

Capturing Entitlement Justice in Algorithmic Decision-Making

Edward Speer
California Institute of Technology
espeer@caltech.edu

Llama-2-7b
Meta Platforms, Inc.

April, 2025

1 Introduction

The rise of algorithmic decision making in the public sector has caused significant public concern. Algorithms increasingly make decisions that affect individuals' lives, from determining creditworthiness to predicting criminal recidivism, and the public has grown cautious of their potential to perpetuate and exacerbate existing social inequalities. A 2018 study showed that 58% of Americans believe that algorithms will always have some level of bias ([Smith, 2018](#)), and as documented in the famed COMPAS case, these fears are not unfounded ([Angwin et al., 2016](#)).

In response to these concerns, a growing body of research has focused on developing algorithmic fairness measures to evaluate and mitigate the biases in algorithmic decision making. A large number of different measures have been proposed ([Corbett-Davies et al., 2023](#)) and applied to a wide range of problems. However, many questions arise about the foundations of these measures and how to apply them sociotechnical systems. ([Hardt et al., 2016](#)) showed that multiple fairness measures are incompatible with each other and cannot be satisfied simultaneously. This has led to a growing recognition that algorithmic fairness is not a one-size-fits-all solution, and that different contexts may require different fairness measures, but there is not yet a consensus as to how to select the appropriate measure for a given system.

In an effort to develop a more principled approach to algorithmic fairness that will inform how fairness measures are selected and applied, researchers have turned to the field of distributive justice for guidance. A distributive theory of justice is a normative framework that provides principles and criteria for allocating benefits and burdens among individuals or groups within a society, with the aim of achieving a just and fair distribution. The field can be seen as polarized along an axis from liberal egalitarianism to entitlement theory. Under the liberal view, commonly associated with John Rawls, the chief objective of justice is to equalize allocation across all individuals in a population. In contrast, the entitlement view, associated with Robert Nozick, emphasizes the importance of individual property rights and the freedom to exchange goods and services without interference.

The relationship between algorithmic fairness measures and distributive justice is not yet well

understood, but several recent papers have begun to explore this connection. (Binns, 2018), (Hertweck et al., 2024), and (Kuppler et al., 2021) have all examined the relationship between algorithmic fairness and egalitarian concepts of justice, showing that fairness measures that measure disparities in the outcome distribution over social groups are predicated on certain assumptions about equality as a foundation for justice. (Baumann et al., 2023) develops this further, but rightly points out that this approach is limited to one particular view of justice, and that it’s unclear how a structurally different theory of distributive justice such as Nozick’s entitlement theory could be formalized as a fairness measure. Given that entitlement justice is a prominent area of inquiry in political philosophy that addresses a broad range concerns not covered in liberal justice, it is worth investigating how to close this gap. In particular, how do issues of entitlement appear in algorithmic decision making? How can these concerns be encoded by algorithmic fairness measures? And what do we stand to lose or gain by conceptualizing algorithmic fairness through the entitlement lens?

In this paper, we will carefully examine the relationship between algorithmic fairness and libertarian justice, and develop a formalism that clarifies the relationship between the two. We will demonstrate that entitlement justice can be encoded within a measure of algorithmic fairness by formulating the problem of fairness on an individual level rather than a group level. We will show that doing so offers a nuanced and context-sensitive means of understanding algorithmic fairness, and that it can be used to inform the selection of appropriate fairness measures for decision making systems.

The rest of this paper is organized as follows. In Section 2, we provide an overview of the existing literature on algorithmic fairness and distributive justice. We draw on the formalism from (Kuppler et al., 2021) and (Corbett-Davies et al., 2023) to create a unified model for understanding algorithmic fairness and distributive justice consistently with each other. In Section 3, we introduce the concept of entitlement justice and discuss its historical development. We contrast entitlement theory with liberal egalitarianism to identify the critical elements of entitlement which must be represented in account of algorithmic fairness, and confront the traditional objections to entitlement theory. In Section 4, we propose a new framework for understanding algorithmic fairness through the lens of entitlement justice. We analyze the implications of this framework for existing algorithmic fairness measures and show an example of how it can be applied to a real-world case study. Finally, in Section 5, we conclude with a discussion of the broader implications of our work and suggest directions for future research.

2 Background

In this paper, we will center our attention on a class of decision problems corresponding to the following formalism. There exists some population of individuals $I = \{i_1, i_2, \dots, i_n\}$ over whom we must distribute some finitely divisible amount of a resource, R . A *decision rule* is a mapping $d : I \rightarrow \{0, 1\}$, under which i_n receives R iff $d(i_n) = 1$. As an example, consider the process of allocating loans over a pool of applicants — for each individual in the applicant pool, we have a binary choice to either approve or deny the loan, and only on acceptance does the individual receive capital.

An algorithmic decision maker in this setup is a system, particularly a technical system,

under which a decision rule is implemented as a set of steps which are applied identically to all individuals. This broadly consists of two tasks. First, a set of covariates $X = \{x_1, x_2, \dots, x_n\}$ must be collected for each individual. Then, a classification scheme $f : X \rightarrow \{0, 1\}$ must be used to map each individual to an outcome. Returning to the loan case, a simple example would be the following: Approve a loan to each individual whose household income exceeds the projected cost of living in the geographic location of their residence by at least \$10,000 per year. Then our algorithm implementing the decision rule for each applicant takes the following shape:

- Collect covariates $X = \{\text{household income, geographic location of residence}\}$
- Let $C(x_2) = \text{projected cost of living in location of residence}$, $x_1 = \text{household income}$
- $f(X) = x_2 - C(x_1) \geq 10000$

It is important to note that not all problem domains to which algorithmic decision making is applied can be formulated in this way. For example, applications of AI in natural language translation may not be easily formulated in terms of resource allocation, but the reader may still be concerned with the perpetuation of social biases through the decisions made in translation — for example, issues of underrepresentation of social groups in translated media. While these cases warrant further study, we will not consider them here.

2.1 Algorithmic Fairness Measures

Algorithmic fairness measures as they're presented in the literature operate on the classification scheme f and a set of morally protected characteristics such as race or gender $A \subseteq X$, attempting to enforce constraints on how decisions can be sensitive to such characteristics. An overview of measures commonly discussed in the literature is presented below.

A critical point to note in the discussion of algorithmic fairness is the distinction between individual and group fairness. In group fairness, we form groups of individuals based on the protected attributes A , and measure fairness as a statistical property on the distribution of outcomes across these groups. For example, we may measure the proportion of individuals from each group who received a positive outcome from the classification function. In individual fairness, on the other hand, we attempt to judge whether like individuals are being treated alike by the algorithm. For example, two individuals who are identical but for their race should receive the same classification outcome in decision-making domains where race is not relevant.

Both group fairness and individual fairness measures have conceptual shortcomings. The coarse-grained nature of group fairness measures means that there is no guarantee that individuals within a group are treated fairly, only that the algorithm is statistically unbiased in the aggregate, discarding the potential for bias against particular individuals in the group. Individual fairness measures, on the other hand, require a notion of similarity between individuals which is often difficult to define. Never are two individuals truly identical but for their race, and the choice of a similarity metric can encode normative assumptions about the importance of different attributes in determining the outcome of the classification function.

To set the scene for further discussion of how algorithmic fairness is connected to distributive justice, a survey of common algorithmic fairness measures together with their strengths and

weaknesses is presented below. We will discuss these measures in the context of the decision problem formulated above. The goal of this discussion is to highlight the limitations of existing algorithmic fairness measures in capturing the normative concerns of entitlement justice. For a more exhaustive presentation of existing measures, see [Corbett-Davies et al. \(2023\)](#).

As a running example, consider the case of a hiring algorithm. The decision problem is to map a set of applicants to hiring decision. Our algorithm must implement a function to do so based only on the personal information presented in their resumes, which is their education, experience, and cover letter. We want to ensure that our hiring practices are fair with respect to race, and gender.

Definition 1. *Fairness Through Unawareness* — f satisfies fairness through unawareness iff

$$A = \emptyset$$

In other words, the classification function may not receive any morally protected covariates as inputs. So, for our hiring admissions case, we would be forced to remove race and gender from X and only consider education and experience.

This notion is intuitively appealing — how can the algorithm discriminate against one based on their race if it doesn't know their race? However, it is clear that this measure is not sufficient to ensure fairness. For example, if one attended a historically black college, their education may function as a proxy for their race, allowing bias to remain in the algorithm.

Definition 2. *Demographic Parity* — f satisfies demographic parity iff

$$P[f(X) = 1|A = a] = P[f(X) = 1] \quad \forall a \in A$$

([Dwork et al., 2012](#)).

Demographic parity holds that the probability of a positive output ($f(X) = 1$) should be statistically independent of the protected attributes. This is an easily understandable and measurable criteria for fairness. At a first look, it is appealing — the same number of individuals from each race will be successful in seeking jobs at a particular company. Indeed, in a situation like selection of individuals for a representative committee this measure is appropriate. However in other circumstances such as hiring or parole, difficulties arise.

Under demographic parity, the probability of success between groups must be balanced, but if the base rates of success are unequal between groups, this leads to poor outcomes. For example, if women are much more qualified for a job on average than men, then demographic parity will require that we hire less qualified men in place of more qualified women in order to balance the probability of being selected between gender groups. In other words, the false positive rate will be very high for men while the false negative rate will be very high for women, creating a severely unfair practice ([Barocas et al., 2017](#)).

Definition 3. *Equalized Odds* — f satisfies equalized odds if given a true outcome o for each individuals, we have

$$P[f(X) = 1|A = a, o = 1] = P[f(X) = 1|A = b, o = 1] \quad \forall a, b \in A$$

$$P[f(X) = 1|A = a, o = 0] = P[f(X) = 1|A = b, o = 0] \quad \forall a, b \in A$$

([Hardt et al., 2016](#)).

In cases where the output value of the classifier f disagrees with the true outcome o , we have either a false position or false negative. Equalized odds requires that the true positive and false positive rates be balanced between groups. This is often thought to ensure there is no disparate mistreatment across groups. In our hiring case, for example, equalized odds ensures that no one racial group is more likely to be falsely rejected or erroneously hired than another. Given one has been rejected, the probability it was a wrongful rejection is equal regardless of their race or gender.

Equalized odds has two concerns associated with it. Firstly, the condition becomes difficult to satisfy when the base rates of success are unequal between groups, as differing base rates imply differing false positive and false negative rates. In our hiring case, once again imagine that women are much more qualified for a job on average than men. This condition in the population would lead to a model distributing more false positives to women and more false negatives to men, and in order to rectify this imbalance, we would need to use a different threshold of qualification for women than for men, sacrificing the accuracy of the model.

A second concern with equalized odds is that it is easily manipulable. Consider the following example from criminal justice. We have a distribution rule which says to allocate a parole to a prisoner if they are very unlikely to recidivate. Due to a history of discriminatory practices and social marginalization, black prisoners are judged as having higher a risk of recidivism much higher than white defendants (Crime and Alliance, 2023). As a result, allocating a parole to a white prisoner has a lower likelihood of being a false positive. Therefore, one could achieve equal false positive rates by *adding* false positives to the white portion of the dataset, resulting in an increase in the number of undeserving white prisoners receiving parole.

As an example of equalized odds gone wrong, consider the COMPAS algorithm (Angwin et al., 2016). COMPAS was calibrated to have equal predictive accuracy across racial groups, but this resulted in a higher false positive rate for black defendants and a much lower false positive rate for white defendants due to unequal base rates.

Definition 4. *Counterfactual Fairness* — f satisfies counterfactual fairness iff

$$P[f_{A \leftarrow a}(X) = 1 | X = x, A = a] = P[f_{A \leftarrow b}(X) = 1 | X = x, A = a] \quad \forall a, b \in A$$

(Kusner et al., 2018). Where $P(f_{A \leftarrow a})$ is the counterfactual value of f if A were set to a .

Borrowing from the language of causal inference, counterfactual fairness posits that the protected attributes may not have any causal effect on the outcome of the classification function. This measure operates on the individual counterfactual—would the output of the classification function have been different if the individual had been different according to some protected attribute? This is a highly appealing notion of fairness. If their protected characteristics do not in any way cause their outcome, then it is difficult to argue that one has been discriminated against. However, this measure is difficult to implement in practice.

Counterfactual fairness is often critiqued based on the difficulty and potential subjectivity of detecting causal links between variables. Recent work on the social construction of demographic variables reveals that causal modeling may have an inherently normative basis (Hu, forthcoming), and even if these issues are set aside, the computational expense of causal discovery can create issues of practicality.

This discussion of dominant algorithmic fairness measures and their critiques reveal that there is no one-size-fits-all solution to the problem of algorithmic fairness. Each measure has its own strengths and weaknesses, and the choice of measure will depend on the specific context in which the algorithm is being applied. Intuitively, one would like to measure a given classifier against a range of measures to ensure it is fair in a variety of ways. However, as previously mentioned, it is impossible to satisfy multiple measures simultaneously, and so the choice of measure becomes a critical one. How should one select a measure? What are the philosophical considerations that should guide this choice? In hopes of developing a more nuanced and structured approach to these questions, we turn to work in the philosophy of distributive justice.

2.2 Theories of Distributive Justice

Distributive justice is a philosophical field of inquiry that examines how to define a fair allocation of goods and resources across a society. A fully fledged account of distributive justice must answer a number of questions. Who should receive those resources which are highly scarce? When and why is it allowed for one person to have more of something than another? By what mechanism can resources be redistributed to achieve justice?

Given that distributive justice defines how fair decisions about allocations can be made, within the formalism we've presented, its role is to broadly define the decision rule which may then be implemented algorithmically. As described in section 1, the dominant theory of distributive justice used in connection with algorithmic decision making is John Rawls' theory of liberal egalitarianism, which we will present here.

Rawls introduces his account of justice through a thought-experiment called the veil of ignorance. In this thought-experiment, one is asked to imagine themselves in a pre-societal world, working in collaboration with a number of others to determine how resources should be allocated across society once it begins. Critically, all those involved in designing this distribution of goods are unaware of what their own position and endowments in society will be. One may find themselves endowed with a high level of intelligence, or a valuable skill, or wealth at birth, or one may find themselves with none of these things, or the opposite. Without knowing which of these positions one will occupy, Rawls argues that one will be motivated to design a society in which the following two principles are satisfied:

1. Each person has an inalienable right to the most extensive basic liberties compatible with equal liberty for all.
2. Social and economic inequalities are to be arranged so that they are both to the greatest benefit of the least advantaged, to offices open under fair equality of opportunity. (dubbed the *difference principle*).

Rawls refers to the group of individuals designing the society from behind the veil of ignorance as the *original position*. The argument from the original position results in citizens living under a social contract which is guided by the two principles given above. The principles allow us to measure, for any given distribution of resources across society, whether or not the distribution is fair. If the distribution is not fair, then Rawls endorses a program of redistribution to bring

the distribution into line with the principles. For example, a society with a high-level of wealth inequality is in violation of the difference principle — the wealth gap represents an economic inequality which is not to the greatest benefit of the least advantaged. In this case, Rawls would endorse a program of redistribution to balance the wealth across the society in accord with the principles given above.

This type of distributive justice theory is what we refer to as an *end-state* theory of justice. The distribution of goods across society represents a discretely evolving state of affairs, and the role of the theory is to determine whether or not each state is just. Let us consider this view in light of the decision problem we formalized above. Liberal egalitarianism tells us that the decision rule d must be such that either all individuals receive the resource of allocation equally, or that inequalities in the allocation of resources must be to the benefit of the least advantaged. Several of the fairness criteria in the literature on algorithmic fairness can be seen as implementing the first condition in terms of equality of opportunity — by regulating the extent and manner in which protected attributes can influence the outcome of the decision, we attempt to ensure that all individuals are receive equal basic rights of opportunity in the decision process. However, whether or not the difference principle is satisfied by these measures is less clear.

In our loan case, for example, we may be concerned with ensuring that every individual receives equal opportunity to a loan. However, if the decision rule is such that only individuals with a household income above \$100,000 per year, or those who are members of a particular race, are able to receive a loan, then this clearly doesn't provide equal opportunity to all. A measure like demographic parity ensures that individuals from each protected group are equally likely to receive a loan, therefore balancing opportunity across groups. However, whether or not this satisfies the difference principle depends on the base rates of success across groups. If the base rate of loan default is higher for men than for women, demographic parity may require that we give loans to some men who are less likely to pay back their loans rather than qualified women. Absent from this metric is any consideration of who is considered to be the least advantaged in the situation. By Rawls, we should be allocating capital to those individuals who have the least access to capital, a condition which is not satisfied by demographic parity. (Hertweck et al., 2023) shows how Rawls' theory can be more fully captured in a group-based fairness metric, which we will not go through the details of here. However, the key point echoed throughout the literature is that the common measures of algorithmic fairness stem from egalitarianism, and it is unclear how to extend these measures to a theory like Nozick's which does not operate on an end-state theory of justice.

3 Entitlement Justice

An entitlement theory of justice is a distributive theory of justice which posits the following distribution rule: Allocate amount R of resource X to agent A if and only if A is entitled to R of X . An entitlement in this context is a *property right* held by the agent over the resource. Different entitlement theories of justice differ in the criteria they use to determine entitlements, and the concept of property rights they endorse. Here we will detail the entitlement theory of justice as proposed by Nozick (1974), its issues, and how it compares to the Rawlsian theory of justice, then discuss more recent efforts at reconciling the theory with the demands of justice.

3.1 Nozick's Entitlement Theory of Justice

Nozick's entitlement theory of justice, often called the concept of libertarian justice, is a theory of justice that was developed as a fundamental challenge to Rawls's liberal egalitarianism. On the liberal egalitarian view, ensuring justice is an inherently redistributive task. The justice of a distribution of resources is determined by the extent to which it is equal over individuals, and there is an implied moral responsibility to redistribute resources to those who lack them to increase the overall equality of the distribution. This ideology provides a strong defense of taxation and welfare programs, which redistribute resources in order to flatten the distribution of wealth (Rawls, 1971).

Libertarian justice takes issue with the consequences of adopting this view. Nozick asks us to consider a thought experiment. Suppose we began with an equal distribution of resources across society. People in this society have the freedom to choose how to use their resources, and to exchange them with others as they feel is fair. Many people are willing to pay to see Wilt Chamberlain play basketball, and so they each pay him a small amount of money to see him play. Over time, Chamberlain will accumulate a large sum of money through his efforts. The distribution of resources in the society will no longer be equal, but will be skewed towards Chamberlain. On the egalitarian account, this excess wealth that Chamberlain has accumulated is unjust, and must be taken and redistributed across society. On the libertarian view, however, Chamberlain has gained an entitlement to his accumulated wealth, and to take it away from him is akin to stealing. After all, if this wealth is taken away from him, then he will have received nothing for his efforts, and enjoyed no fruits of his labor.

In Nozick's theory, people gain entitlements over resources in accordance with 3 principles:

1. The principle of justice in acquisition: A person who acquires a resource through a just process is entitled to that resource. A process of acquisition is just if the acquisition is in accordance with Lockean proviso (discussed below).
2. The principle of justice in transfer: A person who acquires a resource through a just transfer is entitled to that resource. A transfer is just if the transfer is voluntary and the resource is transferred from someone who is entitled to it.
3. The principle of rectification: A person who acquires a resource through the rectification of a prior injustice is entitled to that resource. Rectification must be proportional to the injustice which is being rectified.

On analysis, one will see that a key difference between this libertarian view and the liberal egalitarian view is the fundamental unit of justice. For the liberal egalitarian, justice is realized in the distribution of resources itself. This approach is referred to as a patterned or *end-state* view of justice. For the libertarian, justice is realized in the process by which resources are acquired and transferred. This approach is referred to as a *historical* theory of justice. In order to determine if the current state of affairs is just with respect to a particular holding, one must trace the history of that holding back to its original acquisition, and ensure that each step in the process was just. For Nozick, any end-state theory of justice is inherently flawed, as it requires the restriction of individual liberties (Hendberg, 1977). It is plain that this view of justice hinges strongly on being

able to identify and justify the initial acquisition of resources, else the theory can say nothing about the justice of the current distribution of resources.

3.2 The Justification of Acquisition

For Nozick the Lockean proviso underscored the principle of justice in acquisition. The proviso contains two parts. The first part is a mechanism for justifying the initial acquisition of resources. It begins with the inherent right of self-ownership that all individuals possess. Locke argued that when an individual mixed their own labor with a resource, they transferred some of themselves into the resource, and so extended their right of self-ownership over the resource, thereby obtaining an entitlement to it. The second part of the proviso, almost as an afterthought, is a restriction on the extent to which resources can be acquired. It states that a person can only acquire a resource if there is enough and as good left over for others. This restriction is necessary to ensure that the acquisition of resources does not infringe on the rights of others to acquire resources.

Other accounts of entitlement justice have used different mechanisms to justify the acquisition of resources. (Mack, 1990) proposed that the acquisition of resources could be justified as a separate unalienable right that all individuals possess. (van der Veen and Van Parijs, 1985) proposed that the acquisition of resources could be justified consequentially by the net utility that the acquisition brings to society. In general, Van Der Veen showed that given a particular type of holding, one can specify a theory of entitlement justice with a corresponding utilitarian theory of acquisition that can be used as a basis for determining entitlements.

3.3 Critiques of Entitlement

Nozick's entitlement theory is heavily criticized for its foundation in the Lockean proviso. The final clause of the proviso provides a restriction on the extent to which resources can be acquired, but is a weak restriction that makes it difficult to justify the acquisition of resources in practice. There are two mechanisms by which the proviso as it pertains to Nozick's entitlement theory breaks down.

Firstly, the proviso is a weak and vague restriction. It was written in an era when it seemed plausible that individuals would frequently be staking claim over new possession in the wilderness, in particular parcels of land. However, in the modern setting, there are few unclaimed natural resources, and those that exist come under heavy contention for acquisition. The proviso does not provide a clear mechanism for dividing up the resources in this case, and it seems entirely unlikely that one can satisfy both aspects of the proviso concurrently (Fried, 2004).

Secondly, the proviso has a problem dealing with the issue of surplus value. According to the proviso, when an individual acquires a resource, they acquire it by instilling some valuable portion of themselves into the resource. There is thus a fixed amount of value transferred onto the resource through the person's labor. However, in a free market like the one Nozick describes in his theory of entitlement justice, the value of a resource is not fixed, it is dictated by market forces. If an individual acquires a resource and then the value of that resource increases due to scarcity or high demand, then the individual can trade their resource and gain entitlement over property with a value greater than that which they instilled into their original acquisition (Fried, 1995).

These issues provide a strong challenge to Nozick's entitlement theory as they can result in disastrous consequences. Besides the proviso, Nozick's theory may be criticized for its potential to justify unacceptable outcomes through transfer. For example, someone who is starving may "voluntarily" agree to trade property for food whose value is far below the value of the property, and per Nozick, this trade might be considered just. Critically, this does not spell the end for the entitlement theory of justice, but it does suggest that the underlying theory of property rights for a successful entitlement theory of justice as well as the restrictions on the types of transfers it can justify must be more nuanced than what Nozick proposed.

3.4 Instrumental Property Rights

Successors of Nozick have sought to address these issues by replacing the Lockean proviso with an alternate theory of property rights. [van der Veen and Van Parijs \(1985\)](#) showed that entitlement systems existed on a spectrum such that the theory of property rights at the base could be tailored to the resource being distributed. For example, the precise theory of property rights for land might be different than that for money, or for a scarce natural resource. This observation suggests that a successful entitlement theory of justice must be based on a theory of property rights that is situated in the context of the resource being distributed. This observation is echoed by [Fried \(2004\)](#) who argues that the theory of property rights must be tailored to the resource being distributed, and that the theory must be created with the full scope of its consequences in mind.

Regardless of the theory of property rights used, to overcome the challenges of Nozick's theory, one must prevent an entitlement theory from justifying morally unacceptable outcomes as [Fried \(1995\)](#) worried about in the case of Nozick's theory.. [Sen \(1988\)](#) shows that while the interpretation of property rights as inherently valuable and inalienable leads to severe issues of poverty and hunger, the interpretation of property rights as *instrumental* rights, which are valuable only insofar as they lead to particular desired outcomes, can be used to develop systems of entitlement without such issues. Instrumental property rights cannot supersede the demands of basic necessity for all agents, and so can be used to develop systems of entitlement which protect property rights while avoiding issues that arise alongside emergent wealth disparity from free market transactions.

Combining these lessons, we realize that a successful modern theory of entitlement justice is one situated atop a theory of domain specific and instrumental property rights. For a given type of holding or resource, the theory of property rights must be tailored to the resource, and must be created and enforced with the full scope of its consequences in mind. This approach allows for the development of a theory of entitlement justice that is both normatively justifiable and practically applicable in the modern world, and thus could be used to inform the design of algorithmic fairness measures. A critical result here is that the theory of property rights used to define fairness over a particular domain must be able to be superseded by higher order concerns and should be defended to the population it is applied to to gain acceptance.

3.5 Contrast with Liberal Approach

To begin to craft an account of algorithmic fairness through the lens of entitlement justice, it is useful to contrast entitlement theory and the liberal approach to see the critical dimensions along

which they differ. These differences provide a clear set of concerns of entitlement justice that must be addressed by a fairness measure designed to implement entitlement justice.

- Historical vs. end-state — Under an entitlement theory, the justice of a distribution is determined by the history of how the distribution came to be through acquisition and transfer. In contrast under a liberal egalitarian approach, the justice of a distribution is determined by the current state of the distribution itself.
- Individual and collective responsibility — Under an entitlement theory, there is a heavy focus on the actions and properties of individuals which give rise to their entitlements. It is the individual who acquires or trades for resources, and thus it is the individual who is responsible for their own state of affairs. In contrast, under liberal egalitarianism, the only relevant properties of an individual are their current holdings or status in society, and it is a collective responsibility to ensure that resources are distributed according to the demands of justice.
- Redistribution — Under an entitlement theory, redistribution of resources from the more fortunate to the less fortunate is not a moral imperative. In fact, redistribution is unjust if it is not done willingly on the part of the more fortunate. In contrast, Rawls' difference principle explicitly mandates redistribution of resources to the less fortunate.

4 Entitlement Fairness

In order to understand the relationship between algorithmic fairness and entitlement justice, it is important to first analyze the role of the algorithmic decision-maker in the context of entitlement systems. On the entitlement approach, decisions about allocations are made entirely based on property rights. Therefore the task of a decision-maker within our decision problem is clear — the decision-maker becomes a *property rights oracle*. Given a resource and information about a population of individuals, the job of the decision-maker is to determine which individuals hold property rights over the resource. The role of fairness is thus a bit different than under other theories of justice, because we do not start from the assumption of any sort of equality across our population. What type of assumption should we start with instead? To understand this, we will analyze the points of contrast we have drawn between the Rawlsian and entitlement theories of justice.

First, though, we must elucidate the meaning of property as it will be used in this investigation. In the modern legal system, property is typically restricted to a somewhat narrow range of physical objects, financial assets, and intellectual creations of the mind. Here, though, we will take a broader view of property that encompasses any resource one might be entitled to claim and which algorithmic decision-making might be used to allocate. This includes all of the resources we typically think of as property, but also includes things like admission to a particular college. In brief, property will refer to any finite desirable resource that can be allocated to an individual. This is a much broader understanding of property to allow the entitlement framework to be applied across a wide range of decision-making domains.

4.1 Fairness in Process

One critical dimension of entitlement justice is that it governs the process that gives rise to a distribution of resources rather than the distribution itself. This means that decisions about allocation which govern whether a particular individual is able to acquire a given resource must be made in accordance with a fair process. Rather than answering whether or not a given decision-maker outputs a fair distribution, we must instead ask whether the manner in which decisions about allocations are made is fair. This is a subtle but important distinction. Under the entitlement approach, the output distribution is allowed to be heavily skewed in favor of some individuals or groups if it represents the true distribution of property rights, but the process by which the output distribution is arrived at must be in accordance with the principles of just acquisition.

As an example, consider the case of college admissions. Under the liberal egalitarian approach, we might ask whether the output distribution of our algorithm is fair by asking whether it results in a roughly equal number of students admitted between race A and race B. Under the entitlement approach, however, there is no reason to ask this question—we might find that the students entitled to admissions are 90% members of race B. Instead, the relevant questions are how race is used in the admissions process — are students from race A subject to the same rules of acquisition? Drawing on the framework of the Lockean proviso, given that two students, one from each race, have expended equal effort in their applications and studies, are the values endowed in their applications treated as equal? It is evident that to encode entitlement fairness, we cannot simply measure the outputs of an algorithm, but rather must analyze the treatment of protected features of individuals within the algorithm itself.

4.2 Individual Responsibility

Our second critical dimension of entitlement justice is that it places emphasis on the actions and properties of individuals, while de-emphasizing the role of group membership. Under the entitlement approach, individuals perform actions that give rise to or forfeit their property rights. Fairness must therefore be fundamentally based on a set of features of individuals which are relevant to the process of determining property rights. These features will generally not be simple demographic features or features encoded in the input covariates in a straightforward way, but rather will be a set of more nuanced features that must be predicted on the input data by a complicated high-dimensional model. For example, in the case of college admissions you may want to predict a feature like “academic potential” which will be difficult to extract. These features can be sorted into two broad categories:

- *Positive Entitlement Features*: These are features of an individual that give rise to property rights. For example, an individual’s effort in an application process might be a positive entitlement feature that gives rise to their right to be admitted to a college.
- *Negative Entitlement Features*: These are features of an individual that forfeit their property rights. For example, if an individual cheats on their application, or performs very poorly on an entrance exam, these might be negative entitlement features that forfeit their right to be admitted to a college.

Notice that given the full set of relevant positive and negative features of an individual that determine their property rights, we should be able to fully determine the individual's entitlement to a resource and complete the decision problem. In other words, once we have identified the features and computed them on an individual, the output of the decision maker can be fully specified by these features alone. This lends itself to a natural understanding of fairness in process — the process of mapping an individual to their decision should be fully decided through the morally relevant features identified. These features should be explicitly justified and made transparent to the individuals affected by the decision.

Returning to our discussion of modern entitlement theories, we can recognize that the identification of relevant features is how context-specificity will enter into our account of algorithmic fairness. In each problem domain for algorithmic decision-making, there will be a different set of positive and negative features that are relevant to the entitlement being decided. By identifying and justifying the set of features relevant in each domain, we allow our account of fairness to be sensitive to the context and nuances of the problem at hand. This is a powerful contrast to typical approaches to fairness under which we attempt to identify a universal measure of fairness that can be applied across all domains.

In application, this implies a particular set of structural conditions that must be followed to implement an algorithm which is fair under the entitlement approach. A classification scheme must be developed that first computes the value of each of the relevant positive and negative entitlement features for each individual in the population. Then a decision must be reached through only those computed features, in isolation from the full input data to the algorithm. This may seem to stand in stark contrast to the way that many machine learning algorithms are currently developed. For example, in a typical supervised learning setting, a model is trained to map a set of inputs to a set of outputs according to a high-dimensions loss function, with little regard for the manner in which the inputs are processed. However, the internal structure of neural networks and other machine learning models gives rise to a set of features computed by the model, constituting a lower-dimensional representation of the input data (Liu, 2018). The approach we suggest here can be thought of as a way of manually specifying a lower dimensional set of features that are morally relevant to implement individual fairness over as a way of exercising control over *how* the model makes decisions in order to implement fairness as a process.

What would selection of these features look like in practice? Consider again the case of college admissions, and in particular, how decision are made about admissions of students who have less access to resources and academic opportunities. There are several features that are not straightforwardly encoded in the input data but are certainly relevant to the entitlement

- Firstly, we might want to extract a feature that captures the notion of current academic performance. This is likely a function of GPA, test scores, and other typical academic indicators, and is justified by an appeal to the idea that students who perform well academically are more likely to succeed in college. Likelihood of success is made relevant by the fact that individuals who are likely to succeed at the university return gains to the university, and therefore justify their admissions in a free market of talent. This positive feature represents an action of the individual that gives rise (“earns”) their right to be admitted. This stands in contrast to the Rawlsian approach, under which there is no broad support for meritocratic systems. Academic performance for Rawls is likely a product of natural talent and luck,

and therefore by Rawlsian standards, is not a valid basis to allocate admissions inequitably.

- In contrast, we should also say that a student is entitled to admissions if they have demonstrated a stronger work ethic and commitment level than their peers, even if they attended a lower income school and thereby has less access to advanced classes and tutoring resources. Effort and commitment demonstrate a greater value endowed into an application, and therefore a higher degree of entitlement. Here the difference with Rawls is more subtle. Rawls would agree that the student from a lower income school and with less access to resources should be promoted in admissions, but as a redistributive effort rather than as a recognition of the value endowed in the application. On the entitlement approach, two students who have expended equal work and effort should be treated equally in this respect, regardless of their background.
- Finally, we might also want to consider a student's cultural fit and addition to the campus community. This feature reflects a student's entitlement on the basis of more than just academic merit — a student who provides cultural value to a university provides similar, though reduced value to the university as a student who provides academic value, and therefore has a similar, though lesser, entitlement. This feature can be subdivided into both a positive and negative feature. A student who has demonstrated positive cultural value to the community gains an entitlement, while a student who has demonstrated negative cultural value forfeits their entitlement earned through academic merit and effort. Here Rawls again pulls apart from the entitlement approach. Under the Rawlsian view, cultural fit, so far as it is a product of natural tendencies and socialization, is not a valid basis for admissions.

An algorithm meant to implement entitlement fairness in college admissions would then consist of predictors which extract each of these values for each individual, and then a system of mapping these values to a decision about admissions.

Now, one critical dimension of fairness remains to be discussed — how do we ensure that the predictors which compute the value of the relevant features for each individual are themselves fair? To understand this question, we should delve into the basis of the entitlement features themselves. Positive features of an individual give rise to property rights through the mixing of one's self with the subject of the entitlement itself. If one's college application demonstrates a strong work ethic, the individual has spent significant effort to endow their application with value derived from their own self-ownership. The right to self-ownership is a fundamental principle that does represent one form of equality baked into the entitlement approach — we are required to treat the endowment of value derived from self-ownership as equal across all individuals, irrespective of their characteristics or group membership. Notice that I care about this on an individual basis — if for just one person I value their self-ownership less than another, I have failed to respect their property rights. As a result, it seems clear that the appropriate way to ensure that predictors of positive entitlement features are fair is to ensure that they satisfy counterfactual independence. Given two individuals do not differ significantly in any manner relevant to the positive feature being predicted, the predictor should output the same value for both individuals.

For negative entitlement features, the situation is a bit different. Rather than gaining an entitlement through the mixing of one's self with the subject of the entitlement, the basis for

negative entitlement features is the failure to respect the inherent rights of others. For example, if an individual cheats on entrance exam scores that they include in their application, they are attempting a coercive act that violates the rights of other students who are applying honestly and in doing so, they forfeit their ability to gain entitlement over admissions. In this case, our chief fairness concern is with errors in the predictor. Any error in the predictor of a negative feature will result in an infraction on the property rights of some individual in the population—therefore we must ensure that the predictor is as accurate as possible. Given that the predictor will necessarily be imperfect, we are also required by entitlement theory to grant individuals the right to appeal the decision of the predictor, so that if they’ve been wrongfully accused of a negative entitlement feature, they may present their case for rectification. Thus fairness in the negative feature predictor is a matter of ensuring there is a low error rate, individuals are aware of the features being predicted, and that they have the right to appeal the decision of the predictor.

4.3 Entitlement Fairness and Redistribution

A third critical dimension of entitlement justice is that it rejects the notion of redistribution. Under the entitlement approach, once property rights have been established, they must be respected and protected, and no attempt should be made to redistribute resources in order to achieve a more equal distribution. Note that this is already successfully encoded in the system we have described. Once the set of morally relevant features have been identified, the decision problem output must be entirely separated from the input data given the relevant features. This means that no redistributive scheme may be implemented after the property rights decision that attempts to balance the distribution of resources over less fortunate individuals.

We are, however, offered a mechanism through entitlement theory for improving outcomes for those who have been wrongfully disadvantaged through the principle of rectification. For example, if we find that a particular group of individuals has been systematically excluded from a resource due to a historical injustice, rectification allows us to then consider membership in that group to be a positive entitlement feature, and to thereby account for it within the decision problem. This is a powerful mechanism for addressing historical injustices and allows us to consider the broader social context in which our algorithm is situated.

4.4 Measuring Fairness

Having developed a qualitative account of entitlement fairness, we may now formalize our account in a way that allows us to measure it. We can define a measure of entitlement fairness as follows. Given our typical decision problem, identify a set of morally relevant features of each individual in the population, $V = \{v_1, v_2, \dots, v_j\}$. These features should be partitioned into two sets, V^+ and V^- where V^+ contains the positive entitlement features and V^- contains the negative entitlement features relevant to the problem domain. Implement a predictor for each feature, $v_j = \hat{p}_j(X)$. Now, these are the conditions for entitlement fairness:

1. $P[f(X_1) = 1|V_1] = P[f(X_2) = 1|V_1] \forall X_1, X_2$ (Independence Condition)
2. $P[p_{j,A \leftarrow a}(X) = 1|X = x, A = a] = P[p_{j,A \leftarrow b}(X) = 1|X = x, A = a] \forall a, b \in A, v_j \in V^+$ (Positive Feature Counterfactual Condition)

3. $p_j = \underset{p}{\operatorname{argmin}} L(p_j, v_j)$ where L is a loss function giving the error of the predictor p_j (Negative Feature Accuracy Condition)

In addition to these mathematical conditions, we also require that

1. The values of the features predicted by the model be made transparent to the individuals affected by the decision.
2. The individuals affected by the decision have suitable opportunities to appeal the decision of the predictor for rectification.
3. The moral importance of the features used in the decision problem be justified and defended in a public forum.

These conditions together provide a context-sensitive and principled approach to algorithmic fairness derived from the entitlement approach to justice, and provide a framework for understanding how to select and apply fairness criteria in algorithmic decision-making.

5 Discussion

The entitlement fairness measure derived here is a novel approach to algorithmic fairness and represents a first attempt to formalize an account of justice that is consistent with a historical rather than end-state theory of justice. Here we will briefly discuss the benefits and limitations of this approach, and suggest some directions for future research.

5.1 Benefits of the Entitlement Fairness Measure

The entitlement fairness framework has several advantages over existing fairness measures. Firstly, it is a measure which has baked in context sensitivity. In selected a set of positive and negative features which are morally relevant to the decision problem at hand, the fairness measure is able to capture concerns about the decision problem that could not be captured by simple disparity measures across groups.

In each problem domain the framework is applied to, the system designer must make conscious choices about the features that are allowed to be used in the decision process, and must defend those choices in light of the particulars of the domain. There can be no “lift and shift” of the entitlement fairness measure from one system to another, meaning that the application of the measure is always done consciously and with moral justification. This stands in contrast to existing measures of algorithmic fairness, which can be applied to any system without extensive consideration of the idiosyncrasies of the problem domain. For example, consider the case of the COMPAS algorithm, in which the original designers of the algorithm measured its fairness using predictive parity across races, which is a measure of accuracy across groups. predictive parity is appropriate for some problem domains, particularly those where the base rates of the predicted outcome are the same across groups. The COMPAS system was later revealed to be discriminatory based on differing false positive rates across groups due to a large discrepancy in the base rate of recidivism between racial groups. When using the entitlement fairness measure,

there may be no simple lift and shift of the measure from one system to another as there was for the predictive parity metric into the COMPAS system, and as a result cases where an unsuitable measure is applied to a system are less likely to occur.

Secondly, the entitlement fairness measure is one that is well-aware of its own limitations and assumptions. Part of the requirement of entitlement fairness is the baked-in necessity for a system of rectification for entitlement violations. This holds the system implementer accountable for the decisions made by the system and for providing remedies for those harmed by the system. A system can not be entitlement fair without a way to appeal and rectify the decision made by the system, meaning that the system is always open to scrutiny and correction. One key result is that one can't simply measure their system against one definition of fairness, declare it to be fair, and then put the system into production without oversight and recourse. This is a key difference from existing fairness measures in that it increases the accountability of the system implementer to the system's users to maintain the system's fairness over time and provide support for their users.

The final benefit of the entitlement fairness measure is that it gives a high level of transparency to the population it is employed on. Not only do individuals who decisions are made on have recourse to appeal their decisions, but they have a clear understanding of the features that are being used to make the decision and the moral justification for those features. This is a key result of the entitlement fairness measure, as it allows some small level of explanation to be provided to the individuals who are affected by the system. This is a key difference from existing fairness measures, which may help to soothe concerns about how some particular variables are used in the decision process, but not actually provide or identify the features that are used in the decision process.

5.2 Limitations of the Entitlement Fairness Measure

The entitlement fairness measure is not without its limitations. The framework is much more difficult to apply than existing fairness measures, as it requires a great deal more effort up front to identify morally relevant features to the problem domain and train separate prediction models for each features. This level of complexity is only further compounded by the need to maintain a rectification system for the users of the system, and the need to publish the features used in the decision process. This is a key current limitation of the system, as it is unlikely that many organizations will be willing to put the framework into practice due to the high cost of implementation and the need for a high level of transparency.

The entitlement fairness measure is also highly limited by the requirement of moral justification of the features used in the decision process. While an organization may provide justifications for the features they use in their decision process, it is highly unlikely that the justifications will be universally accepted or uncontroversial. The result is that by publishing the moral reasoning behind the features used to make the decision, the organization is likely to open itself up to a great deal of scrutiny and criticism from the public—more than it would have received if it had simply used a standard fairness and not published extensive justifications. While this is an unfortunate consequence of the framework, it is also a typical consequence of any change that promotes transparency and accountability and is a necessary step in the development of a more principled approach to algorithmic fairness.

5.3 Future Directions

While this work has made significant strides in understanding the relationship between algorithmic fairness and entitlement justice, several important questions remain unanswered.

First, the framework presented here is limited to the case of binary classification problems. While the framework can be extended to multi-class classification problems, the extension is not straightforward and requires significant additional work. The extension of the framework to multi-class classification problems is a key area for future research, as many real-world applications of algorithmic decision making are multi-class classification problems.

Second, the framework presented here is limited to the case of discrete features. The extension to continuously valued features opens new questions around how to define the moral relevance of features and how to measure the success of a particular feature in the decision process. Similarly, rectification becomes much more difficult in the case of continuously valued features, as the system implementer must now decide what threshold values represent a mistake by the algorithm that requires rectification.

Finally, though the framework presented here is a coherent way to understand entitlement justice in the setting of algorithmic decision-making, a large step in the development of the approach would be to develop a test system which implements the framework and demonstrates its utility in a real-world setting. Though many examples are provided in this paper to demonstrate the capabilities of the framework, a real-world system would validate and strengthen the arguments made in this paper.

References

- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias: There's software used across the country to predict future criminals. and it's biased against blacks. *ProPublica*, May 2016. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>. Accessed: 2025-01-12.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness and machine learning. In *Fairness and Machine Learning*. fairmlbook.org, 2017. Online book, available at <https://fairmlbook.org>.
- Joachim Baumann, Corinna Hertweck, Michele Loi, and Christoph Heitz. Distributive justice as the foundational premise of fair ml: Unification, extension, and interpretation of group fairness metrics, 2023. URL <https://arxiv.org/abs/2206.02897>.
- Reuben Binns. Fairness in machine learning: Lessons from political philosophy. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pages 149–159. PMLR, 2018. URL <https://proceedings.mlr.press/v81/binns18a.html>.
- Sam Corbett-Davies, Johann D. Gaebler, Hamed Nilforoshan, Ravi Shroff, and Sharad Goel. The measure and mismeasure of fairness, 2023. URL <https://arxiv.org/abs/1808.00023>.
- Crime and Justice Research Alliance. Black men have higher rates of recidivism despite lower risk factors, 2023. URL <https://crimeandjusticeresearchalliance.org/black-men-have-higher-rates-of-recidivism-despite-lower-risk-factors/>. Accessed: 2025-01-21.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference (ITCS)*, pages 214–226. ACM, 2012.
- Barbara Fried. Wilt chamberlain revisited: Nozick's "justice in transfer" and the problem of market-based distribution. *Philosophy & Public Affairs*, 24(3):226–245, 1995. doi: <https://doi.org/10.1111/j.1088-4963.1995.tb00030.x>.
- Barbara Fried. Left-libertarianism: A review essay. *Philosophy & Public Affairs*, 32(1):66–92, 2004. URL <https://law.stanford.edu/publications/left-libertarianism-a-review-essay/>.
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3315–3323. Curran Associates, Inc., 2016.
- M. C. Henberg. Nozick and rawls on historical versus end-state distribution. *The Southwestern Journal of Philosophy*, 8(2):77–84, 1977. ISSN 0038481X, 21541043. URL <http://www.jstor.org/stable/43155157>.

- Corinna Hertweck, Joachim Baumann, Michele Loi, Eleonora Viganò, and Christoph Heitz. A justice-based framework for the analysis of algorithmic fairness-utility trade-offs, 2023. URL <https://arxiv.org/abs/2206.02891>.
- Corinna Hertweck, Christoph Heitz, and Michele Loi. What’s distributive justice got to do with it? rethinking algorithmic fairness from the perspective of approximate justice, 2024. URL <https://arxiv.org/abs/2407.12488>.
- Lily Hu. Normative facts and causal structure. *The Journal of Philosophy*, forthcoming. To appear.
- Matthias Kuppler, Christoph Kern, Ruben L. Bach, and Frauke Kreuter. Distributive justice and fairness metrics in automated decision-making: How much overlap is there?, 2021. URL <https://arxiv.org/abs/2105.01441>.
- Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness, 2018. URL <https://arxiv.org/abs/1703.06856>.
- Yu Han Liu. Feature extraction and image recognition with convolutional neural networks. *Journal of Physics: Conference Series*, 1087(6):062032, sep 2018. doi: 10.1088/1742-6596/1087/6/062032. URL <https://dx.doi.org/10.1088/1742-6596/1087/6/062032>.
- Eric Mack. Self-ownership and the right of property. *The Monist*, 73(4):519–543, 1990. doi: 10.5840/monist19907343.
- Robert Nozick. *Anarchy, State, and Utopia*. Basic Books, New York, 1974. ISBN 978-0465097203. Proposes the entitlement theory of justice as a response to distributive justice theories like Rawls’s.
- John Rawls. *A Theory of Justice*. Harvard University Press, Cambridge, MA, revised edition edition, 1971. Original edition published in 1971, revised edition in 1999.
- Amartya Sen. Property and hunger. *Economics and Philosophy*, 4, 1988.
- Aaron Smith. Public attitudes toward computer algorithms, November 2018. URL <https://www.pewresearch.org/internet/2018/11/16/public-attitudes-toward-computer-algorithms/>. Accessed: 2025-01-12.
- Robert J. van der Veen and Philippe Van Parijs. Entitlement theories of justice: From nozick to roemer and beyond. *Economics and Philosophy*, 1(1):69–81, 1985. doi: 10.1017/S0266267100001899.