulting memory-augmented NMFDP, U' is Markovian. Then U is value-regular.

Some obvious choices of the stakeholder status function are not value-regular in general—in particular any function that can take an unbounded number of values, for instance where U gives the (possibly discounted) return for each stakeholder. In practice, however, policies may not be run long enough for memory limitations to come into play, and so we could still store the stakeholders’ returns in memory.

5.2. Counterfactual Memories for RL

In this section we consider applying RL to learn policies for (certain) memory-augmented NMFDPs, and present an RL

²To augment states with some form of memory is an idea with a long history (e.g., Bacchus et al., 1996; Peshkin et al., 1999; Thiébaux et al., 2006; Toro Icarte et al., 2020).

algorithm that takes advantage of the memory structure.

The idea of a *reward-like* fairness scheme (defined below) is that we can treat a NMFDP with one like an MDP with a single (Markovian) reward function.

Definition 5.6 (Reward-like). Given an NMFDP $\langle\langle S, s_{\text{init}}, A, P, R_1, \dots, R_n, \gamma \rangle, \langle U, W_{\text{ex}} \rangle\rangle$, we will say that the fairness scheme $\langle U, W_{\text{ex}} \rangle$ is *reward-like* if

- W_{ex} is timepoint-first, so $W_{\text{ex}}(u_1, \dots, u_T) = W_{\text{temporal}}(W(u_1), \dots, W(u_T))$,
- $W_{\text{temporal}}(w_1, \dots, w_T) = \sum_{t=1}^T \gamma^{t-1} w_t$,
- and $W \circ U$ is Markovian (in the sense of Definition 5.1) (this will always be the case if U is Markovian).

If we define the (Markovian) reward function as $R(s_t, a_t, s_{t+1}) = R(s_{t+1}) = W(U(s_{t+1}))$, the objective of a policy is to maximize the expected discounted sum of rewards (here, we will assume an infinite time horizon). For a NMFDP with a reward-like fairness scheme, we can therefore apply standard RL algorithms. When the NMFDP is memory-augmented, we can also do a bit more.

A standard reinforcement learner in an MDP collects experiences of the form $(s_t, a_t, s_{t+1}, R(s_t, a_t, s_{t+1}))$, where s_t is the state at time t and R is the reward function, with which to learn a policy. In a memory-augmented NMFDP where the fairness scheme is reward-like, the learner can do the same, and we can write an experience as

$$(\langle s_t, m_t \rangle, a_t, \langle s_{t+1}, m_{t+1} \rangle, R(\langle s_{t+1}, m_{t+1} \rangle)) \quad (1)$$

where $\langle s_t, m_t \rangle$ is the memory-augmented state at time t . However, whenever the learner has an experience as in Equation (1), we can use our knowledge of the memory augmentation $\langle M, m_{\text{init}}, \mu \rangle$ (and of R) to generate additional experiences of the form

$$(\langle s_t, m'_t \rangle, a_t, \langle s_{t+1}, m'_{t+1} \rangle, R(\langle s_{t+1}, m'_{t+1} \rangle)) \quad (2)$$

where $m'_t \in M$ is any memory state and $m'_{t+1} = \mu(m'_t, a_t, s_{t+1})$. That is, we generate the experience the learner would have had, had it been in the counterfactual memory state m'_t instead of m_t . Generating this requires only the ability to pick $m'_t \in M$ and query μ and R .

Such counterfactual experiences can then be used in training, at least for *off-policy* RL methods that do not have to learn from experience collected using the current policy. For instance, tabular Q-learning (Watkins & Dayan, 1992) could be modified to, after each real experience in the environment, generate (some subset of) the possible counterfactual experiences as in Equation (2), and use all of the experiences to update the Q-function. We call the resulting algorithm **Fair Q-Learning with Counterfactual Memories (FairQCM)** (see Algorithm 1). For methods like DQN (Mnih et al., 2015) that use a replay buffer, counterfactual experiences could be added to the buffer. Action selection, similar to classic Q-learning, can work with a variety of choices. A

commonly adopted strategy is epsilon-greedy.

Algorithm 1 Tabular FairQCM

Input: A memory-augmented NMFDP $\langle\langle S \times M, \langle s_{\text{init}}, m_{\text{init}} \rangle, A, P', R'_1, \dots, R'_n, \gamma \rangle, \langle U', W_{\text{ex}} \rangle\rangle$ with reward-like fairness scheme and $\gamma < 1$, and learning rates $(\alpha_1, \alpha_2, \dots)$.
 Let $R = W \circ U'$
 Initialize the Q-function
 Let $\langle s_1, m_1 \rangle = \langle s_{\text{init}}, m_{\text{init}} \rangle$
for each time $t = 1, 2, 3, \dots$ **do**
 Select some subset $M'_t \subseteq M \setminus \{m_t\}$.
 Select some action a_t , and execute a_t in the current state $\langle s_t, m_t \rangle$ to acquire the experience $(\langle s_t, m_t \rangle, a_t, \langle s_{t+1}, m_{t+1} \rangle, R(\langle s_{t+1}, m_{t+1} \rangle))$
 for each $m'_t \in M'_t$ **do**
 Construct the counterfactual experience $(\langle s_t, m'_t \rangle, a_t, \langle s_{t+1}, m'_{t+1} \rangle, R(\langle s_{t+1}, m'_{t+1} \rangle))$
 where $m'_{t+1} = \mu(m'_t, a_t, s_{t+1})$.
 end for
 for each (real or counterfactual) experience $\langle x_t, a_t, x_{t+1}, r_{t+1} \rangle$ **do**
 /* a Q-learning update with learning rate α_t */
 $Q(x_t, a_t) \leftarrow^{\alpha_t} r_{t+1} + \gamma \max_{a \in A} Q(x_{t+1}, a)$
 end for
end for

This approach of generating counterfactual experiences is adapted from the CRM algorithm from Toro Icarte et al. (2022). There, instead of memory states, they were dealing with states of a “reward machine”, an automaton used to define a non-Markovian reward function.

For tabular FairQCM, we prove that (under some conditions) given a memory-augmented NMFDP with reward-like fairness scheme, it converges to the optimal Q-function, yielding the optimal policy³ (proof in Appendix A.1).

Theorem 5.7 (Convergence of tabular FairQCM). *Let $n^i(x, a)$ be the time at which action a is executed in state x for the i th time (including “executions” in counterfactual experiences) in a run of FairQCM. If the selection of actions and counterfactual experiences and of the learning rates $(\alpha_1, \alpha_2, \dots)$ are such as to satisfy the conditions from Watkins & Dayan (1992, p. 282) that*

$$0 \leq \alpha_n < 1, \quad \sum_{i=1}^{\infty} \alpha_{n^i(x,a)} = \infty, \quad \sum_{i=1}^{\infty} [\alpha_{n^i(x,a)}]^2 < \infty$$

for all x and a , then the Q-function computed by FairQCM will converge (as $t \rightarrow \infty$) to the optimal with probability 1.

³The optimal Q-function is a function $Q^*(x, a)$ which gives the expected discounted return in state x if action a is taken and afterwards the best possible action is taken in all future time steps. For an optimal policy π^* , $\pi^*(a | x) = 0$ if $a \notin \arg \max_{a'} Q^*(x, a')$.

5.3. Computations with Filter Functions

Recall that, in general, a fairness scheme (Definition 4.5) includes a filter function B that selects which points should be assessed for fairness. We can generalize the notion of a *reward-like* fairness scheme (Definition 5.6) to allow for non-trivial filter functions.

Definition 5.8 (Bounded reward-like). Given an NMFDP $\langle \mathcal{M}, \langle U, W_{\text{ex}}, B \rangle \rangle$, we will call the fairness scheme $\langle U, W_{\text{ex}}, B \rangle$ *bounded reward-like* if

- W_{ex} is timepoint-first, so $W_{\text{ex}}(u_1, \dots, u_k) = W_{\text{temporal}}(W(u_1), \dots, W(u_k))$,
- $W_{\text{temporal}}(w_1, \dots, w_k) = \sum_{t=1}^k w_t$,
- $W \circ U$ is Markovian (this will always be the case if U is Markovian),
- and B is also Markovian.

Importantly, B could be Markovian as a result of memory augmentation, in much the same way as we considered making U Markovian in Section 5.1. For example, consider periodic fairness with period p , which corresponds to having $B(s_1, a_1, \dots, a_k, s_k) = 1$ iff $k \equiv 0 \pmod{p}$. That could be made Markovian by having a mod p counter as memory.

In contrast to Definition 5.6, we assume a finite horizon, and set $\gamma = 1$. Therefore, there is no discounting in the sum computed by W_{temporal} . For a bounded reward-like fairness scheme, it can be seen that the fairness score of a finite trace is the same as its undiscounted return using the following reward function R , allowing for standard algorithms for finite horizon MDPs:

$$R(s_t, a_t, s_{t+1}) = R(s_{t+1}) = W(U(s_{t+1})) \cdot B(s_{t+1})$$

6. Experiments

In this section we compare different methods for designing the augmented memory via two simulation studies. Experiments show how generating counterfactual memories during RL can improve the overall fairness and sample efficiency of training in dynamic settings with multiple stakeholders. In each experiment, let $\langle \mathcal{M}, \mathfrak{F} \rangle$ be an NMFDP where \mathcal{M} is a multi-stakeholder MDP and $\mathfrak{F} = \langle U, W_{\text{ex}}, B \rangle$ is a fairness scheme defined w.r.t. \mathcal{M} , where W_{ex} is timepoint-first, and $B(\tau) = 1$. Given the NMFDP, we implement FairQCM and several baseline memory-augmented agents within a Deep RL framework. Analogous tabular results are in Appendix A.2, a continuous experiment is in Appendix A.3, and technical details of experiments are in Appendix A.4.

Baselines. To evaluate how memory affects the ability to learn a fair policy, we evaluate several baselines based on Deep Q Networks (DQN) (Mnih et al., 2015) with different types of augmented memory units M : *Full* stores the entire stakeholder status $U(\tau_t)$ while *Min* stores $U(\tau_t) - \min_i U(\tau_t)_i$; and *RNN* does not have a separate

memory, instead it has an extra layer of GRU (Cho et al., 2014) to remember the past. *FairQCM* is a deep, DQN-based, version of our proposed method from Section 5.2, which we use with different types of memory.

6.1. Resource Allocation

We modeled our running example of doughnut allocation as a stochastic environment, featuring a doughnut shop that bakes one doughnut in each step. There are $n = 5$ customers, with customer i being in front of the counter with probability p_i in each step. The decision to allocate the freshly baked doughnut is made by the server, who selects one customer for the allocation. If the chosen customer is not at the counter, the doughnut goes to waste. A state encapsulates the presence of individuals at the counter, and there are n associated actions, each corresponding to the allocation of a doughnut to a specific customer. At step t , $U(\tau_t)_i$ is the number of doughnuts allocated to customer i so far. We use Nash welfare to define the aggregation function for each time point as $W(U(\tau_t)) = \log(\text{Nash}(U(\tau_t) + 1)) = \sum_{i \in [n]} \log(U(\tau_t)_i + 1)$. The server’s goal is to maximize the discounted sum of $W \circ U$, treating it as the reward (except that on steps where the server wastes a doughnut, the reward is considered to be 0). The episode length is 100.

In this setting FairQCM stores $U(\tau_t)$ in the memory (like the *Full* baseline) and generates counterfactual experiences based on it. For this environment we have some extra memory-augmented baselines. *Reset* stores $U(\tau_t)$ but if in a time t all $U(\tau_t)_i$ are equal, it sets them to zero and counts from there for the next steps. *Oracle* and *Random* are hard-coded solutions, with *Oracle* realizing the optimal turn taking algorithm, and *Random* taking random actions.

Figure 2 illustrates the accumulated Nash welfare scores at the end of the episode in different phases of training across five approaches with different types of augmented memory. FairQCM outperforms other algorithms. The RNN baseline learns an approximation of the history and achieves higher accumulated Nash welfare scores compared to some of the other baselines. The RNN also demonstrates faster adaptation to avoid wasting doughnuts, as it focuses on the recent past rather than the entire history.

6.2. Simulated Lending

Consumer lending is an established test bed for fair machine learning algorithms (Dwork et al., 2012; Hardt et al., 2016), owing in part to US regulations on credit scoring and banking practices (Barocas & Selbst, 2016). As we are interested in fairness over time, we adapt the dynamic lending environment of Liu et al. (2018). We consider a finite pool of loan applicants. Each applicant is characterized by their credit score C , representing the probability of loan repayment, and belongs to a protected group within the

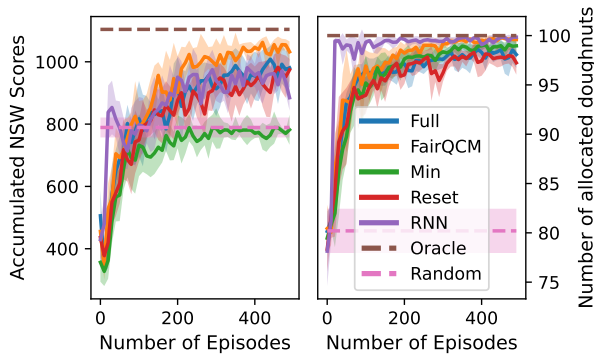


Figure 2. Resource Allocation: In simulations of our doughnut allocation task, (deep) FairQCM achieves higher Nash welfare than competing memory-augmented RL agents (left), while learning to allocate doughnuts effectively near the end of training (right).

population (two groups in total: A and B). The credit score distribution differs between the two groups. In each step, a subset of applicants applies for a loan, and the bank must decide which applicant to grant the loan. Each applicant i applies for a loan in each step with a probability of p_i . Successful loan repayment increases the bank’s utility by r and the applicant’s credit score by c , while defaulting decreases the bank’s utility by r and the applicant’s credit score by c . In our experiments, we set $r = 1$ and $c = 0.1$. The states encode the subset of applicants applying for a loan, their credit score, and profit margin so far. To assess the fairness of the loan granting process, we calculate the difference in the number of loans allocated to each protected subgroup. At step t , $U(\tau_t)_i$ is the number of loans allocated to person i so far. The aggregation function for each time point, which is inspired by demographic parity (Dwork et al., 2012), we call the *Relaxed Demographic Parity (DP) Score* and is set to $W(U(\tau_t)) = -|\sum_{i \in A} U(\tau_t)_i - \sum_{i \in B} U(\tau_t)_i|$. To ensure fairness, the bank’s goal is to maximize the discounted sum of $W \circ U$, treating that as the reward function – with a couple exceptions. The bank aims to achieve a profit margin of at least 10 percent. If at the end of the episode it does not make the targeted profit, it incurs a substantial negative reward; furthermore, there is a negative reward for granting a loan to someone who didn’t apply. In the simulated lending environment, for the *Full* baseline in this experiment, instead of storing each stakeholder status, we store $U(\tau_t)_X = \sum_{i \in X} U(\tau_t)_i$ in the memory for each protected subgroup $X \in \{A, B\}$. Similarly, the *Min* baseline’s memory stores $\langle U(\tau_t)_A, U(\tau_t)_B \rangle - \min_{i \in \{A, B\}} U(\tau_t)_i$. We use FairQCM with the same memory as the *Min* baseline.

Figure 3 illustrates the accumulated relaxed DP scores at the end of the episode in different phases of training across four approaches with different types of augmented memory. Similar to resource allocation, FairQCM outperforms the other approaches. The RNN baseline still learns an approximation of history through 40 steps of an episode and

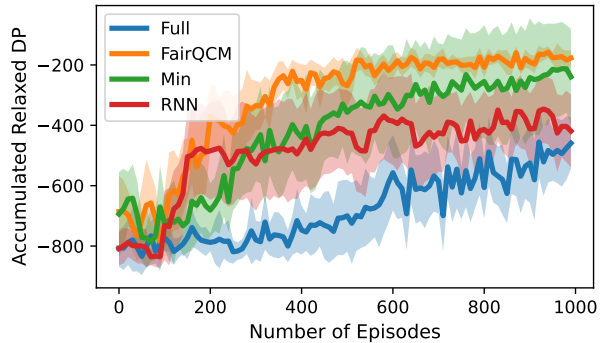


Figure 3. Simulated Lending: Accumulated Relaxed Demographic Parity scores for different approaches of augmenting memory during different phases of training.

achieves higher accumulated relaxed DP scores compared to some of the other baselines. All approaches learn policies that achieve the desired profit margin of 10 percent.

The outcomes from experiments in distinct environments emphasize the complexity of selecting and engineering an appropriate memory solution. Algorithmically, FairQCM emerges as a standout performer with superior results, while RNN-based methods provide viable solutions, particularly where engineering a memory proves to be challenging.

7. Conclusion

We have explored the notion of multi-stakeholder fairness in the context of sequential decision making. We have argued that the fairness of sequential decision making is inherently non-Markovian since the assessment of fair traces relies on the history of states and actions. We have also argued that fairness of processes is naturally assessed at significant time points, defining long-term, periodic, anytime, and bounded fairness. We have observed that in a number of circumstances, memory can be used to convert a non-Markovian fairness problem into a Markovian one, making it amenable to standard Markovian solutions—useful in the context of policy generation. We have studied the performance of various methods (Markovian and non-Markovian) and various memory models. We also proposed a method called FairQCM to generate counterfactual experience that expedite learning a fair policy, proved convergence properties in the tabular case, and demonstrated its effectiveness.

We hope the contributions of this paper will provide the rich foundations for future work on fairness in sequential decision making—both in the area of policy generation and in auditing. A limitation of this work is that we have only scratched the surface in terms of exploring different fairness scoring functions and how to compute policies for them, and much work is left to do to formulate and assess our regime in practical settings in collaboration with domain experts.

Impact Statement

Automated sequential decision making systems can affect many stakeholders, and thus have the potential to induce a variety of complex societal impacts. In this paper we focus on how the fairness and overall social welfare of these systems evolves over time. While our approach emphasizes the importance of historical decisions on near- and long-term fairness, we ultimately propose a flexible framework that is compatible with many existing fairness criteria. For the proposed research to be socially beneficial, our methodological approach underscores the importance of a trained domain expert who can act in good faith to design appropriate fairness criteria that reflect normative commitments appropriate to the problem at hand.

Acknowledgements

We wish to acknowledge funding from the Natural Sciences and Engineering Research Council of Canada (NSERC), the Canada CIFAR AI Chairs Program (Vector Institute), and Microsoft Research. The second author also received funding from Open Philanthropy. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute for Artificial Intelligence.

References

- Ala, A., Alsaadi, F. E., Ahmadi, M., and Mirjalili, S. Optimization of an appointment scheduling problem for healthcare systems based on the quality of fairness service using whale optimization algorithm and NSGA-II. *Scientific Reports*, 11:19816, 2021.
- Bacchus, F., Boutilier, C., and Grove, A. J. Rewarding behaviors. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pp. 1160–1167, 1996.
- Barocas, S. and Selbst, A. D. Big data’s disparate impact. *California Law Review*, 104(3):671–732, 2016.
- Binns, R. What can political philosophy teach us about algorithmic fairness? *IEEE Security and Privacy*, 16(3): 73–80, 2018.
- Boehmer, N. and Niedermeier, R. Broadening the research agenda for computational social choice: Multiple preference profiles and multiple solutions. In *AAMAS ’21: 20th International Conference on Autonomous Agents and Multiagent Systems*, pp. 1–5, 2021.
- Bouveret, S. and Lemaître, M. Characterizing conflicts in fair division of indivisible goods using a scale of criteria. *Autonomous Agents and Multi-Agent Systems*, 30(2):259–290, 2016.
- Caragiannis, I., Kurokawa, D., Moulin, H., Procaccia, A. D., Shah, N., and Wang, J. The unreasonable fairness of maximum Nash welfare. *ACM Transactions on Economics and Computation (TEAC)*, 7(3):12:1–12:32, 2019.
- Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- D’Amour, A., Srinivasan, H., Atwood, J., Baljekar, P., Sculley, D., and Halpern, Y. Fairness is not static: deeper understanding of long term fairness via simulation studies. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 525–534, 2020.
- Deng, Z., Sun, H., Wu, Z. S., Zhang, L., and Parkes, D. C. Reinforcement learning with stepwise fairness constraints. *arXiv preprint arXiv:2211.03994*, 2022.
- Du, H., Saiyed, S., and Gardner, L. M. Association between vaccination rates and COVID-19 health outcomes in the United States: a population-level statistical analysis. *BMC Public Health*, 24:220, 2024.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, ITCS ’12*, pp. 214–226, New York, NY, USA, 2012. Association for Computing Machinery.
- Erdoğan, G., Yücel, E., Kiavash, P., and Salman, F. S. Fair and effective vaccine allocation during a pandemic. *Socio-Economic Planning Sciences*, 93:101895, 2024.
- Fan, Z., Peng, N., Tian, M., and Fain, B. Welfare and fairness in multi-objective reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2023*, pp. 1991–1999, 2023.
- Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- Hashimoto, T., Srivastava, M., Namkoong, H., and Liang, P. Fairness without demographics in repeated loss minimization. In *International Conference on Machine Learning*, pp. 1929–1938. PMLR, 2018.
- Hu, L. and Chen, Y. A short-term intervention for long-term fairness in the labor market. In *Proceedings of the 2018 World Wide Web Conference*, pp. 1389–1398, 2018.
- Ibaraki, T. and Katoh, N. *Resource Allocation Problems: Algorithmic Approaches*. MIT press, 1988.

- Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., and Roth, A. Fairness in reinforcement learning. In *International Conference on Machine Learning*, pp. 1617–1626. PMLR, 2017.
- Kash, I., Procaccia, A. D., and Shah, N. No agent left behind: Dynamic fair division of multiple resources. *Journal of Artificial Intelligence Research*, 51:579–603, 2014.
- Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. Delayed impact of fair machine learning. In *International Conference on Machine Learning*, pp. 3150–3158. PMLR, 2018.
- Mandal, D. and Gan, J. Socially fair reinforcement learning. *arXiv preprint arXiv:2208.12584*, 2022.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Narayanan, A. Translation tutorial: 21 fairness definitions and their politics. In *Proc. Conference on Fairness, Accountability and Transparency*, 2018.
- Peshkin, L., Meuleau, N., and Kaelbling, L. P. Learning policies with external memory. In *Proceedings of the Sixteenth International Conference on Machine Learning (ICML 1999)*, pp. 307–314. Morgan Kaufmann, 1999.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- Qi, J. Mitigating delays and unfairness in appointment systems. *Management Science*, 63(2):566–583, 2017.
- Roijers, D. M., Vamplew, P., Whiteson, S., and Dazeley, R. A survey of multi-objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113, 2013.
- Shams, P., Beynier, A., Bouveret, S., and Maudet, N. Minimizing and balancing envy among agents using ordered weighted average. In *Algorithmic Decision Theory - 7th International Conference, ADT 2021*, volume 13023 of *Lecture Notes in Computer Science*, pp. 289–303. Springer, 2021.
- Siddique, U., Weng, P., and Zimmer, M. Learning fair policies in multi-objective (deep) reinforcement learning with average and discounted rewards. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 8905–8915. PMLR, 2020.
- Sipser, M. *Introduction to the Theory of Computation*. PWS Publishing Company, 1997.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Thiébaux, S., Gretton, C., Slaney, J. K., Price, D., and Kabanza, F. Decision-theoretic planning with non-Markovian rewards. *Journal of Artificial Intelligence Research*, 25:17–74, 2006.
- Toro Icarte, R., Valenzano, R., Klassen, T. Q., Christoffersen, P., Farahmand, A.-m., and McIlraith, S. A. The act of remembering: A study in partially observable reinforcement learning. *arXiv preprint arXiv:2010.01753*, 2020.
- Toro Icarte, R., Klassen, T. Q., Valenzano, R., and McIlraith, S. A. Reward machines: Exploiting reward function structure in reinforcement learning. *Journal of Artificial Intelligence Research*, 73:173–208, 2022.
- Usher, A. D. A beautiful idea: how COVAX has fallen short. *The Lancet*, 397(10292):2322–2325, 2021.
- Watkins, C. J. C. H. and Dayan, P. Q-learning. *Machine Learning*, 8:279–292, 1992.
- Wen, M., Bastani, O., and Topcu, U. Algorithms for fairness in sequential decision making. In *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS*, pp. 1144–1152. PMLR, 2021.
- World Health Organization. Allocation logic and algorithm to support allocation of vaccines secured through the COVAX facility, 2021. URL <https://www.who.int/publications/m/item/allocation-logic-and-algorithm-to-support-allocation-of-vaccines-secured-through-the-covax-facility>.
- Xinying Chen, V. and Hooker, J. A guide to formulating fairness in an optimization model. *Annals of Operations Research*, pp. 1–39, 2023.
- Zhang, C. and Shah, J. A. Fairness in multi-agent sequential decision-making. In *Advances in Neural Information Processing Systems 27*, 2014.

A. Appendix

A.1. Proofs

Proof of Theorem 5.4. Let us say that the set of possible output values of U is $V = \{v_1, \dots, v_k\}$. Since U is value-regular, for each value $v_i \in V$ there exists a deterministic finite automaton (DFA) (see, e.g., Sipser, 1997, Chapter 1)

$$A_U^i = \langle Q^i, \Sigma, \delta^i, q_0^i, Acc^i \rangle$$

that accepts the string $\sigma = a_1, \dots, a_{t+1}$ iff $U(s_{\text{init}}, \sigma) = v_i$. Note that Q^i is the finite set of automaton states, $\Sigma = A \times S$ is the set of input symbols, $\delta^i : Q^i \times \Sigma \rightarrow Q^i$ is the transition function, q_0^i is the automaton's initial state, and $Acc^i \subseteq Q^i$ is the automaton's set of accepting states.

We define a memory augmentation $\langle M, m_{\text{init}}, \mu \rangle$ as follows:

$$\begin{aligned} M &= Q^1 \times Q^2 \times \dots \times Q^k \\ m_{\text{init}} &= \langle q_0^1, q_0^2, \dots, q_0^k \rangle \\ \mu(\langle q^1, \dots, q^k \rangle, a, s) &= \langle \delta^1(q^1, \langle a, s \rangle), \dots, \delta^k(q^k, \langle a, s \rangle) \rangle \end{aligned}$$

That is, the memory keeps track of the states of the automata that correspond to each possible output value v_1, \dots, v_k . Then it can be seen that U' is Markovian in the memory-augmented NMFDP, because we can define

$$U'_M(\langle s, \langle q^1, \dots, q^k \rangle \rangle) = v_i \text{ iff } q^i \in Acc^i$$

This is well-defined since the languages of the different automata are necessarily a partition of $(A \times S)^*$, so exactly one automaton will always be in an accepting state. \square

Proof of Theorem 5.5. U can only take the same values as U' , and since U' is Markovian, it can only take finitely many values (one for each augmented state $\langle s, m \rangle$). Let us say that $V = \{v_1, \dots, v_k\}$ is the set of all possible output values of U . For each v_i , we can construct a DFA $A_U^i = \langle Q^i, \Sigma, \delta^i, q_0^i, Acc^i \rangle$ where

- $Q^i = S \times M$ is the set of automaton states,
- $q_0^i = \langle s_{\text{init}}, m_{\text{init}} \rangle$ is the initial state,
- $\Sigma = A \times S$ is the alphabet,
- $\delta^i(\langle s, m \rangle, \langle a, s' \rangle) = \langle s', \mu(m, a, s') \rangle$ is the transition function, and
- $Acc^i = \{ \langle s, m \rangle \in S \times M : U'_M(\langle s, m \rangle) = v_i \}$ is the set of accepting states.

It can be shown by induction that the DFA A_U^i accepts the sequence $\sigma = \langle a_1, s_2 \rangle, \langle a_2, s_3 \rangle, \dots, \langle a_t, s_{t+1} \rangle$ if and only if $U(s_{\text{init}}, \sigma) = v_i$. It follows from the equivalence of DFAs and regular expressions that U is value-regular. \square

Proof of Theorem 5.7. The proof of convergence of Q-learning by Watkins & Dayan (1992) mostly carries over; the only thing we have to be careful about is whether the generated counterfactual experiences are biased relative to the environment's transition probabilities.⁴ However, since the set M'_t of counterfactual memory states to use in constructing the counterfactual experiences is chosen *before* observing the outcome of action a_t , that issue does not arise: all action outcomes used in training are sampled according to the transition probabilities. \square

A.2. Tabular Q-Learning Experiments

We carry out the experiments described in Section 6.1 in a tabular setting. The environment setting is similar to Section 6.1, where there are $n = 3$ people in the doughnut shop in the tabular version, and the episode length is 12. We evaluate several baselines based on Q-Learning with different types of augmented memory units M :

- *Full* stores the entire stakeholder status $U(\tau_t)$
- *FairQCM* stores the entire stakeholder status $U(\tau_t)$ as the memory m_t at time t and uses Algorithm 1 to generate counterfactual experiences. We limit it to 8 counterfactual experiences per time step, each such experience corresponding to a counterfactual memory m'_t where for each i , $(m_t)_i < (m'_t)_i \leq (m_t)_i + 2$.
- *Min* stores $U(\tau_t) - \min_i U(\tau_t)_i$
- *Reset* stores $U(\tau_t)$ but if in a time t all $U(\tau_t)_i$ are equal, it sets them to zero and counts from there for the next steps
- *Oracle* and *Random* are hard-coded solutions, with *Oracle* realizing the optimal turn taking algorithm, and *Random* taking random actions at each step.

⁴For example, if counterfactual experiences were only generated at time t if the action a_t resulted in a lottery being won, that could make winning the lottery seem more likely than it really is.

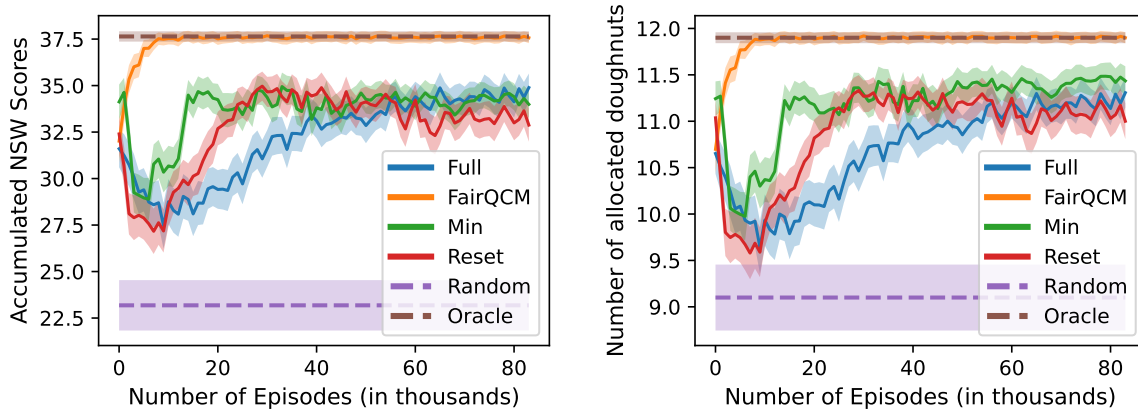


Figure 4. **Tabular Q-Learning for Resource Allocation:** The left plot shows accumulated Nash welfare scores at the end of the episode for different approaches of augmented memory in different phases of training. The right plot shows the number of doughnuts that are not wasted during the process for each approach.

Similar to the results in Section 6, FairQCM outperforms other approaches in terms of both NSW scores, and number of samples needed to train.

Technical details We set $\gamma = 0.99$, $\alpha = 0.1$, and use epsilon-greedy for exploration. $\epsilon = 1.0$ at the beginning for each state, and every time we visit a state s , ϵ_s is multiplied by 0.95 (diminishing factor), and remains greater than 0.2. We ran each method 10 times, and Figure 4 shows the average and variance of those runs across 100000 time steps.

A.3. Extra Experiments with Simulated Lending Environment

We extend the experiments described in Section 6.2 in the simulated lending domain, to a scenario where the credit scores of applicants are continuous variables. The change in the credit scores (whether positive or negative) is sampled from a Gaussian distribution ($\mu = 0.05$, $\sigma = 0.1$ when the applicant successfully repays their loan, and $\mu = -0.05$, $\sigma = 0.1$ when the applicant defaults). The rest of the environment and the baselines are the same as Section 6.2. Figure 5 illustrates the accumulated relaxed DP scores at the end of the episode in different phases of training across three approaches with different types of augmented memory. We observe the analogous behavior that FairQCM was more sample efficient, yielding superior fairness measures for a given number of episodes.

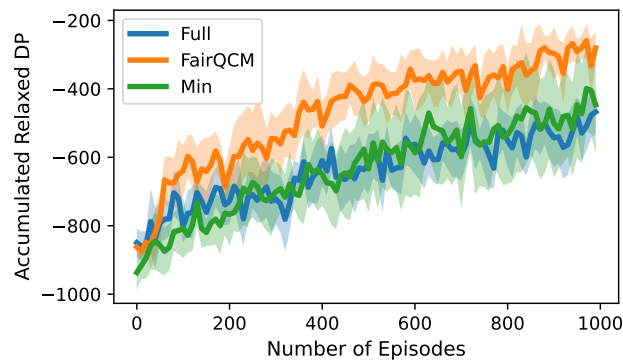


Figure 5. **Simulated Lending with Gaussian Credit Score Changes:** Accumulated Relaxed Demographic Parity scores for different approaches of augmenting memory during different phases of training.

A.4. Experimental Details

We ran the experiments on a system with the following specification: 2.3 GHz Quad-Core Intel Core i7 and 32 GB of RAM. The total running time is less than 24 hours. The code for all experiments is available at <https://github.com/praal/remembering-to-be-fair>.

A.4.1. RESOURCE ALLOCATION

State Representation We consider $n = 5$ people at the doughnut shop. Each person is at the counter with $p = 0.8$. The state is a binary sequence of length n showing the people present at the counter.

Actions There are n actions which represent allocating a doughnut to person i for each i .

Rewards $U(\tau_t)_i$ represents the number of doughnuts person i got so far. At time step t , if the doughnut is not wasted $r_t = W(U(\tau_t)) = \log(\text{Nash}(U(\tau_t) + 1)) = \sum_{i \in [n]} \log(U(\tau_t)_i + 1)$ otherwise, $r_t = 0$.

Neural Network Architectures Our DQN consists of 4 fully connected layers with ReLU activation function: states $\times 32$, 32×16 , 16×8 , $8 \times$ actions

RNN approach consists of 3 fully connected layers and one GRU layer: states $\times 32$, 32×16 , GRU (hidden states = 256), $16 \times$ actions

FairQCM FairQCM stores the entire stakeholder status $U(\tau_t)$ as the memory m_t at time t and uses Algorithm 1 to generate counterfactual experiences. We limit it to 32 counterfactual experiences per time step, each such experience corresponding to a counterfactual memory m'_t where for each i , $(m_t)_i < (m'_t)_i \leq (m_t)_i + 2$. FairQCM stores the counterfactual experiences in the replay buffer, and later samples from them during training.

Results We ran each method 10 times, and plot the average and variance of results across 1000 episodes in Figure 2.

Hyperparameters See Table 1.

A.4.2. SIMULATED LENDING

State Representation We consider $n = 4$ people in total, where each group has two people in it. The initial credit score of people in group A is 0.5, and 0.9 for group B . Each applicant applies for a loan at each step with $p = 0.9$. The state consists of a binary sequence of length n showing the applicants applying for a loan, credit score of each applicant, and profit margin so far. If an applicant gets a loan and repays it, their credit score increases by 0.1, and similarly it decreases by 0.1 if the applicant defaults. However, a credit score of an applicant always remains in the range of $[0.2, 0.9]$.

Actions There are n actions which represent granting the loan to person i for each i .

Rewards The rewards in this encoding of the problem are serving two purposes. One purpose is to incentivize the system to learn to do the right thing—in this case, to only give loans to people who applied, and to maintain the stipulated profit margin. In service of this, (i) at the final time step, if the bank didn't make the profit margin of 10 percent then $r_t =$

Hyperparameter	Full, Min, and Reset approaches	RNN approach	FairQCM
Episode Length	100	100	100
Learning Rate	0.0001	0.002	0.0001
Discount Factor (γ)	0.95	0.95	0.95
Min Exploration Rate (ϵ)	0.2	0.2	0.2
Replay Buffer Size	400	1000	6400
Batch Size	64	256	2048

Table 1. Resource Allocation Hyperparameters

$-10 \times$ episode length; and (ii) if the bank grants a loan to someone who didn't apply for a loan then $r_t = -$ episode length, where episode length is 40.

The reward is also used to reflect the fairness score. $U(\tau_t)_i$ represents the number of loans person i has received so far. At time step t , if the bank grants a loan to an applicant that applied for the loan then $r_t = W(U(\tau_t)) = -|\sum_{i \in A} U(\tau_t)_i - \sum_{i \in B} U(\tau_t)_i|$.

Neural Network Architectures Our DQN consists of 3 fully connected layers with ReLU activation function: states $\times 32$, 32×8 , $8 \times$ actions.

RNN approach consists of 3 fully connected layers and one GRU layer: states $\times 32$, 32×16 , GRU (hidden states = 256), $16 \times$ actions.

FairQCM For each protected subgroup $X \in \{A, B\}$, FairQCM stores $\sum_{j \in X} U(\tau_t)_j - \min_{k \in \{A, B\}} \sum_{j \in k} U(\tau_t)_j$ as the memory m_t at time t and uses Algorithm 1 to generate counterfactual experiences. As each m_t can be written as $m_t = \langle 0, x \rangle$ or $m_t = \langle x, 0 \rangle$ ($x \in [n]$), we limit the number of counterfactual experiences to at most 10 per time step, each such experience corresponding to a counterfactual memory m'_t where, if $m_t = \langle 0, x \rangle$, $m'_t \in \{\langle 0, x - 5 \rangle, \langle 0, x - 4 \rangle, \dots, \langle 0, x - 1 \rangle, \langle 0, x + 1 \rangle, \dots, \langle 0, x + 5 \rangle\}$, and similarly if $m_t = \langle x, 0 \rangle$. FairQCM stores the counterfactual experiences in the replay buffer, and later samples from them during training.

Results We ran each method 10 times, and plotted the average and variance of results across 1000 episodes in Figure 3.

Hyperparameters See Table 2.

Hyperparameter	Full and Min approaches	RNN approach	FairQCM
Episode Length	40	40	40
Learning Rate	0.0001	0.005	0.0001
Discount Factor (γ)	0.95	0.95	0.95
Min Exploration Rate (ϵ)	0.2	0.2	0.2
Replay Buffer Size	1000	2000	8000
Batch Size	64	512	512

Table 2. Simulated Lending Hyperparameters