
Revisiting Fair-PAC Learning and the Axioms of Cardinal Welfare

Cyrus Cousins

University of Massachusetts Amherst

Abstract

Cardinal objectives like welfare and malfare have recently enjoyed increased attention in fair machine learning as computationally, statistically, and philosophically sound alternatives to constraint-based methods. Welfare summarizes utility over a population of g groups, whereas malfare measures overall disutility. Subject to standard axioms, both welfare and malfare functions are restricted to the family of w -weighted p -power-means, i.e., $M_p(\mathbf{u}; \mathbf{w}) \doteq \sqrt[p]{\sum_{i=1}^g w_i u_i^p}$, with $p \leq 1$ for welfare (utility \mathbf{u}), or $p \geq 1$ for malfare (disutility \mathbf{u}). This work revisits said axioms, finding that a weaker basis is sufficient to show the same, and furthermore that strengthening these axioms partition the welfare half of the spectrum (i.e., $p \leq 1$) into a few cases by further limiting p . It is known that $p \geq 1$ power means (i.e., malfare functions) are Lipschitz continuous, and thus statistically easy to estimate or learn (i.e., each \mathbf{u}_i can be approximated with a sample estimate). We show that all power means are at least locally Hölder continuous, i.e., $|M(\mathbf{u}; \mathbf{w}) - M(\mathbf{u}'; \mathbf{w})| \leq \lambda \|\mathbf{u} - \mathbf{u}'\|^\alpha$ for some constants $\lambda > 0$, $\alpha \in (0, 1]$, and some norm $\|\cdot\|$. Furthermore, λ and $\frac{1}{\alpha}$ are bounded except as $p \rightarrow 0$ or $\min_i w_i \rightarrow 0$, and via this analysis we bound the sample complexity of estimating or optimizing welfare functions. This yields a novel concept of fair-PAC learning, with dependence on the quantities $\frac{1}{|p|}$ and/or $\frac{1}{w_{\min}}$ (which measure closeness to the challenging $p = 0$ case and inverse minimum group weight, respectively), wherein fair welfare functions are only polynomially harder to optimize than fair malfare functions, except when $p \approx 0$ or $\min_i w_i \approx 0$, which may be exponentially harder. These challenging cases may be avoided by assuming the appropriate strengthened axioms. In all cases, we show that if a bounded quantity is learnable with finite sample complexity, then so too is the welfare of said quantity. This takes estimat-

ing and learning welfare objectives to near-parity with malfare objectives, as although sample complexity may be larger, all such objectives are *uniformly learnable*.

Keywords

Fair Machine Learning \diamond Algorithmic Fairness
Fair PAC Learning \diamond Statistical Learning Theory
Social Planner’s Problem \diamond Cardinal Welfare Theory
Weighted Power Means \diamond Hölder Continuity Analysis

1 Introduction

The recent resurgence of cardinal welfare and malfare based methods in group-based fairness settings has led to increased attention as to how to *objectively quantify* fairness. Welfare summarizes utility across a population, and is thus suitable for fair utility-maximization tasks (e.g., bandit problems, reward-based reinforcement learning [Siddique et al., 2020, Cousins et al., 2022a], and recommender systems, as well as utility-based finance and economic settings), whereas malfare measures overall disutility, and is thus directly applicable to fair loss minimization tasks (arguably most machine learning tasks). The promise of statistical and computational efficiency differentiates such approaches from constraint-based fairness methods [Dwork et al., 2012, Zemel et al., 2013], which often yield hard non-convex optimization problems, requiring convex relaxations (potentially sacrificing fairness), as well as inducing statistical quandaries in estimating whether such fairness constraints generalize beyond the training set [Yona and Rothblum, 2018, Thomas et al., 2019]. The axiomatic justification for cardinal welfare and malfare functions also gives them a sense of objectivity, whereas fairness constraints are often intuitively motivated, and at times mutually incompatible [Kleinberg et al., 2017, Friedler et al., 2021]. We find a basis of cardinal welfare axioms that is weaker than the standard basis, and we then propose stronger axioms to further specify such functions, and explore the resulting classes of fair learnability.

It is now well-understood that many problems of unfairness in machine learning are caused at least in part by a lack of training data for relevant groups [Chen et al., 2018a, Mehrabi et al., 2021], e.g., an overrepresentation of images

of whites and males in training data for facial recognition systems [Buolamwini and Gebru, 2018, Cook et al., 2019, Cavazos et al., 2020]. While addressing such data biases is certainly a step in the right direction, we argue that *the objective itself* must be selected to appropriately compromise between the wants and needs of various groups, which may not mutually align.¹ Collecting sample sizes proportional to population frequencies of each group and minimizing empirical risk (average loss) or maximizing empirical utility seems reasonable, but this is implicitly a utilitarian perspective, as this strategy is equivalent to maximizing *empirical weighted utilitarian welfare or malfare* (as a proxy for their expected values), with the weight of each group i proportional to the number of training points from group i . Furthermore, even if one accepts the tenets of utilitarianism (a perfectly valid choice, but one that should be made *consciously*, rather than *by default*), the result is still not necessarily fair, as we risk overfitting to the proportional, but overall small, amount of data collected for smaller groups. Our axiomatic approach to fairness rigorously treats the issue of selecting and optimizing fair objectives, in particular issues of sample complexity and overfitting to fairness, whether one’s chosen fairness objective be strictly utilitarian or of a more prioritarian nature. In particular, any given fair objective effectively specifies how *tradeoffs* must be made between groups when *disagreements* arise between groups during training (discussed further in section 4.3), and the utilitarian welfare or malfare is just one of many axiomatically justifiable choices.

Section 3.1 shows that some of the cardinal welfare or malfare axioms of Cousins [2021] can be relaxed or reorganized to form a *piecewise-weaker equivalent basis* (i.e., each of our axioms is no stronger than an existing standard axiom, yet their collective effect is equivalent) that is more elegant and more concise, thus our axiomatization is a more convincing premise upon which to develop a theory of fair learning. Subject to these axioms, the *w-weighted p-power-mean family* arises as the only axiomatically justified class of fair aggregator functions, however the parameter space of this class is quite large, thus the theory does not uniquely specify an “ideal fairness concept.” Many have argued that exact human-desirable fairness concepts cannot be fully specified without unjustifiable assumptions, and that variation in feasible aggregator function concepts reflects variation in human morality and social values [Awad et al., 2018, Schneider and Leland, 2021]. We do not reject this claim, however we do show in section 3.3 that additional axioms can further restrict the family of malfare or welfare functions, although such axioms may be less universal than the standard basis of cardinal welfare or malfare axioms.

¹For example, a learned spellchecker may encounter both UK English and US English, and the decision to correct, e.g., either “color” or “colour” is a tradeoff that explicitly harms writers of one dialect while benefiting writers of the other.

In particular, we propose a stronger variant of *strict monotonicity*, which precludes both egalitarian and utilitarian welfare, as well as a relaxation of the Pigou-Dalton [Pigou, 1912, Dalton, 1920] transfer principle, an axiom which characterizes fairness by promoting equitable redistribution of (dis)utility, as well as a strengthening of it, which specifically characterizes utilitarianism as *neutral*, rather than equitable or fair. We also introduce the *zero barrier* and *finite ceiling* axioms, both of which promote a weak form of egalitarianism (prioritarianism) by restricting the behavior of welfare functions under extreme inequality.

While welfare maximization and malfare minimization appear to be two sides of the same coin, salient mathematical differences arise. We find in section 3.4 that, unlike malfare functions, welfare functions are not always Lipschitz continuous (though they are at least Hölder continuous), and not uniformly fair-PAC (FPAC) learnable in the sense of Cousins [2021]. However, section 4.3 shows that under a slightly more lenient definition of FPAC learnability, in which sample complexity (i.e., the sufficient sample size to approximately optimize an objective over some class \mathcal{H} , for any data distributions, with high probability) is allowed to depend on the welfare function through the *minimum group weight* reciprocal $\frac{1}{w_{\min}}$ and/or the quantity $\frac{1}{|p|}$ (which quantifies how close the welfare function is to the challenging $p = 0$ Nash social welfare case), then if \mathcal{H} is uniformly convergent with polynomial sample complexity, \mathcal{H} is also FPAC learnable with polynomial sample complexity.

We then split the power-mean spectrum into regions, the boundaries of which are defined by our extended axioms, and show that each is in some sense FPAC learnable. This work culminates in theorem 4.4, which shows that if a bounded utility maximization task is PAC-learnable with polynomial sample complexity, then the corresponding *welfare maximization* task is FPAC-learnable, with possible exponential dependence on $\frac{1}{w_{\min}}$ or $\frac{1}{|p|}$, and otherwise fully-polynomial sample complexity. We close section 4 by showing that our sample-complexity bounds can be incorporated into progressive-sampling routines, which adapt their sample consumption to the difficulty of the task at hand, with only *logarithmic multiplicative overhead*.

Our study of fair-PAC learnability is based on the Lipschitz and Hölder continuity of power means, which only coarsely describe their behavior. We thus bound *worst-case* sample complexity, as such analysis is necessarily focused on behavior under extreme inequality of per-group sentiment, where small changes to (dis)utility values are most impactful (as per the *Pigou-Dalton transfer principle*). Under equality, power-mean functions behave roughly linearly, and may thus be much easier to estimate or optimize. We argue that our worst-case bounds for sample complexity are theoretically interesting and practically significant, as they lend confidence to the proposition that a machine learning system accurately and fairly learns an objective. In section 4.4,

we argue that the actual inefficiency (excess sample consumption) stemming from our analysis is not necessarily dependent on the ratio of the best-case to worst-case sample complexity, as efficient progressive sampling algorithms (which draw successively larger samples, and terminate when the objective is estimated or optimized to within some error tolerance ϵ) generally depend primarily on the *true sufficient sample size*, with an excess factor of the *double-logarithm* of this ratio. Consequently, a polynomial gap between best-case mean estimation and worst case welfare-estimation bounds translates to only *double-logarithmic overhead*, although the exponential gap (dependence on $\frac{1}{w_{\min}}$ and/or $\frac{1}{|p|}$ for $p \approx 0$) may translate to a *logarithmic* increase in sample complexity (relative to methods that are given knowledge of the optimal sample size *a priori*), even with efficient progressive sampling methods.

For the sake of brevity, all proofs are meticulously derived in appendix A. In summary, the primary contributions of this paper are as follows.

- 1) In section 3.1 we derive a simplified axiomatic basis for cardinal welfare and malfare functions.
- 2) We extend the core axiomatic basis in section 3.3 with optional axioms, which lead to convenient computational and statistical properties, while enjoying intuitive real-world interpretations, thus guiding welfare function selection.
- 3) In section 3.4 we analyze the Lipschitz and Hölder continuity of power means. We argue that local behavior is crucial to algorithm analysis, particularly statistical behavior and response to small input perturbations.
- 4) Section 4 explores the impact of additional axioms on FPAC learnability. In particular, if we allow sample complexity to depend on $\frac{1}{w_{\min}}$ and $\frac{1}{|p|}$, additional axioms naturally split the power-mean spectrum into several regions, and show that each is in some sense FPAC learnable.

2 Related Work

In cardinal fairness learning tasks, we generally receive training data or feedback from *multiple groups*, which represents the needs or wants of *each group*, and we seek to maximize welfare or minimize malfare to *fairly compromise* among groups. Objective choice is a crucial modelling decision, as it mathematically encodes the values of the social planner [Sen, 1977, Roberts, 1980, Moulin, 2004]. Due to variation in human values and value systems, we can't uniquely characterize fairness with mathematics alone, however analysis does help to elucidate the limitations and properties of cardinal objectives. Axiomatic reasoning and analysis from the econometrics literature informs us as to the limitations and properties of cardinal objectives [Pigou, 1912, Dalton, 1920, Debreu, 1959, Gorman, 1968]. The moral philosophy literature also provides insight into social objectives, from classical *utilitarian theory* [Bentham, 1789, Mill, 1863], in which all parties are treated as equals, to *prioritarianism*, [Parfit, 1997, Arneson, 2000], where

the less-fortunate are given more weight, to *egalitarian* or *Rawlsian theory* [Rawls, 1971, 2001], which considers the least-fortunate before all others. We also draw from various critiques of utilitarianism [Nozick, 1974, Hurka, 1982] in our analysis.

Prior fair learning work in computer science primarily handles the malfare case. In particular, group-DRO (distributionally robust optimization) methods minimize worst-case (over groups; i.e., egalitarian) risk [Hu et al., 2018, Oren et al., 2019, Sagawa et al., 2019, Dong and Cousins, 2022], which is also known as minimax-fair learning [Diana et al., 2021, Shekhar et al., 2021, Abernethy et al., 2022] and by other names [Martinez et al., 2020, Lahoti et al., 2020, Cortes et al., 2020, Shekhar et al., 2021], and fair-PAC learning generalizes this idea by optimizing arbitrary power-mean malfare objectives [Cousins, 2021, 2022, Cousins et al., 2022a], which derive from an axiomatic welfare theory perspective. Similar algorithmic and statistical concerns arise in the *multi-group agnostic PAC learning* setting [Blum and Lykouris, 2020, Rothblum and Yona, 2021, Cousins, 2022], wherein the goal is to minimize *regret* (over groups) w.r.t. needing to compromise on a *shared model*, rather than each group selecting their own preferred model. Some authors, e.g., Hu and Chen [2020], do discuss direct welfare optimization, however they do not treat the resulting statistical questions or bound generalization error or sample complexity, and thus the issues we identify with the statistical difficulty of welfare optimization are not addressed.

As for welfare-theoretic approaches to fair learning, Heidari et al. [2018] employ axiomatic cardinal welfare theory to introduce fairness constraints for machine learning tasks, and Cousins [2021, 2022] generalizes the axioms of cardinal welfare to allow for per-group weight values, explores computational and statistical learnability, and bounds sample complexity for the malfare ($p \geq 1$) case. Thomas et al. [2019] also introduce a concept of fair statistical learnability, termed the *Seldonian learner*, which encapsulates both constraint-based and cardinal objective-based fair learning settings, however, this framework is so general that it is difficult to establish classes of learnability. In contrast, by adopting an *axiomatic* cardinal-welfare centric viewpoint, we operate over *specific classes* of welfare or malfare functions, which allows us to broadly analyze the sample complexity of various learning problems.

Outside the sphere of fair learning, similar questions and tradeoffs arise in *fair allocation*. In particular, welfare-based fair allocators must explicitly decide how to allocate a pool of limited items (resources) to a set of agents given their valuation (utility) functions over sets of items. In a sense, a machine learning model is also a resource allocator, where the resource is the sentiment achieved by each group as a result of decisions made by the model. Fair allocation problems generally sidestep the issues of generalization error and PAC learnability by assuming known utility, but the

core problems of selecting a fair objective and optimizing a given fair objective remain. Power-means have recently seen use as fair objectives in various fair allocation settings. Barman et al. [2020] show efficient approximation algorithms for unweighted power-means under subadditive valuations, in more restricted classes of submodular valuations, Viswanathan and Zick [2023] give efficient algorithms to maximize weighted power-means, and in broader valuation classes, Cousins et al. [2023b,c] give efficient algorithms to maximize unweighted power-means.

3 On Axiomatizations of Cardinal Fairness Objectives

Herein we define *aggregator functions*, which summarize overall sentiment (either utility or disutility), expressed as a *sentiment vector* $\mathbf{u} \in \mathbb{R}_{0+}^g$ over a population of g groups weighted by *weights vector* $\mathbf{w} \in \Delta_g$, where Δ_g is the *nondegenerate probability simplex* over g groups, i.e., $\mathbf{w} \in (0, 1)^g$ and $\|\mathbf{w}\|_1 = 1$. We use $M(\mathbf{u}; \mathbf{w})$ to denote generic aggregator functions, $W(\mathbf{u}; \mathbf{w})$ for welfare functions (positively-connoted sentiment), and $\Lambda(\mathbf{u}; \mathbf{w})$ for malfare functions (negatively-connoted sentiment). At times, we are also interested in *unweighted aggregator functions*, which take the form $M(\mathbf{u})$, but in all cases, these will be equivalent to some weighted aggregator function under a uniform weighting, i.e. $M(\mathbf{u}) = M'_p(\mathbf{u}; \langle \frac{1}{g}, \dots, \frac{1}{g} \rangle)$ for some weighted aggregator function $M'(\cdot; \cdot)$. We present in section 3.1 an axiomatic basis of desirable properties that is weaker (i.e., makes fewer assumptions) than the standard basis, and in section 3.3, we give additional axioms and strengthened variants of our basis axioms, which further constrain the space of fair aggregator functions. In section 3.2, we explore the resulting classes of aggregator functions, and in section 3.4 we further explore various continuity properties and applications of these classes of aggregator functions.

To concisely denote sentiment, we use functional and vector notation interchangeably, e.g., the *logarithmic utility transformation* of concave utility theory, applied to some \mathbf{u} , could be written either as $i \mapsto \ln(1 + \mathbf{u}_i)$ or as $\langle \ln(1 + \mathbf{u}_1), \ln(1 + \mathbf{u}_2), \dots, \ln(1 + \mathbf{u}_g) \rangle$. Furthermore, indicator functions, i.e., $\mathbb{1}_a(b)$ is 1 if $b = a$ or $b \in a$, and 0 otherwise, can also be interpreted as *indicator vectors*, where $\mathbb{1}_i = j \mapsto \mathbb{1}_i(j)$ is the i th standard basis vector. This and all other potentially confusing notation pertaining to aggregator functions is summarized in table 1.

Although this work explicitly considers only group-level fairness and decision making, this is predominantly for reasons of *statistical learnability*. Section 4 assumes we can draw many samples *for each group* to estimate their sentiment while learning across a class \mathcal{H} of possible models, which is reasonable when samples correspond to the *individuals* that comprise the group (thus their experiences

are averaged to determine the group’s sentiment), but it is often unrealistic to assume a learner has access to a large amount of individual-level data. Philosophically, operating at the group level also allows us to consider only expected outcomes over populations, which circumvents the need to reason about randomness and probabilities of individual events. These issues aside (e.g., if the outcomes of any feasible decision were known deterministically), there is no reason our framework could not be applied to individual-level fairness; we would need only substitute the word “individual” for “group,” or equivalently, define the population of groups as *singletons* each consisting of exactly one individual. Note also that the *weighted aggregator functions* are implicitly motivated by differences in the relative sizes of groups, whereas at an individual level, it is generally desirable to treat all individuals equally, and thus instead operate with *unweighted aggregator functions*, but as we shall see, the axiomatic justification behind unweighted and weighted aggregator functions are quite similar.

It is important to acknowledge the implicit assumptions of this setup. Generally, sentiment vectors arise from *real-world grounded situations* that impact each group. By construction, aggregator functions are cardinal, and therefore impose a *preference ordering* over sentiment vectors (and the grounded situations that give rise to said sentiment vectors). The social planner seeks to select a grounded situation (i.e., learn a model) to optimize this preference ordering, and is thus impartial towards the grounded situation, except insofar as it impacts each group’s sentiment value. This factorization simplifies the question of how to define fairness by avoiding *objective characterization*, and instead defining an *intersubjective concept* (shared welfare or malfare function) based on the *subjective experience* (sentiment value) of each group.

This is inherently an “intersubjective consequentialist” perspective on ethics, valuing the impact of decisions as perceived by the groups affected by them and agreed upon systems of compromise, and it thus draws inspiration from the philosophy of *altruistic hedonism* (i.e., it seeks to maximize the pleasure and minimize the pain of everyone). While we do not explicitly consider other moral and ethical systems, it is not difficult to modify our framework to accommodate many such philosophies. In particular, more paternalistic or prescriptive ethical systems, such as deontological or virtue ethical systems, can be modeled by modifying the definition of sentiment to reflect the feelings or opinions of the social planner (acting as a “morality arbitrator”), rather than those of the groups impacted by decisions (for example, if *truth* is a moral virtue, then a classifier should strive for accuracy, even if the decisions cause harm and leave each group in a worse position than would a less accurate model). Similarly, more solipsistic moral frameworks, e.g., *egotistical hedonism* (wherein one’s concept of morality aligns with that which brings them pleasure and avoids bringing them pain),

Object	Definition or Space	Description
Δ_g	$\doteq \{\mathbf{w} \in (0, 1)^g \mid \ \mathbf{w}\ _1 = 1\}$	Nondegenerate probability simplex over g atoms
$i \mapsto f_i(\mathbf{u}_i)$	$\doteq \langle f_1(\mathbf{u}_1), \dots, f_g(\mathbf{u}_g) \rangle$	Functional vector notation
$\mathbb{1}_a(b)$	$\doteq 1$ if $b \in a$, 0 otherwise	Indicator vector (or function)
$\mathbb{1}_i$	$\doteq j \mapsto \mathbb{1}_{\{i\}}(j)$	One-hot indicator vector
\mathcal{G}	Usually $\{1, \dots, g\}$	Group identity space (e.g., race, gender, or language categories)
g	$\doteq \mathcal{G} $	Group count or cardinality (usually finite)
r	$\in [0, \infty]$	Maximum possible sentiment value (range)
\mathbf{u}	$\in [0, r]^g$	Per-group sentiment (utility or disutility), i.e., \mathbf{u}_i pertains to group i
\mathbf{w}	$\in \Delta_g$	Probability weighting over groups (probability simplex)
$M(\mathbf{u}), W(\mathbf{u}), \Lambda(\mathbf{u})$	$\in [0, r]^g \rightarrow \mathbb{R}_{0+}$	Unweighted aggregator, welfare, or malfare function
$M(\mathbf{u}; \mathbf{w}), W(\mathbf{u}; \mathbf{w}), \Lambda(\mathbf{u}; \mathbf{w})$	$\in ([0, r]^g \times \Delta_g) \rightarrow \mathbb{R}_{0+}$	Weighted aggregator, welfare, or malfare function
$M_p(\mathbf{u}; \mathbf{w})$	See definition 3.4	Weighted power-mean aggregator function

Table 1: Functions, spaces, and common variables for aggregator functions.

can be modelled by changing the social planner from an unbiased abstract outsider not impacted by their decisions into a member of the population being impacted that selfishly makes decisions influenced by their own sentiment.

3.1 A Fundamental Basis of Cardinal Welfare Axioms

We now present a reduced basis of cardinal welfare axioms, which are componentwise-weaker than the standard basis, yet we find that the same Debreu–Gorman [Debreu, 1959, Gorman, 1968] type theorems and Pigou–Dalton [Pigou, 1912, Dalton, 1920] characterizations of fairness still hold.

As a prelude to the dense axiomatic mathematics that will soon follow, we briefly discuss the goals of such reasoning and analysis. It is tempting to think that the value in these axioms derives from their consequent properties, but this line of reasoning is dangerous, as it may lead to “baxiomatic” quasireligious thinking, i.e., the pressure to accept the axioms comes from the convenient or wishful thinking that, once we do so, we can reason over a convenient space, and thus normatively refrain from considering any situation outside their purview. To remain grounded, we thus argue that *the axioms themselves* must be *inherently reasonable*. This “I know it when I see it” criterion, however, does lead to a metric for of the efficacy of an axiomatization: namely, is it *simple enough* to be *easily understood* by rational thinkers, and are they likely to agree with *all parts* of it (emphasis on all, as “most” is not good enough: a single faulty premise can destroy an entire chain of reasoning)?

For these reasons, in settings of axiomatic reasoning such as this, elegance and simplicity are paramount, as they lead to interpretable and uncontroversial axioms. Formal logic plays a role as well, as showing that $A \implies B$ and $A \neq B$, or in other words, B is strictly weaker than A , tells us that B is less likely to be rejected by any rational thinker than A . We thus seek to ground our reasoning on a *minimal basis* of axioms, each of which cannot be weakened individually without changing their collective effect. To emphasize the axiomatization itself, rather than its consequent properties,

we present our axioms in isolation, before stating any of their resultant properties. We encourage the reader to evaluate them in and of themselves, though occasionally we do mention specific consequences of various assumptions, in particular to illustrate concrete differences between weaker or stronger variants of some axioms.

Classical econometric theory primarily describes the *unweighted case*, wherein an aggregator function $M(\mathbf{u})$ aggregates sentiment \mathbf{u} over an *unweighted finite discrete* population \mathcal{G} , i.e., $\mathcal{G} = \{1, \dots, g\}$. We generalize this setting to the *weighted discrete* case, wherein an aggregator function $M(\mathbf{u}; \mathbf{w})$ operates on a *w-weighted discrete* (possibly countably infinite) population \mathcal{G} , i.e., $\mathcal{G} = \{1, \dots, g\}$, for $g \in \mathbb{Z}_+ \cup \{\infty\}$. We modify the axiomatization of Cousins [2021], showing that a reduced axiomatic basis is equivalent.

Axiomatization 3.1 (Weighted Aggregator Axioms). We define axioms for aggregator function $M(\mathbf{u}; \mathbf{w})$ below. For each item, assume (if necessary) that the axiom applies for all $\mathbf{u}, \mathbf{u}' \in \mathbb{R}_{0+}^g$ scalars $\alpha, \varepsilon \in \mathbb{R}_+$, indices $i, j \in \mathcal{G}$, and discrete probability measures $\mathbf{w} \in \Delta_g$ over \mathcal{G} .

- 1) *Strict Monotonicity* (SM): Suppose $\mathbf{u} \succ 0$, i.e., each $u_j > 0$. Then $M(\mathbf{u}; \mathbf{w}) < M(\mathbf{u} + \varepsilon \mathbb{1}_i; \mathbf{w})$.
- 2) *Weighted Symmetry* (WS): For all permutations π over \mathcal{G} , it holds that $M(\mathbf{u}; \mathbf{w}) = M(\pi(\mathbf{u}); \pi(\mathbf{w}))$.
- 3) *Weighted Decomposability* (WD): Suppose $\alpha \in (0, 1)$. Then $M(\mathbf{u}; \mathbf{w}) = M(\langle \alpha \mathbf{u}_1, \mathbf{u}_1, \mathbf{u}_2, \dots \rangle; \langle \alpha \mathbf{w}_1, (1 - \alpha) \mathbf{w}_1, \mathbf{w}_2, \dots \rangle)$.
- 4) *Continuity*: $M(\mathbf{u}; \mathbf{w})$ is a continuous function in \mathbf{u} .
- 5) *Independence of Unconcerned Agents* (IUA): If $\mathbf{u}_i = \mathbf{u}'_i$, then $M(\mathbf{u}; \mathbf{w}) \leq M(\mathbf{u}'; \mathbf{w}) \iff M(\mathbf{u} + \varepsilon \mathbb{1}_i; \mathbf{w}) \leq M(\mathbf{u}' + \varepsilon \mathbb{1}_i; \mathbf{w})$.
- 6) *Multiplicative Linearity*: $M(\alpha \mathbf{u}; \mathbf{w}) = \alpha M(\mathbf{u}; \mathbf{w})$.
- 7) *Unit Scale*: $M(\mathbf{1}; \mathbf{w}) = M(i \mapsto 1; \mathbf{w}) = 1$.
- 8) *Weak Transfer Principle* (WTP): Let $i \doteq \operatorname{argmin}_i \mathbf{u}_i$ & $j \doteq \operatorname{argmax}_j \mathbf{u}_j$. If $\mathbf{u}_i \neq \mathbf{u}_j$, then *there exists* some $\varepsilon > 0$ s.t. $\mathbf{u}_i + \mathbf{w}_j \varepsilon < \mathbf{u}_j - \mathbf{w}_i \varepsilon$, and $W(\mathbf{u} + \varepsilon \mathbf{w}_j \mathbb{1}_i - \varepsilon \mathbf{w}_i \mathbb{1}_j; \mathbf{w}) \geq W(\mathbf{u}; \mathbf{w})$ for welfare or $\Lambda(\mathbf{u} + \varepsilon \mathbf{w}_j \mathbb{1}_i - \varepsilon \mathbf{w}_i \mathbb{1}_j; \mathbf{w}) \leq \Lambda(\mathbf{u}; \mathbf{w})$ for malfare.

Axioms 1–8 are generally assumed in this work, but we

present several alternatives below, to which we compare.

9) *Weighted Additivity* (WA): Suppose $g' \in \mathbb{Z}_+ \cup \{\infty\}$, $\mathbf{u}' \in \mathbb{R}_{0+}^{g'}$, and weights vector $\mathbf{w}' \in \Delta_{g'}$ such that for all $u \in \mathbb{R}_{0+}$, it holds that $\sum_{i \in \mathcal{G}} w_i \mathbb{1}_u(\mathbf{u}_i) = \sum_{i \in \mathcal{G}'} w'_i \mathbb{1}_u(\mathbf{u}'_i)$. Then $M(\mathbf{u}; \mathbf{w}) = M(\mathbf{u}'; \mathbf{w}')$.

10) *Independence of Common Scale* (ICS): $M(\mathbf{u}; \mathbf{w}) \leq M(\mathbf{u}'; \mathbf{w}) \implies M(\alpha \mathbf{u}; \mathbf{w}) \leq M(\alpha \mathbf{u}'; \mathbf{w})$.

11) *Pigou-Dalton Transfer Principle* (PDTP)²: If $\mathbf{u}_i + \mathbf{w}_j \varepsilon \leq \mathbf{u}_j - \mathbf{w}_i \varepsilon$, then $W(\mathbf{u} + \mathbf{w}_j \varepsilon \mathbb{1}_i - \mathbf{w}_i \varepsilon \mathbb{1}_j; \mathbf{w}) \geq W(\mathbf{u}; \mathbf{w})$ for welfare or $\Lambda(\mathbf{u} + \mathbf{w}_j \varepsilon \mathbb{1}_i - \mathbf{w}_i \varepsilon \mathbb{1}_j; \mathbf{w}) \leq \Lambda(\mathbf{u}; \mathbf{w})$ for mal-
fare.

We now pause to discuss the rationale behind each axiom. Axioms 1–5 & 11 generalize the standard basis of *axioms of cardinal welfare* to *weighted discrete populations*, and together they imply that any aggregator function can be expressed as $M(\mathbf{u}; \mathbf{w}) \doteq F(\sum_{i=1}^g w_i f(\mathbf{u}_i))$ for strictly monotonically increasing functions f, F . A similar basis of axioms pertains to the class of *unweighted aggregator functions*, i.e., those of the form $M(\mathbf{u})$, where implicitly all weights are equal (which translate to the form $M(\mathbf{u}; (\frac{1}{g}, \dots, \frac{1}{g}))$ in the weighted nomenclature). In particular, unweighted variants of axioms 1 & 4–7 simply drop all weights terms, and are otherwise identical to their weighted counterparts. Unweighted variants of the remaining axioms are slightly more involved, and are discussed in turn below.

On Weighted Additivity In the unweighted case, it is standard to define symmetry as simply $M(\mathbf{u}) = M(\pi(\mathbf{u}))$ for all permutations π over \mathcal{G} . With weights, *weighted symmetry* (axiom 2), i.e., $M(\mathbf{u}; \mathbf{w}) = M(\pi(\mathbf{u}); \pi(\mathbf{w}))$, only requires *equal treatment* given *equal weights*. *Weighted decomposability* (axiom 3) then codifies the relative impact of weights by requiring that a group can be decomposed into two groups of equal sentiment and total weight without changing the aggregate. Prior work [Cousins, 2021, 2022] assumes *weighted additivity* (axiom 9) directly, but we argue that this axiom seems rather contrived and unintuitive, whereas axioms 2 & 3 are so natural that it would be perverse not to assume them. We now show that, despite their vastly simpler form, together axioms 2 & 3 equate to axiom 9.

Lemma 3.2 (Equivalence of Weighted Axioms). Consider some aggregator function $M(\cdot; \mathbf{w})$. It always holds that WS (axiom 2) \wedge WD (axiom 3) \Leftrightarrow WA (axiom 9).

On Units, Scale, and Canonical Forms Axiom 6 (*multiplicative linearity*) is a natural and useful property, and ensures that *dimensional analysis* on aggregator functions is possible; in particular, *units* of aggregator functions match

²Cousins [2021, 2022] adopts a seemingly more complicated variant of this axiom, but it is equivalent for countable populations by repeated application. In particular, they take the following: suppose $\mu = \mathbb{E}_w[\mathbf{u}] = \mathbb{E}_w[\mathbf{u}']$, and $\forall i \in \mathcal{G} : |\mu - \mathbf{u}'_i| \leq |\mu - \mathbf{u}_i|$. Then $W(\mathbf{u}'; \mathbf{w}) \geq W(\mathbf{u}; \mathbf{w})$ for welfare, or $W(\mathbf{u}'; \mathbf{w}) \leq W(\mathbf{u}; \mathbf{w})$ for mal-
fare.

those of sentiment values. Axiom 6 is also known in the *Constant Elasticity of Substitution* (CES) literature [Arrow et al., 1961, McFadden, 1963] as *homogeneity of degree 1*. Note that axiom 6 implies axiom 10, and is thus a simple strengthening of a more basic traditional cardinal welfare axiom. Axiom 7 (*unit scale*) furthers this theme, as it ensures that not only do *units* of aggregates match those of \mathbf{u} , but *scale* does as well, thus axiom 7 accords with *average utilitarianism* Hurka [1982], rather than *sum utilitarianism*, as we do not depend on the *size* of \mathcal{G} . Cousins [2021] shows that these axioms lead to the *power-mean* characterization of aggregator functions. Weakening axioms 6 & 7 to just axiom 10, the Debreu-Gorman theorem still implies a *monotonic transform* of the power-mean, so axioms 6 & 7 don't impact *comparisons between* aggregator values, but merely specify a convenient and elegant *canonical form* for their cardinal values.

On Equitable Redistribution and Transfer Principles

The *Pigou-Dalton transfer principle* (PDTP, axiom 11) is also standard in cardinal welfare theory. It essentially states that transferring (dis)utility between two groups is *not harmful*, up to the point where the two groups have equal (dis)utility, thus it incentivizes *equitable redistribution of "wealth."* This codifies the intuition that redistributing (dis)utility towards equitability is not harmful to society.

One could argue that, while a general trend towards equality may be good, this characterization of radical equality is too strong. The *weak transfer principle* (WTP, axiom 8) is less impeachable in this regard, as it weakens the quantifier over transfer magnitude from *universal* to *existential*, i.e., it states only that transferring *some nonzero amount* of (dis)utility between the (dis)utility maximizing and minimizing groups is not harmful, making no claim about the remaining groups, or the magnitude of the transfer. We now show that, subject to the standard Debreu-Gorman axioms, the WTP and PDTP are equivalent, thus the radical equality characterization of the PDTP is not necessary in this context.

Lemma 3.3 (Transfer Principle Equivalencies). Consider some aggregator function $M(\cdot; \mathbf{w})$. The following relate properties (axioms) that $M(\cdot; \mathbf{w})$ obeys.

- 1) PDTP (11) \implies WTP (8); &
- 2) Suppose axioms 1–7. Then WTP (8) \implies PDTP (11).

Note cautiously that both the PDTP and WTP are careful to claim that equitable transfers are *not harmful* to society, rather than *beneficial*. Section 3.3 presents strong variants of these axioms, which require *strict benefit* to equitable transfers. To obtain unweighted variants of the transfer principle axioms, we again simply substitute in uniform weights, but here this results in some elegant simplifications. In particular, the WTP and PDTP assume (either existentially or universally) some $\varepsilon > 0$ such that $\mathbf{u}_i + \mathbf{w}_j \varepsilon$ and $\mathbf{u}_j - \mathbf{w}_i \varepsilon$ obey some relationship. Observe that all weighting terms essentially cancel, and we now recover the guarantee for the

change-of-variables $\varepsilon' = \frac{1}{g}\varepsilon$, since $w_j\varepsilon = w_i\varepsilon = \varepsilon'$.

3.2 The Power Mean

We now define the class of weighted power-mean aggregator functions, and show that our aggregator function axioms are uniquely satisfied by this class.

Definition 3.4 (Power-Mean Welfare and Malfare). Suppose $p \in \mathbb{R} \cup \pm\infty$. The *weighted power-mean*, given *sentiment vector* $\mathbf{u} \in \mathbb{R}_{0+}^g$ and *weights vector* $\mathbf{w} \in \Delta_g$, is defined as

$$M_p(\mathbf{u}; \mathbf{w}) \doteq \lim_{\rho \rightarrow p} \lim_{\varepsilon \rightarrow 0^+} \sqrt[\rho]{\sum_{i=1}^g w_i (\mathbf{u}_i + \varepsilon)^\rho}. \quad (1)$$

Note that in most cases the limits can safely be ignored. The inner $\lim_{\varepsilon \rightarrow 0^+}$ simply avoids indeterminate forms for $p \leq 0$ while some $\mathbf{u}_i = 0$ while preserving continuity. The outer $\lim_{\rho \rightarrow p}$ resolves to the *weighted geometric mean* for $p = 0$, i.e., $M_0(\mathbf{u}; \mathbf{w}) = \lim_{\varepsilon \rightarrow 0^+} \prod_{i=1}^g (\mathbf{u}_i + \varepsilon)^{w_i}$, generally termed *Nash social welfare* in this context. Similarly, $\lim_{\rho \rightarrow p}$ resolves to the unweighted *maximum* and *minimum* operators for $p = \pm\infty$, also known as *egalitarian malfare* or *welfare*, respectively. We omit the weighting \mathbf{w} to denote unweighted power-means, which are equivalent to weighted power-means under a uniform weighting, i.e. $M_p(\mathbf{u}) = M_p(\mathbf{u}; \langle \frac{1}{g}, \dots, \frac{1}{g} \rangle)$.

We now characterize the class of fair aggregator functions in a result similar to theorem 2.4 of Cousins [2021], albeit under our reduced axiomatic basis.

Theorem 3.5 (Aggregator Function Properties). Suppose aggregator function $M(\mathbf{u}; \mathbf{w})$, and assume arbitrary sentiment vector $\mathbf{u} \in \mathbb{R}_{0+}^g$ and weights vector $\mathbf{w} \in \Delta_g$. The following then hold.

- 1) *Power-Mean Factorization*: Axioms 1–7 imply there exists some $p \in \mathbb{R}$ such that $M(\mathbf{u}; \mathbf{w}) = M_p(\mathbf{u}; \mathbf{w})$.
- 2) *Fair Welfare and Malfare*: Axioms 1–8 imply $p \in (-\infty, 1]$ for welfare and $p \in [1, \infty)$ for malfare.

3.3 Extended and Contextual Axioms

We now present new axioms and stronger variants of the axioms thus far stated. These generally extend the themes and justifications of weaker axioms, and thus require a larger concession to accept, but they have greater descriptive power and reduce the space of admissible fair aggregator functions.

Axiomatization 3.6 (Strong Axioms). Suppose as in axiomatization 3.1. The following two axioms strengthen various components of axiomatization 3.1.

- 12) *Strict Monotonicity at 0* (SM0): $M(\mathbf{u}; \mathbf{w}) < M(\mathbf{u} + \varepsilon \mathbb{1}_i; \mathbf{w})$.
- 13) *Strict Weak Transfer Principle* (SWTP): Let $i \doteq \operatorname{argmin}_i \mathbf{u}_i$, $j \doteq \operatorname{argmax}_j \mathbf{u}_j$. If $\mathbf{u}_i \neq \mathbf{u}_j$, then *there exists* some $\varepsilon > 0$ s.t. $\mathbf{u}_i + w_j\varepsilon < \mathbf{u}_j - w_i\varepsilon$ and $W(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w}) > W(\mathbf{u}; \mathbf{w})$ for welfare, or $\Lambda(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w}) < \Lambda(\mathbf{u}; \mathbf{w})$ for malfare.

- 14) *Strict PDTP* (SPDTP): Suppose $\mathbf{u}_i + w_j\varepsilon < \mathbf{u}_j - w_i\varepsilon$. Then $W(\mathbf{u} + w_j\varepsilon \mathbb{1}_i - w_i\varepsilon \mathbb{1}_j; \mathbf{w}) > W(\mathbf{u}; \mathbf{w})$ for welfare, or $\Lambda(\mathbf{u} + w_j\varepsilon \mathbb{1}_i - w_i\varepsilon \mathbb{1}_j; \mathbf{w}) < \Lambda(\mathbf{u}; \mathbf{w})$ for malfare.

Lemma 3.7 (Consequences of Strong Axioms). Suppose power-mean aggregator function $M_p(\cdot; \mathbf{w})$. The following then hold:

- 1) Strengthening the SM axiom (i.e., $1 \rightarrow 12$) implies $p > 0$.
- 2) Strengthening the WTP axiom (i.e., $8 \rightarrow 13$) implies $p \neq 1$, thus $p < 1$ for welfare and $p > 1$ for malfare.
- 3) Strict PDTP (i.e., $11 \rightarrow 14$) implies $p \neq 1$ and $p \neq \pm\infty$, thus $p \in (-\infty, 1)$ for welfare and $p \in (1, \infty)$ for malfare.

We now comment on the philosophical implications of lemma 3.7. The consequences of SM0 are immense: the “brand name” Nash social welfare ($p = 0$) is now inadmissible as a fair welfare function, and moreover we reduce the unbounded spectrum of p to just $p \in (0, 1]$. Intuitively, the “strict” aspect of SM0 encodes the idea that gains to utility should *always* be relevant, and consequently prevents a form of “minority rule,” wherein a group with utility 0 ensures that welfare can not possibly improve without benefiting said group. Note that for any $p < 1$, the weighted relative impact of helping disadvantaged groups is still higher than privileged groups (as can be seen by inspecting the power-mean gradient, see lemma 3.10), but SM0 puts a sharp limit on the strength of this effect by preventing $p \leq 0$.

Under SWTP, pure utilitarianism ($p = 1$) is inadmissible: intuitively, transferring any ε utility would, by linearity, not change welfare, thus not yield strict improvement. In this sense, SWTP incentivizes equitable redistribution of wealth more strongly than does WTP. Strengthening of PDTP to strict inequality is also interesting, but it necessarily precludes both the utilitarian ($p = 1$) and egalitarian cases ($p \in \pm\infty$), and thus does not actually represent a strict preference towards egalitarianism.

The following axioms control the “degree of prioritarianism” more precisely than do our transfer principles. Each describes the behavior of a welfare function under situations of *extreme inequality*, where some group receives 0 or ∞ utility.

Axiomatization 3.8 (Extreme Axioms). Suppose as in axiomatization 3.1. We now define two additional axioms.

- 15) *Zero Barrier* (OB): $\lim_{\mathbf{u}_i \rightarrow 0^+} W(\mathbf{u}; \mathbf{w}) = 0$.
- 16) *Finite Ceiling* (FC): $\lim_{c \rightarrow \infty} W(\mathbf{u} + c \mathbb{1}_i; \mathbf{w}) < \infty$.

Nozick [1974] criticizes utilitarianism via *reductio ad absurdum* by positing a “utility monster,” which derives extremely high utility from some good, and thus utilitarian theory dictates we must allocate all resources to the monster. Our *zero barrier axiom* (15) promotes prioritarianism by ensuring that welfare $W(\mathbf{u}; \mathbf{w}) \rightarrow 0$ as any group utility $\mathbf{u}_i \rightarrow 0$, a property shared by the egalitarian welfare, and furthermore, the zero barrier axiom incentivizes aiding the most needy by creating a “barrier” at 0 utility, thus disincentivizing ex-

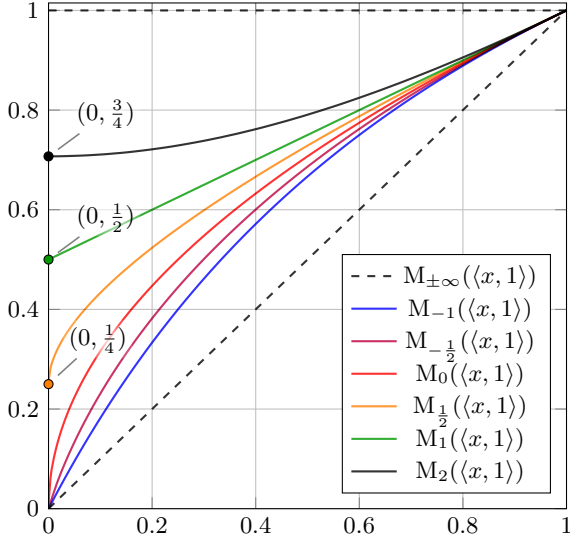


Figure 1: Plots of the unweighted power-mean (i.e., $\mathbf{w} = \langle \frac{1}{2}, \frac{1}{2} \rangle$) for various p . Observe that the region around $x = \mathbf{u}_1 = 0$, wherein Lipschitz discontinuities may occur, exhibits sharp changes to welfare, as $M_p(\langle x, 1 \rangle)$ is very sensitive to small changes to x .

treme harm to any group, and bounding the harm caused by the utility monster. Similarly, the *finite ceiling axiom* (16) ensures that even as some $u_i \rightarrow \infty$, $W(\mathbf{u}; \mathbf{w})$ remains *finite*, i.e., the disadvantaged (finite utility) groups are not “forgotten” in the monster’s rush toward infinite utility (and infinite inequality).

Lemma 3.9 (Consequences of Extreme Axioms). Suppose as in lemma 3.7. The following then hold.

- 1) 0B (axiom 15) $\Leftrightarrow p \leq 0$. 2) FC (axiom 16) $\Leftrightarrow p < 0$.

Observe that, subject to axioms 1–8, $\text{FC} \implies \text{0B} \implies \text{SWTP} \implies \text{WTP}$. The above “egalitarian” framing of the 0B and FC axioms is complemented by a “utilitarian” framing, wherein we would require that taking the utility of any group to 0 *does not* take welfare to 0, or that taking the utility of any group to ∞ *does* unboundedly increase welfare, in each case concluding the complementary set of permissible p in lemma 3.9, though again such axioms codify intuition, but can not dictate the “correct” welfare function.

Some authors also assume variants of the *Independence of Irrelevant Alternatives* (IIA) axiom, which restricts to the Nash social welfare [Roth et al., 1977, Kaneko and Nakamura, 1979], i.e., $p = 0$. We do not discuss this axiom further, but note that in our framework, it is equivalent to jointly assuming axiom 15 and the utilitarian form of 16.

3.4 Continuity Properties of Aggregator Functions

Another axiom that can be strengthened to great effect is the *continuity* axiom. Myriad varieties of continuity exist, and we investigate the Hölder and Lipschitz varieties in great detail here. In particular, while it does make sense in principle

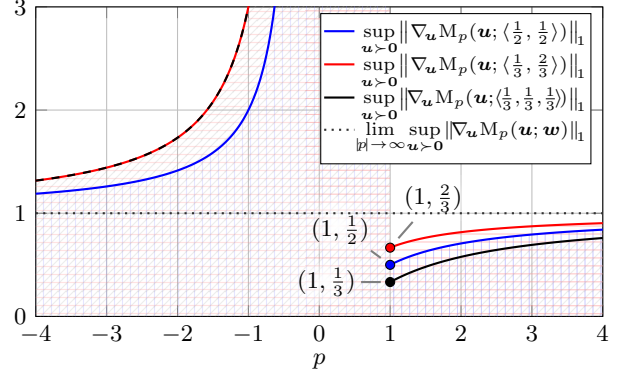


Figure 2: Plots of $\|\cdot\|_1$ Lipschitz constants of weighted power-means for various weightings as functions of p . Note that $\sup_{\mathbf{u} > \mathbf{0}} \|\nabla_{\mathbf{u}} M_p(\mathbf{u}; \mathbf{w})\|_1 = \sup_{\mathbf{u}, \mathbf{u}' > \mathbf{0}} \frac{|M(\mathbf{u}; \mathbf{w}) - M(\mathbf{u}'; \mathbf{w})|}{\|\mathbf{u} - \mathbf{u}'\|_1}$

describes this Lipschitz constant, as $M(\cdot; \mathbf{w})$ is assumed to be continuous, and for any \mathbf{w} , this quantity approaches 1 as $p \rightarrow \pm\infty$ (plotted as a dashed line). The region below each Lipschitz constant plot is shaded and patterned, to emphasize that *higher values* allow for *sharper rates of change* to the power-mean functions.

to axiomatically assume a stronger notion of continuity, due to their parameterized nature, such characterizations lack the elegant simplicity of our axiomatization. We present the material in this manner, so as to emphasize the consequences of choice of aggregator function on continuity, rather than how the act of assuming continuity impacts the aggregator function.

Previous works bound deviations between power-mean welfare functions [Cousins, 2021, 2022], and analyze their Lipschitz continuity [Beliakov et al., 2009]. We extend this analysis to power-mean welfare functions (i.e., $p \leq 1$), showing that they are Lipschitz continuous for $p < 0$, though not for $p \in [0, 1)$, and the Lipschitz constants depend on the *minimum weight* w_{\min} . This is initially surprising, as intuitively, low-weight groups should have little impact on the power mean, however we know that for $p \leq 0$, by lemma 3.9 item 1, as any group’s sentiment $u_i \rightarrow 0$, then so too must the power mean $M_p(\mathbf{u}; \mathbf{w}) \rightarrow 0$, thus as $w_i \rightarrow 0$, this must occur more rapidly, hence the dependence on w_{\min} .

While it is not unreasonable to axiomatically assume a stronger notion of continuity, due to their parameterized nature, such characterizations lack the elegant simplicity of our axiomatization. We thus present continuity properties as *consequent from* choice of welfare function, rather than *vice versa*, to reflect the practical impact of this choice. We now analyze the local behavior of power means, first through their gradients, and then their Lipschitz and Hölder continuity properties. The reader is invited to reference figure 1 throughout, wherein various power-means are plotted, revealing their pathological and highly nonlinear behavior for $p \approx 0$.

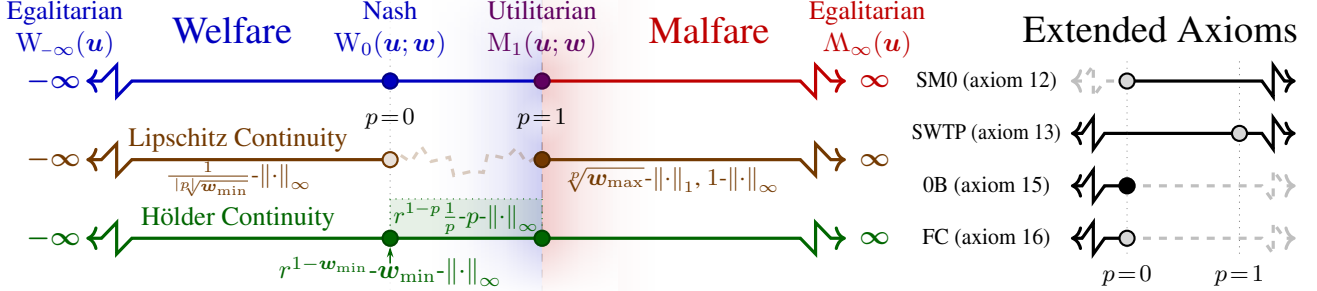


Figure 3: Illustration of power-mean properties and axioms. Solid lines and filled circles denote the values of p that concord with an axiom or property, while dashed lines and unfilled circles denote their complement. Basic properties (theorem 3.5) and continuity (lemmata 3.12 & 3.13) are plotted on the left, and the consequences of our strong and extended axioms (lemmata 3.7 & 3.9) are shown on the right.

Lemma 3.10 (Power-Mean Differentiation). Suppose $\mathbf{u}_{\setminus i} \succ \mathbf{0}$, some weights vector $\mathbf{w} \in \Delta_g$, and $p \in \mathbb{R}$. The power mean then differentiates in \mathbf{u}_i as follows.

- 1) If $\mathbf{u}_i > \mathbf{0}$, then $\frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})}$.
- 2) If $p < 0$, then $\lim_{\mathbf{u}_i \rightarrow 0^+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = -\sqrt[p]{\frac{1}{\mathbf{u}_i}}$.
- 3) If $p \in [0, 1)$, then $\lim_{\mathbf{u}_i \rightarrow 0^+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = \infty$.

Definition 3.11 (Lipschitz and Hölder Continuity). An aggregator function $M(\mathbf{u}; \mathbf{w})$ is *Hölder continuous* in the variable \mathbf{u} w.r.t. some norm $\|\cdot\|_M$ over domain $[0, r]^g$ if there exist some *scale* $\lambda \geq 0$ and *power* $\alpha \in (0, 1]$, such that for all $\mathbf{u}, \mathbf{u}' \in [0, r]^g$, it holds that

$$|M(\mathbf{u}; \mathbf{w}) - M(\mathbf{u}'; \mathbf{w})| \leq \lambda \|\mathbf{u} - \mathbf{u}'\|_M^\alpha. \quad (2)$$

We say that such a function is $\lambda\text{-}\alpha\text{-}\|\cdot\|_M$ Hölder continuous, and if $\alpha = 1$, it is $\lambda\text{-}\|\cdot\|_M$ Lipschitz continuous, and if $\lambda \leq 1$ it is $\|\cdot\|_M$ *nonexpansive*.

Intuitively, Lipschitz continuity bounds *infinitesimal rates of change*, and Hölder continuity bounds *change over small regions*, where the relative size of the change grows larger as $\|\cdot\|_M \rightarrow 0^+$. The relationship between Lipschitz, Hölder, and standard ε - δ limit continuity for range $r \|\cdot\|_M$ is

$$\lambda \text{ Lipschitz} \implies \lambda r^{-\alpha} \text{ Hölder} \implies \varepsilon\text{-}\delta \text{ Limit}.$$

As we assume *continuity* throughout (axiom 4), all aggregator functions of interest are tautologically ε - δ limit continuous, however we shall see that they do not all share the same Hölder and Lipschitz continuity properties, which has crucial implications for sampling and learning from data, as well as privacy and algorithmic stability. The following result (visualized in figure 2) analyzes Lipschitz continuity.

Lemma 3.12 (Power-Mean Lipschitz Continuity). Suppose $p \in \mathbb{R}$, sentiment vectors $\mathbf{u}, \mathbf{u}' \in \mathbb{R}_{0+}^g$, and weights vector $\mathbf{w} \in \Delta_g$. The following then hold.

- 1) Suppose $p \geq 1$. Then $M_p(\cdot; \mathbf{w})$ is $\sqrt[p]{\mathbf{w}_{\max}}\text{-}\|\cdot\|_1, 1\text{-}M_p(\cdot; \mathbf{w})$, and $1\text{-}\|\cdot\|_\infty$ Lipschitz.
- 2) Suppose $p < 0$. Then $M_p(\cdot; \mathbf{w})$ is $\frac{1}{\sqrt[p]{\mathbf{w}_{\min}}}\text{-}\|\cdot\|_\infty$ Lipschitz.

While the $p \in [0, 1)$ power-means are not Lipschitz continuous (see lemma 3.10 item 3), we find that they are still

Hölder continuous, which largely results in similar enough properties for our purposes.

Lemma 3.13 (Power-Mean Hölder Continuity). Suppose $\mathbf{u} \in [0, r]^g$, group index $i \in \mathcal{G}$, weights vector $\mathbf{w} \in \Delta_g$, and assume where appropriate that $\mathbf{u}_i + \varepsilon \leq r$. The power mean then obeys the following Hölder continuity criteria.

- 1) *Generic Welfare Hölder Condition*: Suppose $p \leq 1$. Then $|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})| \leq r^{1-\mathbf{w}_i} \varepsilon^{\mathbf{w}_i}$, and $M_p(\cdot; \mathbf{w})$ is $r^{1-\mathbf{w}_{\min}}\text{-}\mathbf{w}_{\min}\text{-}\|\cdot\|_\infty$ Hölder continuous.
- 2) *Positive Welfare Hölder Condition*: Suppose $p \in (0, 1]$. Then $|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})| \leq r^{1-p} \frac{\mathbf{w}_i}{p} \varepsilon^p$. Furthermore, $M_p(\cdot; \mathbf{w})$ meets the following Hölder conditions:
 - A) $r^{1-p} \frac{\mathbf{w}_{\max}}{p} \text{-} p\text{-}\|\cdot\|_1$;
 - B) $r^{1-p} \frac{1}{p} \text{-} p\text{-}M_1(\cdot; \mathbf{w})$; &
 - C) $r^{1-p} \frac{1}{p} \text{-} p\text{-}\|\cdot\|_\infty$.

Applications Understanding the response of power-mean functions to small changes to sentiment, i.e., their gradients and continuity properties, is highly relevant to *privacy*, *adversarial robustness*, *algorithmic stability*, *strategy proofness*, and *statistical learnability*. Due to their commonalities, we briefly treat the first four here, while sections 4 & 5 theoretically and experimentally explore statistical learning in detail. To clarify the relationship of our axioms to these properties, we visualize them in figure 3.

We first assume that the *parameters* or *decisions* made by some algorithm are robust to small changes to the objective (welfare), and note that Lipschitz or Hölder continuity describe how robust the objective is to small changes to sampled, estimated, or queried per-group utility values.³ Lipschitz continuity is highly relevant to *differential privacy* [Bassily et al., 2019, 2020, Wang et al., 2022, Patel et al., 2022], as differential privacy is sensitive to changes to algorithm output caused by *individual-level changes*, thus sensitivity to infinitesimal change is paramount.

For *algorithmic stability*, *adversarial robustness*, and *strat-*

³Note that this applies to some classes of stable algorithms and objectives. For example, algorithms maximizing *strongly concave* welfare functions can only *boundedly change* the optimal parameters or output under bounded change to the objective function.

egy proofness, if *individual-level change* is too small to cause harm, Hölder continuity is powerful, since while *cumulative change* to utility is linear in the number of colluding parties, its *impact on welfare* is sublinear, due to the α -power. In other words, here individuals may be powerful, but the cumulative effect of collusion is less than the sum of its parts.

In the context of *algorithmic stability* for convex optimization, it is well-known that if the gradient is *bounded* and *Lipschitz continuous* (i.e., the function is *smooth*), we obtain desirable bounds on the generalization error and computational cost of (stochastic) gradient descent methods. Unfortunately, even if the underlying per-group sentiment functions are smooth with bounded gradient, these properties are not necessarily preserved by the power-mean (e.g., egalitarian objectives create nondifferentiabilities where the maximum or minimum sentiment group changes). However, more recent work [Lei and Ying, 2020] shows that a Hölder continuous gradient also suffices: in particular, it suffices to have $\|\nabla_{\theta}M(\mathbf{u}(\theta); \mathbf{w}) - \nabla_{\theta}M(\mathbf{u}(\theta'); \mathbf{w})\|_2 \leq \lambda\|\theta - \theta'\|_2^{\alpha}$. This condition is a bit more subtle than that which we consider, but observe that it may be bounded for appropriate model classes (which define the space of θ) using lemmata 3.10 & 3.13 and the classical chain rule from the calculus of infinitesimals.

Outside of some convenient median-based constructs, e.g., linear regression under absolute loss [Chen et al., 2018b], true strategy-proofness in learning settings can be quite difficult to achieve [Procaccia, 2008]. However, under the ε -*truthfulness* assumption [Meir and Rosenschein, 2011], agents only lie about their labels if they gain at least ε utility (or lose at least ε disutility) by doing so. In this setting, the above robustness analysis is relevant, as we have effectively shown that, assuming bounded loss values, Hölder continuity implies that small colluding groups have small impact on the objective, and thus (assuming, e.g., strong convexity in the neighborhood of the optimal h^*) lying only weakly impacts \hat{h} . Therefore, if the hypothesis space \mathcal{H} is *continuous* (e.g., bounded linear regression, but not discrete linear classification), they can not strongly benefit from doing so. We thus conclude that robustness and stability, which are preserved by Hölder continuous objectives, are sufficient for ε -truthful strategy-proofness in some learning settings. We thus conclude that, via lemmata 3.12 & 3.13, we can analyze the privacy, stability, and statistical properties of many welfare-based algorithms for arbitrary power-mean welfare functions.

4 Generalizing Fair-PAC Learning

Suppose now that we seek to estimate or optimize some welfare function, but do not know the utility values of each group, and must instead estimate them via sampling (i.e., from data). We predominantly study the *plug-in estimator*,

which approximates the welfare of expected utility values with the welfare of *empirical mean* utilities over m samples from each group’s instance distribution \mathcal{D}_i over *labeled instance space* $(\mathcal{X} \times \mathcal{Y})$, i.e., features $x \in \mathcal{X}$ and labels $y \in \mathcal{Y}$. We assume a hypothesis class $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}'$ mapping inputs \mathcal{X} to prediction space \mathcal{Y}' , and a utility function $u : \mathcal{Y}' \times \mathcal{Y} \rightarrow \mathbb{R}_{0+}$ that assesses the quality of prediction $\hat{y} \in \mathcal{Y}'$ given true label $y \in \mathcal{Y}$, and we thus express the true utility of hypothesis $h \in \mathcal{H}$ for group i as $\mathbb{E}_{\mathcal{D}_i}[u \circ h]$ and the empirical estimate of utility as $\hat{\mathbb{E}}_{\mathcal{D}_i}[u \circ h]$, where $(u \circ h)(x, y) \doteq u(h(x), y)$. Notation introduced in this section is summarized in table 2.

We first briefly outline computational concerns and build intuition for welfare optimization and FPAC learning in section 4.1. We then discuss asymptotic convergence and the fundamentals of statistical estimation of $W(i \mapsto \mathbb{E}_{\mathcal{D}_i}[u \circ h]; \mathbf{w})$ for a single function h in section 4.2. Section 4.3 then introduces and analyzes a concept of *fair-PAC learnability*, showing worst-case bounds on the sample complexity of estimating and optimizing welfare functions over a function family \mathcal{H} . Finally, section 4.4 discusses applying progressive sampling techniques to *adaptively query* to avoid worst-case sample complexity, and instead near-optimally sample to address a given optimization task.

4.1 On Compromise and Fair Learning

In standard supervised learning or risk minimization settings, given distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, we generally seek to approximate

$$h^* \doteq \operatorname{argmin}_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [-(u \circ h)(x, y)] ,$$

whereas in FPAC learning, we seek to approximate

$$h^* \doteq \operatorname{argmax}_{h \in \mathcal{H}} W \left(i \mapsto \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [(u \circ h)(x, y)]; \mathbf{w} \right) .$$

The change of sign is entirely superficial, and standard first-order convex optimization techniques can be applied in either setting [Cousins, 2021], assuming appropriate structure of $u \circ \mathcal{H} \doteq \{u \circ h | h \in \mathcal{H}\}$ and the model parameters θ , but the impact of the welfare function and the per-group distributions on both h^* and the optimization dynamics are worth discussing.

Ideally, we would be able to optimize a model by making only local modifications that benefit all groups, but in general this is not realistic. Often decisions must be made during training that benefit some groups while harming others, which we term “disagreement.” Disagreement arises *locally*, on the scale of individual changes made to the model during gradient updates, but also *globally*, in the sense that the overall model learned is a *compromise* between the wants and needs of each group. Some sources of disagreement are obvious, such as two groups disagreeing on how to classify a contiguous region of \mathcal{X} , but others are subtler. For

Object	Definition or Space	Description
$\mathcal{X}, \mathcal{Y}, \mathcal{Y}'$	—	Supervised prediction domain, label space, and codomain
h	$\in \mathcal{X} \rightarrow \mathcal{Y}'$	Hypothesis (possible supervised learning model)
\mathcal{H}	$\subseteq \mathcal{X} \rightarrow \mathcal{Y}'$	Hypothesis class (class of possible supervised learning models)
$u(\hat{y}, y)$	$\in \mathcal{Y}' \times \mathcal{Y} \rightarrow [0, r]$	Utility function (score prediction \hat{y} given true label y)
$u \circ h$	$(u \circ h)(x, y) \doteq u(h(x), y), \in (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, r]$	Utility composition
$u \circ \mathcal{H}$	$\doteq \{u \circ h \mid h \in \mathcal{H}\}, \subseteq (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, r]$	Utility class (composition of utility function and hypothesis class)
\mathcal{D}_i	Over $(\mathcal{X} \times \mathcal{Y})$	Distribution over $(\mathcal{X}, \mathcal{Y})$ pairs pertaining to group i
m	$\in \mathbb{Z}_+$	Sample size

Table 2: Functions, spaces, and common variables in FPAC learning.

example, all finite learning models have learning capacity limitations, be they *finite neurons counts* in artificial neural networks, *norm constraints* on parameter vectors, or *depth bounds* in decision trees, and this paucity begets disagreement between groups, each of whom would prefer structures be learned that benefit themselves. Learners must balance between competing preferences that arise when groups have *direct disagreement*, i.e., conflicting preferences for \mathcal{Y} given \mathcal{X} , or *indirect disagreement*, i.e., agreement on \mathcal{Y} given \mathcal{X} , but disagreement as to what to *prioritize*, due to different distributions over \mathcal{X} .

It is enlightening to consider the process by which first-order optimization trains a model. Ideally, under complete agreement, each group would have the same response $f(\theta)$ in expected or empirical sentiment to the parameters θ . In this case, since $W(i \mapsto f(\theta); \mathbf{w}) = f(\theta)$ (by axioms 6 & 7), it holds that $\nabla_{\theta} W(i \mapsto f(\theta); \mathbf{w}) = \nabla_{\theta} f(\theta)$, thus optimization proceeds as in standard gradient ascent on expected utility. However, complete agreement is unrealistic, as different groups generally respond differently to different models.

Starting from an arbitrary parameterization θ_0 , it is often possible to improve utility for all groups simultaneously by first focusing on aspects of the learning problem relevant to all groups (i.e., initially groups may agree on how to proceed in a mutually beneficial manner). However, once Pareto-dominance is achieved, we enter the realm of *disagreement*, as no decision will benefit any group without harming others. Observe now that for $p \leq 1$ welfare functions, the power-mean derivative $\frac{\mathbf{w}_i \mathbf{u}_i}{W_i^{p-1}(\mathbf{u}; \mathbf{w})}$ (see lemma 3.10) *depends inversely* on group utility \mathbf{u}_i . Therefore, infinitesimal movement along the gradient favors equitable transfer from high-utility to low-utility groups, and gradient ascent steps discretely approximate this process.

The welfare-optimal model is thus an equilibrium between the forces of preferentially increasing low-utility group utilities and those utilities rising relative to their peers. This tension in training dynamics results in *compromise* between groups, and in some sense the whole process has more in common with *adversarial learning* than standard maximization tasks; indeed for $p = -\infty$, we have an explicit maximin objective, albeit one with a finite inner minimization set. For these reasons, the equilibrium model h^* that results from

welfare maximization generally differs greatly from any of the models each group would select for themselves. Thus far, we have discussed only the *computational* aspects of welfare-based learning tasks, but the remainder of this work is focused on the *statistical aspect* of this setting.

4.2 The Fundamentals of Estimation

Before treating the intricacies of learning (optimization) over a *class of functions* \mathcal{H} , we first treat the subject of *estimating* utility values via sampling for a single function h , and we then discuss welfare estimates via the plug-in estimator. Namely, there exists some true expected utility vector $\mathbf{u} \doteq i \mapsto \mathbb{E}_{(x,y) \sim \mathcal{D}_i}[(u \circ h)(x, y)]$, and we have some estimate $\hat{\mathbf{u}}$, and we ask the questions, “How well does $\hat{\mathbf{u}}$ approximate \mathbf{u} ?” and moreover, “How well does the plug-in welfare estimate $W(\hat{\mathbf{u}}; \mathbf{w})$ approximate the true welfare $W(\mathbf{u}; \mathbf{w})$?”

In estimation settings, we generally assume some *consistent estimator* $\hat{\mathbf{u}}$ of true utility $\mathbf{u} \in [0, r]^g$, i.e., some $\hat{\mathbf{u}}$ such that $\lim_{m \rightarrow \infty} \hat{\mathbf{u}} = \mathbf{u}$. For our purposes, the empirical mean over m samples per group will suffice. Here, for finite \mathbf{u} , using only ε - δ *limit continuity* (of welfare) and the *weak law-of-large-numbers*, we have *consistency* of the plug-in welfare estimate, i.e., $\lim_{m \rightarrow \infty} W(\hat{\mathbf{u}}; \mathbf{w}) = W(\mathbf{u}; \mathbf{w})$, but in the grand scheme of things, this is a rather weak guarantee. If we assume finite variance, then for any norm $\|\cdot\|_{\mathbb{W}}$ and any *failure probability* δ , by the central limit theorem, we have

$$\lim_{m \rightarrow \infty} \mathbb{P}_{\hat{\mathbf{u}}} \left(\frac{1}{\sqrt{m}} \|\mathbf{u} - \hat{\mathbf{u}}\|_{\mathbb{W}} \leq \sqrt{2v \ln \frac{1}{\delta}} \right) \geq 1 - \delta ,$$

for *variance proxy* $v \doteq \mathbb{V}_{\hat{\mathbf{u}}}[\|\hat{\mathbf{s}}\|_{\mathbb{W}}]$. In general, given some weighted welfare function $W(\cdot; \mathbf{w})$, we wish to bound the *estimation error* $|W(\hat{\mathbf{u}}; \mathbf{w}) - W(\mathbf{u}; \mathbf{w})|$ of the plug-in estimator.

Table 3 precisely characterizes the asymptotic convergence rates of estimation error under each continuity concept studied in this work. We find root-hyperbolic rates under Lipschitz continuity, and $\frac{\alpha}{2}$ power-law rates under Hölder continuity, though in all cases, the variance proxy v (and thus the norm $\|\cdot\|_{\mathbb{W}}$) also plays a substantial role.

Despite the above positive results, significant challenges arise in estimating the welfare of even a single h . While

Continuity Concept	Deviation Bound	Asymptotic Error Convergence Rate	Sample Complexity
ε - δ Limit	$ W(\hat{\mathbf{u}}; \mathbf{w}) - W(\mathbf{u}; \mathbf{w}) \leq$ —	$\lim_{m \rightarrow \infty} \varepsilon = 0$	$m_W(\varepsilon, \delta) < \infty$
λ - α Hölder	$\lambda \ \mathbf{u} - \hat{\mathbf{u}}\ _M^\alpha$	$\varepsilon \lesssim \lambda \left(\frac{2v \ln \frac{1}{\delta}}{m} \right)^{\alpha/2}$	$m_W(\varepsilon, \delta) \lesssim \frac{2\lambda^{2/\alpha} v \ln \frac{1}{\delta}}{\varepsilon^{2/\alpha}}$
λ Lipschitz	$\lambda \ \mathbf{u} - \hat{\mathbf{u}}\ _M$	$\varepsilon \lesssim \lambda \sqrt{\frac{2v \ln \frac{1}{\delta}}{m}}$	$m_W(\varepsilon, \delta) \lesssim \frac{2\lambda^2 v \ln \frac{1}{\delta}}{\varepsilon^2}$
Nonexpansive	$\ \mathbf{u} - \hat{\mathbf{u}}\ _M$	$\varepsilon \lesssim \sqrt{\frac{2v \ln \frac{1}{\delta}}{m}}$	$m_W(\varepsilon, \delta) \lesssim \frac{2v \ln \frac{1}{\delta}}{\varepsilon^2}$

Table 3: Continuity concepts and asymptotic estimation guarantees. For each type of continuity studied herein, we bound the *estimation error* $\varepsilon = |M - \hat{M}|$, as well as the *sample complexity* $m_W(\varepsilon, \delta)$, which is the number of samples required to attain ε estimation error with probability at least $1 - \delta$. Here $a \lesssim b$ denotes *asymptotic inequality*, which can be formalized as $\lim_{m \rightarrow \infty} \frac{a(m)}{b(m)} \leq 1$ or $\lim_{\varepsilon \rightarrow 0^+} \frac{a(\varepsilon)}{b(\varepsilon)} \leq 1$. The asymptotic approximation stems from the central limit theorem, and upper-bounds are due to continuity-based deviation bounds and Gaussian-Chernoff tail bounds.

the plug-in welfare estimator preserves the *consistency* of the utility estimator $\hat{\mathbf{u}}$, no unbiased estimator for welfare functions (unless $p = 1$ or stringent assumptions are made on the distributions $\mathcal{D}_{1:g}$) essentially because $W(\mathbf{u}; \mathbf{w})$ is a nonlinear operator. There are also deep questions as to how to optimally allocate sampling effort, since the variances of each \hat{u}_i may differ, as may the impact of each \hat{u}_i on $W(\hat{\mathbf{u}}; \mathbf{w})$, both of which directly impact the variance proxy v . These questions are probed in detail by Cousins [2022], but in this work, we generally assume m samples for each group, which greatly simplifies the matter, and never requires more than a factor g more samples than any nonuniform allocation of sampling effort.

The above analysis is straightforward, but it does not scale well from a single h to learning over a large hypothesis class \mathcal{H} , and moreover the behavior of the quantity v depends intimately on $\|\cdot\|_W$ and the sampling distributions $\mathcal{D}_{1:g}$. Going forward, we shall bound the generalization error and sample complexity of welfare estimation by first separately bounding each group’s utility values, and then taking a union bound over groups. This technique introduces $\ln \frac{g}{\delta}$ terms instead of the $\ln \frac{1}{\delta}$ terms of our central limit theorem argument (see table 3), but it allows us to express our results in terms of readily-available generalization error and sample complexity bounds for per-group (scalar) utilities.

4.3 Fair-PAC Learning and Learnability

We first show that the estimation error of welfare $W(\mathbf{u}; \mathbf{w})$ may be bounded in terms of the error of each utility value u_i . From there, we bound the sample complexity of optimization, and describe a notion of fair-PAC (FPAC) learnability for welfare functions, wherein the goal is to uniformly bound the number of samples required to learn in \mathcal{H} .

We abstract away the statistical details of this estimation process by assuming, for each group i , a known bound on the *supremum deviation* of the expected utility for each $h \in \mathcal{H}$, i.e., a bound of the form

$$\forall i : \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{\mathcal{D}_i} [u \circ h] - \hat{\mathbb{E}}_{\mathcal{D}_i} [u \circ h] \right| \leq \varepsilon_i. \quad (3)$$

The details of obtaining such bounds with high probability are well-studied, and Cousins [2022] discusses them under the name *additive error vector* (AEV) bounds in group-fairness settings, showing that they can be obtained via the Chernoff method [Bennett, 1962, Hoeffding, 1963, Boucheron et al., 2013], Rademacher averages [Bartlett and Mendelson, 2002, Shalev-Shwartz and Ben-David, 2014, Mitzenmacher and Upfal, 2017, Cousins and Riondato, 2020], or other such tools.

In particular, using Rademacher averages, bounds on the supremum deviation often take the form

$$\mathcal{O} \left(\frac{r \ln \frac{g}{\delta}}{m} + \sqrt{\frac{C + r^2 \ln \frac{g}{\delta}}{m}} \right),$$

where C depends on \mathcal{D}_i and \mathcal{H} , and can generally be bounded as $v \ln n$, where $v \doteq \sup_{h \in \mathcal{H}} \mathbb{V}_{\mathcal{D}} [s \circ h]$ is the *supremum variance* of utility over the class \mathcal{H} , and n measures the *effective size* of the hypothesis class. The details are beyond the scope of this work, but $n = |\mathcal{H}|$ suffices.⁴ Observe that here v plays the same fundamental role as in single-function estimation (as discussed in section 4.2), n essentially corrects for the *multiple comparisons problem* (i.e., our bound holds across all of \mathcal{H} simultaneously) with similar form to union bounds, and the $\mathcal{O} \frac{r \ln \frac{1}{\delta}}{m}$ term acts as an (asymptotically negligible) *finite-sample correction* to a tail bound that otherwise resembles a Gaussian-Chernoff bound.

With the table set, we now bound the estimation error of welfare objectives. The following result holds essentially by the definition of Hölder continuity, and can immediately be applied to any model class \mathcal{H} for which bounds on the supremum deviation are available.

⁴Furthermore, taking n to bound the cardinality of projecting $u \circ \mathcal{H}$ onto m_i samples also suffices (via, e.g., the Vapnik-Chervonenkis dimension), and more sophisticated metrics of effective size, such as *covering numbers*, i.e., the smallest cover size n such that there exists some $\mathcal{H}' \subseteq \mathcal{X} \rightarrow \mathcal{Y}$ such that $|\mathcal{H}'| = n$ and each $h \in \mathcal{H}$ can be well-approximated by some $h' \in \mathcal{H}'$.

Theorem 4.1 (Hölder Continuity and Welfare Optimality). Suppose $W(\cdot; \mathbf{w})$ is $\lambda\text{-}\alpha\text{-}\|\cdot\|_{\mathcal{W}}$ Hölder continuous w.r.t. some norm $\|\cdot\|_{\mathcal{W}}$, and additive error bounds ε that obey (3). Then

$$\sup_{h \in \mathcal{H}} \left| W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ h]; \mathbf{w}]\right) - W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbb{E}[u \circ h]; \mathbf{w}]\right) \right| \leq \lambda \|\varepsilon\|_{\mathcal{W}}^{\alpha}.$$

Consequently, the *empirical welfare maximizer*

$$\hat{h} \doteq \sup_{h \in \mathcal{H}} W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbb{E}[u \circ h]; \mathbf{w}]\right)$$

approximates the *true welfare maximizer*

$$h^* \doteq \sup_{h^* \in \mathcal{H}} W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ h^*]; \mathbf{w}]\right),$$

in terms of welfare-optimality, as it holds that

$$W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ \hat{h}]; \mathbf{w}]\right) \geq W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ h^*]; \mathbf{w}]\right) - 2\lambda \|\varepsilon\|_{\mathcal{W}}^{\alpha}.$$

This result is easily interpreted and practically relevant. Furthermore, it is readily applied to any power-mean welfare function via lemmata 3.12 & 3.13, whereas prior work [Cousins, 2021, 2022] only handles the case of Lipschitz-continuous aggregator functions.

In PAC-learning theory, it is standard to analyze the difficulty of learning as a function of the *complexity of the model class* \mathcal{H} . To this end, we assume \mathcal{H} is parameterized by D variables $\mathbf{d} \in \mathbb{R}_{0+}^D$, e.g., the dimension of a linear classifier [Shalev-Shwartz and Ben-David, 2014], or for neural network families, a vector of per-layer widths [Anthony and Bartlett, 2009] or per-layer norm constraints [Bartlett et al., 2017], where $\mathcal{H}_{\mathbf{d}}$ denotes the class parameterized by \mathbf{d} .

Furthermore, we seek to know *how much data* we need to probabilistically learn an objective to within a given error tolerance ε , rather than *how well* we can learn a concept with a given sample size. Henceforth, a *sample complexity function* $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d})$ is some function such that, for utility range r , any probability distribution \mathcal{D} , a sample size of at least $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d})$ ensures that $u \circ \mathcal{H}$ is *uniformly convergent*, i.e., with probability at least $1 - \delta$, it holds that

$$\sup_{h \in \mathcal{H}_{\mathbf{d}}} \left| \mathbb{E}_{\mathcal{D}}[u \circ h] - \hat{\mathbb{E}}_{\mathcal{D}}[u \circ h] \right| \leq \varepsilon. \quad (4)$$

Similarly, we express the (per-group) sample complexity of uniformly estimating a *welfare function* $W(\cdot; \mathbf{w})$ as $m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, g, r, \mathbf{d})$, requiring with probability at least $1 - \delta$, it holds that

$$\sup_{h \in \mathcal{H}_{\mathbf{d}}} \left| W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ h]; \mathbf{w}]\right) - W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbb{E}[u \circ h]; \mathbf{w}]\right) \right| \leq \varepsilon. \quad (5)$$

We now analyze the sample complexity of welfare estimation, and subsequently FPAC learning with welfare objectives.

Theorem 4.2 (Welfare Sample Complexity). Suppose sample complexity function $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d})$ for hypothesis class \mathcal{H} , and some welfare function $W(\cdot; \mathbf{w})$ that is $\lambda\text{-}\alpha\text{-}\|\cdot\|_{\infty}$ Hölder continuous. Then the sample complexity function

$$m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, g, r, \mathbf{d}) \leq m_{\mathcal{H}}\left(\sqrt{\frac{\varepsilon}{\lambda}}, \frac{\delta}{g}, r, \mathbf{d}\right)$$

is sufficient, i.e., for at least this many samples from each of the g groups, (5) holds. Moreover, for this sample size, with probability at least $1 - \delta$, the empirical welfare maximizer is 2ε -optimal.

From these uniform generalization error and sample complexity bounds, we can show that *classes of welfare functions* are FPAC learnable, defined as follows.

Definition 4.3 (Fair-PAC Learning). Suppose *hypothesis class* $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}'$ parameterized by $\mathbf{d} \in \mathbb{R}_{0+}^D$, *utility function* $u: \mathcal{Y}' \times \mathcal{Y} \rightarrow \mathbb{R}_{0+}$, and *welfare class* $\mathcal{W} \subseteq \mathbb{R}_{0+}^g \rightarrow \mathbb{R}_{0+}$. We say \mathcal{H} is *fair-PAC-learnable* w.r.t. u and \mathcal{W} if there exists an algorithm \mathcal{A} and sample complexity function $m_{\mathcal{W}, \mathcal{H}}$ such that for all

- 1) class parameterizations \mathbf{d} ;
- 2) group counts g ;
- 3) per-group instance distributions $\mathcal{D}_{1:g}$, each over $(\mathcal{X} \times \mathcal{Y})$;
- 4) (weighted) welfare concepts $W(\cdot; \mathbf{w})$ in \mathcal{W} ;
- 5) additive approximation errors $\varepsilon > 0$; &
- 6) failure probabilities $\delta \in (0, 1)$;

it holds that \mathcal{A} can identify a hypothesis $\hat{h} \in \mathcal{H}_{\mathbf{d}}$, i.e., $\hat{h} \leftarrow \mathcal{A}(\mathcal{D}_{1:g}, W, \varepsilon, \delta, \mathbf{d})$, such that

- 1) for each group, $\mathcal{A}(\mathcal{D}_{1:g}, W, \varepsilon, \delta, \mathbf{d})$ draws no more than $m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, W, g, \mathbf{d})$ samples; &
- 2) with probability at least $1 - \delta$ (over randomness of \mathcal{A} and sampling), \hat{h} obeys

$$W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ \hat{h}]; \mathbf{w}]\right) \geq \sup_{h^* \in \mathcal{H}_{\mathbf{d}}} W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbb{E}[u \circ h^*]; \mathbf{w}]\right) - \varepsilon.$$

Furthermore, if $m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, W, g, \mathbf{d})$ can be uniformly bounded for any $W(\cdot; \mathbf{w}) \in \mathcal{W}$, then we say that \mathcal{H} is *uniformly PAC learnable* over \mathcal{W} w.r.t. u .

With trivial changes to convert the maximization objective to a minimization objective, this definition can also be applied to loss functions and classes of malfare functions. In particular, this definition generalizes the FPAC concept given by Cousins [2021], which was specified for the class of all malfare functions satisfying a set of axioms corresponding to $p \geq 1$ weighted power-means. We also relax the definition to allow the sample complexity function to depend on the (weighted) welfare function $W(\cdot; \mathbf{w}) \in \mathcal{W}$, but our concept of *uniform FPAC-learnability* strictly generalizes that of Cousins [2021].

Theorem 4.4 (Characterizing FPAC Learnability). Suppose some weighted power-mean welfare function $W_p(\cdot; \mathbf{w})$, utility function u with range r , and hypothesis class \mathcal{H} with sample complexity function $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d}) \in \text{Poly}\left(\frac{1}{\varepsilon}, \log \frac{1}{\delta}, r, \mathbf{d}\right)$. We then bound the sample complexity $m \doteq m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, W, g, \mathbf{d})$ of FPAC learning \mathcal{H} w.r.t. welfare class $\mathcal{W} \doteq \{W(\cdot; \mathbf{w})\}$ as

- 1) $m \leq m_{\mathcal{H}}\left(\sqrt{\frac{\varepsilon}{2\lambda}}, \frac{\delta}{g}, r, \mathbf{d}\right) \in \text{Poly}\left(\sqrt{\lambda}, \frac{1}{\sqrt{\varepsilon}}, \log \frac{1}{\delta}, \log g, r, \mathbf{d}\right)$;
- 2) $p \in (0, 1] \Rightarrow m \in \text{Poly}\left(\sqrt[p]{p}, \frac{1}{p\sqrt{\varepsilon}}, \frac{1}{p\sqrt{\varepsilon}}, \log \frac{1}{\delta}, \log g, \mathbf{d}\right)$;
- 3) $p = 0 \Rightarrow m \in \text{Poly}\left(\frac{1}{\sqrt[p]{\lambda}}, \frac{1}{\sqrt[p]{\varepsilon}}, \log \frac{1}{\delta}, \log g, \mathbf{d}\right)$;

- 4) $p < 0 \Rightarrow m \in \text{Poly}\left(\frac{1}{\varepsilon}, \frac{1}{\sqrt{|p|}w_{\min}}, \log \frac{1}{\delta}, \log g, r, \mathbf{d}\right)$; &
 5) for any $c \in (0, 1)$, if $|p| \geq c$ and group weights obey the *nonnegligibility condition* $w_{\min} \geq \frac{c}{g}$, then $m \in \text{Poly}^{\frac{1}{c}}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}, g, \log \frac{1}{\delta}, r, \mathbf{d}\right)$.

Observe that the subfamilies of power-mean welfare functions considered in items 2–4 are induced by specific axiomatic choices. In particular, item 2 follows from SM0 (axiom 12), item 3 follows from either IIA or from 0B (axiom 15) and SM0, and item 4 follows from either Lipschitz continuity and SWTP (axiom 13; to prevent $p = 1$) or from FC (axiom 16) — see figure 3 for visual explication. Similarly, item 5 follows by assuming *nonnegligibility* of weights, which holds, e.g., for *unweighted* aggregator functions, and bounding p away from 0, which may be accomplished in a variety of ways. We thus conclude that the axiomatic choices made to restrict the space of welfare functions directly impact their FPAC-learnability.

On Uniform and Polynomial FPAC-Learnability The bounds of theorem 4.4 items 1–4 imply FPAC learnability, but not uniform FPAC learnability, and exponential dependencies on α , $\frac{1}{w_{\min}}$, or $\frac{1}{|p|}$ do appear. It is only in item 5, for any constant c , that the class is uniformly-FPAC-learnable. In contrast, using only Lipschitz continuity, it is straightforward to show that the entire class of *all fair malfare functions*, i.e., any $\Lambda_p(\cdot; \mathbf{w})$ for $p \geq 1$, for which $\alpha = 1$ by lemma 3.12 item 1, is uniformly FPAC learnable.

In general, if α is bounded away from 0, and $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d})$ is polynomial, then the uniform sample complexity of FPAC learning is also polynomial in all parameters, and thus FPAC learning in some sense *preserves polynomial learnability*. We thus conclude that FPAC learning with malfare concepts is easier than FPAC learning with welfare concepts, however under appropriate axiomatically-motivated conditions, the gap in sample complexity between the two settings is polynomial.

4.4 Sample-Efficient Learning and Estimation with Progressive Sampling

Our study of FPAC learnability is based on the Lipschitz and Hölder continuity of power means, which only coarsely describe their behavior, yielding *worst-case* sample complexity bounds. Such bounds may be improved if *a priori* knowledge regarding relative (dis)utility values is available; for example, under near-equality, power-mean functions behave roughly linearly, and may thus be much easier to estimate or optimize. We now briefly voyage into the world of *progressive sampling* to show that, even without such *a priori* knowledge, efficient learning algorithms can adapt their sample consumption to the *inherent difficulty* of the specific task at hand, which may be substantially less than the *worst-case* sample complexity bounds of theorem 4.4 would suggest.

Without considering the delicate intricacies of probabilistic reasoning, one might naïvely assume they could iteratively draw one sample per group, and terminate when the welfare objective is uniformly estimated or approximately optimized. Unfortunately, this quickly runs into statistical errors via the *multiple comparisons problem*, as the sampling process is inherently probabilistic. Efficient progressive sampling methods take this basic idea and account for these issues, but rather than incrementing the sample size at each step, they instead increase the sample size *geometrically*. Such methods have had great impact in myriad settings, including statistical data science [Riondato and Upfal, 2015, 2018, Cousins et al., 2020, 2023d], where estimators query *a single distribution*, empirical game theoretic analysis [Viqueira et al., 2020, 2021, Cousins et al., 2022b, Mishra et al., 2022, Cousins et al., 2023a], where estimators query a noisy utility oracle at *strategy profiles* to bound various game-theoretic properties, and fair machine learning [Cousins, 2022], wherein fair objectives on model classes are estimated and optimized by sampling from *group-specific distributions*.

At a glance, assume our sample complexity bounds scale as $\Theta(\log \frac{1}{\delta})$. Consequently, a schedule of length T with *uniformly allocated* δ (i.e., one that considers up to T sample sizes, and takes probability $1 - \frac{\delta}{T}$ tail bounds at each) can overshoot the sufficient sample size by a constant factor (due to geometric spacing), and furthermore would need a factor $\mathcal{O}(\log T)$ excess samples to correct for the multiple comparisons problem. However, aside from these factors, it is otherwise as sharp as knowing the (task-specific) minimum sufficient sample size *a priori*. Therefore, in cases where Hölder analysis only loosely bounds sample complexity, progressive sampling methods can still adaptively consume about as many samples as are actually required for the task at hand.

The question remains, “How long must the sampling schedule be?” In other words, “How large must T be to guarantee a sufficient sample size is reached?” Here our sample complexity bounds prove invaluable: a geometric sampling schedule must have length logarithmic in the ratio of maximum to minimum sufficient sample sizes, both of which are $\Theta(\log \frac{Tg}{\delta})$, thus solving for a minimal sufficient T is straightforward, e.g., with Hoeffding’s inequality, a doubling schedule admits

$$T = \left\lceil \log_2 \frac{\left[\frac{1}{2} \left(\frac{2r\lambda}{\varepsilon}\right)^{\frac{2}{\alpha}} \ln \frac{2|\mathcal{H}|Tg}{\delta}\right]}{\left[\frac{1}{2} \ln \frac{2|\mathcal{H}|Tg}{\delta}\right]} \right\rceil \in \Theta\left(\frac{1}{\alpha} \log \frac{r\lambda}{\varepsilon}\right)$$

via theorem 4.4 item 1. In general, this progressive sampling strategy induces an *overhead cost factor* of

$$\mathcal{O} \log \log T \subseteq \mathcal{O} \log \log \frac{m_{\mathcal{W}, \mathcal{H}}\left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}, \frac{\delta}{T}, g, r, \mathbf{d}\right)}{\log \frac{Tg}{\delta}} \subset \log \text{Poly}\left(\frac{1}{|p|}, \frac{1}{w_{\min}}, \log \text{Poly}\left(\frac{1}{\varepsilon}, g, r, \mathbf{d}\right)\right) \quad (6)$$

relative to the (unknown) task-specific sufficient sample size. Note that in (6), even terms *exponential* in $\frac{1}{|p|}$ and $\frac{1}{w_{\min}}$ in

the FPAC sample complexity bound become *logarithmic*, due to the double-logarithm. Thus while welfare functions may be *inherently difficult* to estimate, the *statistical overhead* of progressive sampling, as compared to drawing a task-specific sufficient sample, is quite negligible.

5 Experiments

To demonstrate practical relevance, we present a synthetic experiment in a *welfare maximization 1-armed bandit* setting, and study the sample complexity and sensitivity of welfare estimates to various parameters. In particular, we assume each pull of the bandit arm gives a single *utility sample* for each group, and from *empirical mean utilities* $\hat{\mathbf{u}}$, we wish to estimate the welfare $W(\mathbf{u}; \mathbf{w})$ as $W(\hat{\mathbf{u}}; \mathbf{w})$. Note that this is a key step towards regret-optimally selecting among k arms.

We assume groups $\{\mathbf{g}_1, \mathbf{g}_2\}$, where \mathbf{g}_1 is the majority group and \mathbf{g}_2 is the minority group, and the weight and/or expected utility of \mathbf{g}_1 is no less than that of \mathbf{g}_2 . In this experiment, utility samples are $\text{UNIFORM}(\mathbf{u}_i - \frac{1}{2}, \mathbf{u}_i + \frac{1}{2})$ i.i.d. random variables (mean \mathbf{u}_i) for each group i . Figure 4 varies the key parameters of welfare p , minority mean \mathbf{u}_2 , minority weight w_2 , and sample size m , in order to study welfare estimation around the particularly challenging $p \approx 0$ and $w_{\min} \approx 0$ domains. Here empirical utilities $\hat{\mathbf{u}}_i$ have variance $\frac{1}{12m}$, but similar results are shown for *Bernoulli* and *beta* noise models with more complicated variance structure in appendix B. We present the true welfare and approximate Gaussian $\pm 1\sigma$ (68.27%) confidence intervals on empirical welfare, i.e., $W_p(\mathbf{u}; \mathbf{w}) \pm \lambda(\frac{1}{12m})^{\frac{\alpha}{2}}$, where λ, α are as in lemmata 3.12 & 3.13. Using 5000 trials over sampled utilities, we also plot average empirical welfare, and a 68.27% empirical confidence band on empirical welfare, i.e., $W_p(\hat{\mathbf{u}}; \mathbf{w})$.

Note that our theoretical Gaussian confidence intervals should ideally contain *at least* 68.27% of the empirical welfare samples (i.e., the empirical band), and indeed this is the case, despite the error due to sampling and the Gaussian approximation. We don't necessarily expect the empirical 68.27% confidence intervals to contain the true welfare, although this does always occur in these experiments. Moreover, empirical confidence bands may be substantially smaller than Gaussian confidence bands, as they are estimated via sampling, and thus not susceptible to the looseness of Lipschitz and Hölder continuity bounds. We see exactly this, and the gap between the empirical and theoretical confidence bounds varies with the parameters.

Figure 4a shows the impact of changing p on the welfare. Observe that $W_p(\cdot; \mathbf{w})$ is monotonic in p , and it is thus no surprise that all measures of welfare are increasing in p . The interesting portion of the experiment is that both the 68.27% approximate Gaussian and empirical confidence intervals are very wide for $p \approx 0$, and narrow as $|p|$ increases. This concurs with the theory of section 3.4, as despite the small

variances of per-group utility estimates, $W_p(\cdot; \mathbf{w})$ for $p \approx 0$ remains difficult to estimate, due to high sensitivity to near-0 utility values. Observe that the Hölder and Lipschitz analyzes are not sharp, particularly around $p = 0$, because $\mathbf{u}_2 > 0$, and indeed the approximate Gaussian confidence bands contain the empirical percentiles.⁵ Figure 4b then varies the minority utility \mathbf{u}_2 , and we find that as $\mathbf{u}_2 \rightarrow 0$, empirical confidence intervals sharply diverge, due to high sensitivity to minimum utility, i.e., $\hat{\mathbf{u}}_2 \approx 0$. In figure 4c, we vary the weight of the minority group w_2 , and find extremely wide confidence intervals as $w_2 \rightarrow 0$, since $W_0(\hat{\mathbf{u}}; \mathbf{w})$ is very sensitive to $\hat{\mathbf{u}}_2 \approx 0$ when $w_2 \approx 0$, but as $w_2 \rightarrow \frac{1}{2}$, the estimate of welfare becomes much more stable, as higher w_2 means smaller overall welfare, but less sensitivity to $\hat{\mathbf{u}}_2$.

Figures 4a–4c all show the pathologically large estimation error of welfare functions, which when left unchecked causes models to *overfit* to disadvantaged groups, then exhibit bias against them when applied *ex vitro*. In figure 4d, we study a mitigation to this problem, by observing the impact of sample size m on $W_0(\hat{\mathbf{u}}; \mathbf{w})$. Note that by theorem 4.2 the sample complexity of $W_0(\cdot; \mathbf{w})$ estimation is $\mathcal{O}(\frac{1}{\epsilon})^{\frac{2}{w_2}}$, whereas for $W_1(\cdot; \mathbf{w})$ it is only $\mathcal{O}(\frac{w_2}{\epsilon})^2$, and this asymptotically larger sample complexity is visually manifest as slower convergence rates for both the Gaussian and empirical confidence intervals. In all cases, we conclude that, as the theory suggests, an understanding of the continuity properties of power-mean functions is crucial to understanding the sample complexity and estimation error of practical welfare objectives.

6 Conclusion

We show an alternative axiomatic basis for fair aggregator functions, which we argue is simpler, weaker, and more fundamental than prior art. We also draw interdisciplinary connections to moral philosophy and econometric theory to establish stronger axioms, which intuitively guide modellers on fair objective selection, and theoretically distinguish between natural classes of welfare functions. In particular, our (strict) *weak transfer principle*, *zero barrier*, and *finite ceiling* axioms strengthen arguments for prioritarian (i.e., more egalitarian than utilitarian) fairness concepts by assuming less and/or concluding more, and our axioms handling group weights \mathbf{w} simplify existing theory.

We then perform a detailed analysis of the Lipschitz and Hölder continuity of classes of power-mean welfare functions that satisfy our axioms. In particular, we find that our extended axioms naturally partition the class of power-mean functions into classes, each of which share Lipschitz or Hölder continuity properties, which is visually depicted in

⁵Still, the Lipschitz continuity bounds used for sufficiently negative $p < 0$ closely follow the Lipschitz constants plotted in figure 2. Thus while Lipschitz continuity describes the behavior of power means reasonably well here, our Hölder continuity analysis is needed to understand behavior about $p \approx 0$ and for $p > 0$.

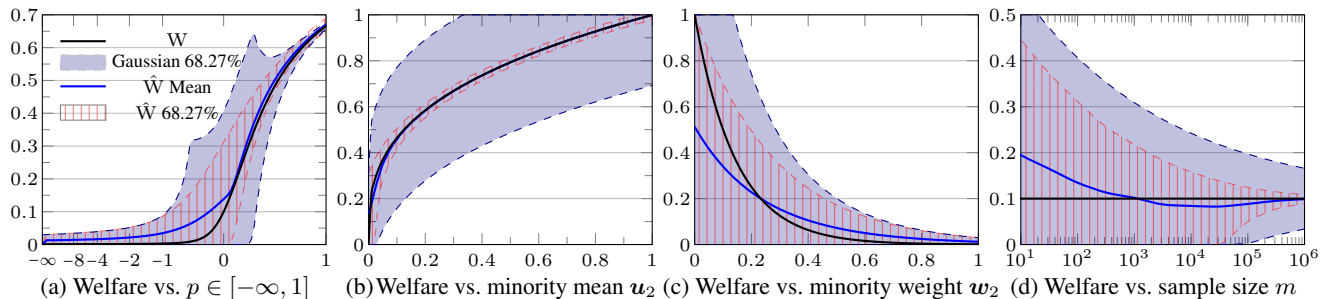


Figure 4: Estimating the Welfare of a 1-Armed Bandit with Uniform Noise. Each plot studies the response of welfare to one parameter, and the remaining parameters are selected from $p = 0$, $\mathbf{u} = \langle 0.999, 0.001 \rangle$, $\mathbf{w} = \langle \frac{2}{3}, \frac{1}{3} \rangle$, and $m = 100$. All axes are linear, except 4a, which plots $p \in [-\infty, 1]$ by transforming $x = \frac{1}{\pi} \arctan(1 - p)$, and 4d, which is logarithmic in x .

figure 3. We also argue that Lipschitz continuity itself may be viewed as a form of *stability* axiom, with consequent resistance to *minority rule*. We follow with a discussion of applications in privacy, algorithmic stability, adversarial robustness, and strategy proofness, finding that Lipschitz and/or Hölder continuity of welfare are often sufficient to show these properties, and we later experimentally study the relationship between choice of axioms, welfare function, data distributions, and the difficulty of estimation.

Finally, we generalize the concept of fair-PAC learning to arbitrary families of welfare functions. We then show conditions under which fair-PAC learning welfare objectives has polynomial sample complexity, and is nearly as efficient as fair-PAC learning malfare objectives, improving the state of the art in utility-based and econometric learning settings. Moreover, prior work handles only the continuity and sample-complexity analysis of the $p \geq 1$ case (malfare), and we show that while the $p < 1$ case is more challenging, the difficulty of learning actually increases smoothly as $p \rightarrow 0$ from both directions, yielding intuitive, mathematically rigorous, and practically actionable understanding of learning and estimation problems over the entire power-mean spectrum. Furthermore, specific *axiomatic choices* regarding the class of welfare functions specify *discrete classes* with interpretable properties and desirable fair-PAC learning guarantees, thus establishing a hierarchy of fair learning settings.

We hope also that these results will be generalized to related settings. For instance, similar analysis is clearly beneficial in *regret minimization* (over groups), also termed “*multi-group agnostic PAC learning*” [Blum and Lykouris, 2020, Rothblum and Yona, 2021, Cousins, 2022], wherein the task is to minimize malfare of *differences between* utility or risk of each group’s preferred model and some compromise model. We envision similar applications in *improvement maximization*, wherein we seek to maximize the welfare of *differences between* utility or risk of the learned compromise model and some fixed reference model, which addresses issues raised by Thomas et al. [2019], Estornell et al. [2023] as to how fair intervention models may be perceived as unfair by groups that preferred the reference model over the fair intervention.

Acknowledgments

We thank Yair Zick for insightful discussion regarding esoteric axioms for aggregator functions, and Indra Kumar, Justin Payan, and Vignesh Viswanathan for helpful feedback on drafts of this work. The authors also wish to acknowledge the support of the Center for Data Science at the University of Massachusetts Amherst, where part of this work was performed under a postdoctoral fellowship.

References

- Jacob D Abernethy, Pranjal Awasthi, Matthäus Kleindessner, Jamie Morgenstern, Chris Russell, and Jie Zhang. Active sampling for min-max fairness. In *International Conference on Machine Learning*, volume 162, 2022.
- Martin Anthony and Peter L Bartlett. *Neural network learning: Theoretical foundations*. Cambridge University Press, second edition, 2009.
- Richard J Arneson. Luck egalitarianism and prioritarianism. *Ethics*, 110(2):339–349, 2000.
- Kenneth J Arrow, Hollis B Chenery, Bagicha S Minhas, and Robert M Solow. Capital-labor substitution and economic efficiency. *The review of Economics and Statistics*, pages 225–250, 1961.
- Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.
- Siddharth Barman, Umang Bhaskar, Anand Krishna, and Ranjani G Sundaram. Tight approximation algorithms for p-mean welfare under subadditive valuations. *Leibniz International Proceedings in Informatics, LIPIcs*, 173, 2020.
- Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural net-

- works. *Advances in neural information processing systems*, 30, 2017.
- Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Guha Thakurta. Private stochastic convex optimization with optimal rates. *Advances in neural information processing systems*, 32, 2019.
- Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic gradient descent on non-smooth convex losses. *Advances in Neural Information Processing Systems*, 33:4381–4391, 2020.
- Gleb Beliakov, Tomasa Calvo, and Simon James. Some results on Lipschitz quasi-arithmetic means. In *European Society for Fuzzy Logic and Technology Conference*, pages 1370–1375, 2009.
- George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.
- Jeremy Bentham. An introduction to the principles of morals and legislation. *University of London: the Athlone Press*, 1789.
- Avrim Blum and Thodoris Lykouris. Advancing subgroup fairness via sleeping experts. In *Innovations in Theoretical Computer Science Conference*, volume 11, 2020.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.
- Jacqueline G Cavazos, P Jonathon Phillips, Carlos D Castillo, and Alice J O’Toole. Accuracy comparison across face recognition algorithms: Where are we on measuring race bias? *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2020.
- Irene Chen, Fredrik D Johansson, and David Sontag. Why is my classifier discriminatory? *Advances in Neural Information Processing Systems*, 31, 2018a.
- Yiling Chen, Chara Podimata, Ariel D Procaccia, and Nisarg Shah. Strategyproof linear regression in high dimensions. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 9–26, 2018b.
- Cynthia M Cook, John J Howard, Yevgeniy B Sirotnin, Jerry L Tipton, and Arun R Vemury. Demographic effects in facial recognition and their dependence on image acquisition: An evaluation of eleven commercial systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 1(1):32–41, 2019.
- Corinna Cortes, Mehryar Mohri, Javier Gonzalez, and Dmitry Storcheus. Agnostic learning with multiple objectives. *Advances in Neural Information Processing Systems*, 33, 2020.
- Cyrus Cousins. An axiomatic theory of provably-fair welfare-centric machine learning. In *Advances in Neural Information Processing Systems*, 2021.
- Cyrus Cousins. Uncertainty and the social planner’s problem: Why sample complexity matters. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 2022.
- Cyrus Cousins and Matteo Riondato. Sharp uniform convergence bounds through empirical centralization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Cyrus Cousins, Shahrzad Haddadan, and Eli Upfal. Making mean-estimation more efficient using an MCMC trace variance approach: DynaMITE. *arXiv:2011.11129*, 2020.
- Cyrus Cousins, Kavosh Asadi, and Michael L. Littman. Fair E³: Efficient welfare-centric fair reinforcement learning. In *5th Multidisciplinary Conference on Reinforcement Learning and Decision Making (RLDM)*, 2022a.
- Cyrus Cousins, Bhaskar Mishra, Enrique Areyan Viqueira, and Amy Greenwald. Computational and data requirements for learning generic properties of simulation-based games. *arXiv:2208.06400*, 2022b.
- Cyrus Cousins, Bhaskar Mishra, Enrique Areyan Viqueira, and Amy Greenwald. Learning properties in simulation-based games. In *Proceedings of the 22nd International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, 2023a.
- Cyrus Cousins, Vignesh Viswanathan, and Yair Zick. Dividing good and better items among agents with submodular valuations. *arXiv:2302.03087*, 2023b.
- Cyrus Cousins, Vignesh Viswanathan, and Yair Zick. The good, the bad and the submodular: Fairly allocating mixed manna under order-neutral submodular preferences. *arXiv:2307.12516*, 2023c.
- Cyrus Cousins, Chloe Wohlgemuth, and Matteo Riondato. BAVarian: Betweenness centrality approximation with variance-aware Rademacher averages. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 17(6):1–47, March 2023d. ISSN 1556-4681. doi: 10.1145/3577021.
- Hugh Dalton. The measurement of the inequality of incomes. *The Economic Journal*, 30(119):348–361, 1920.
- Gerard Debreu. Topological methods in cardinal utility theory. *Cowles Foundation Discussion Papers*, 76, 1959.
- Emily Diana, Wesley Gill, Michael Kearns, Krishnaram Kenthapadi, and Aaron Roth. Minimax group fairness: Algorithms and experiments. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 66–76, 2021.
- Evan Dong and Cyrus Cousins. Decentering imputation: Fair learning at the margins of demographics. In *Queer in AI Workshop @ ICML*, 2022.

- 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*, pages 214–226, 2012.
- Andrew Estornell, Sanmay Das, Brendan Juba, and Yevgeniy Vorobeychik. Popularizing fairness: Group fairness and individual welfare. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 7485–7493, 2023.
- Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. The (im)possibility of fairness: Different value systems require different mechanisms for fair decision making. *Communications of the ACM*, 64(4):136–143, 2021.
- William M Gorman. The structure of utility functions. *The Review of Economic Studies*, 35(4):367–390, 1968.
- Hoda Heidari, Claudio Ferrari, Krishna Gummadi, and Andreas Krause. Fairness behind a veil of ignorance: A welfare analysis for automated decision making. In *Advances in Neural Information Processing Systems*, pages 1265–1276, 2018.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- Lily Hu and Yiling Chen. Fair classification and social welfare. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 535–545, 2020.
- Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR, 2018.
- Thomas M Hurka. Average utilitarianisms. *Analysis*, 42(2):65–69, 1982.
- Mamoru Kaneko and Kenjiro Nakamura. The Nash social welfare function. *Econometrica: Journal of the Econometric Society*, pages 423–435, 1979.
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. In *8th Innovations in Theoretical Computer Science Conference (ITCS)*, volume 67, page 43. Schloß Dagstuhl–Leibniz-Zentrum für Informatik, 2017.
- Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. Fairness without demographics through adversarially reweighted learning. *Advances in neural information processing systems*, 33:728–740, 2020.
- Yunwen Lei and Yiming Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *International Conference on Machine Learning*, pages 5809–5819. PMLR, 2020.
- Natalia Martinez, Martin Bertran, and Guillermo Sapiro. Minimax Pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, pages 6755–6764. PMLR, 2020.
- Daniel McFadden. Constant elasticity of substitution production functions. *The Review of Economic Studies*, 30(2):73–83, 1963.
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35, 2021.
- Reshef Meir and Jeffrey S Rosenschein. Strategyproof classification. *ACM SIGecom Exchanges*, 10(3):21–25, 2011.
- John Stuart Mill. *Utilitarianism*. Parker, Son, and Bourn, London, 1863.
- Bhaskar Mishra, Cyrus Cousins, and Amy Greenwald. Regret pruning for learning equilibria in simulation-based games. *arXiv:2211.16670*, 2022.
- Michael Mitzenmacher and Eli Upfal. *Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis*. Cambridge university press, second edition, 2017.
- Hervé Moulin. *Fair division and collective welfare*. MIT Press, 2004.
- Robert Nozick. *Anarchy, state, and utopia*. Basic Books, 1974.
- Yonatan Oren, Shiori Sagawa, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust language modeling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4227–4237, 2019.
- D Parfit. Equality and priority. *Ratio (Oxford)*, 10(3):202–221, 1997.
- Neel Patel, Reza Shokri, and Yair Zick. Model explanations with differential privacy. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1895–1904, 2022.
- Arthur Cecil Pigou. *Wealth and welfare*. Macmillan and Company, limited, 1912.
- Ariela D Procaccia. Towards a theory of incentives in machine learning. *ACM SIGecom Exchanges*, 7(2):1–5, 2008.
- John Rawls. *A theory of justice*. Harvard University Press, 1971.
- John Rawls. *Justice as fairness: A restatement*. Harvard University Press, 2001.
- Matteo Riondato and Eli Upfal. Mining frequent itemsets through progressive sampling with Rademacher averages. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1005–1014. ACM, 2015.

- Matteo Riondato and Eli Upfal. ABRA: Approximating betweenness centrality in static and dynamic graphs with Rademacher averages. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 12(5):1–38, 2018.
- Kevin WS Roberts. Interpersonal comparability and social choice theory. *The Review of Economic Studies*, pages 421–439, 1980.
- Alvin E Roth et al. Independence of irrelevant alternatives, and solutions to Nash’s bargaining problem. *Journal of Economic Theory*, 16(2):247–251, 1977.
- Guy N Rothblum and Gal Yona. Multi-group agnostic PAC learnability. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18–24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 9107–9115. PMLR, 2021.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2019.
- Mark Schneider and Jonathan W Leland. Salience and social choice. *Experimental Economics*, 24(4):1215–1241, 2021.
- Amartya Sen. On weights and measures: Informational constraints in social welfare analysis. *Econometrica: Journal of the Econometric Society*, pages 1539–1572, 1977.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge University Press, 2014.
- Shubhanshu Shekhar, Greg Fields, Mohammad Ghavamzadeh, and Tara Javidi. Adaptive sampling for minimax fair classification. *Advances in Neural Information Processing Systems*, 34, 2021.
- Umer Siddique, Paul Weng, and Matthieu Zimmer. Learning fair policies in multi-objective (deep) reinforcement learning with average and discounted rewards. In *International Conference on Machine Learning*, pages 8905–8915. PMLR, 2020.
- Philip S Thomas, Bruno Castro da Silva, Andrew G Barto, Stephen Giguere, Yuriy Brun, and Emma Brunskill. Preventing undesirable behavior of intelligent machines. *Science*, 366(6468):999–1004, 2019.
- Enrique Areyan Viqueira, Cyrus Cousins, and Amy Greenwald. Improved algorithms for learning equilibria in simulation-based games. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, pages 79–87, 2020.
- Enrique Areyan Viqueira, Cyrus Cousins, and Amy Greenwald. Learning competitive equilibria in noisy combinatorial markets. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, 2021.
- Vignesh Viswanathan and Yair Zick. A general framework for fair allocation under matroid rank valuations. In *Proceedings of the 24th ACM Conference on Economics and Computation*, pages 1129–1152, 2023.
- Puyu Wang, Yunwen Lei, Yiming Ying, and Hai Zhang. Differentially private SGD with non-smooth losses. *Applied and Computational Harmonic Analysis*, 56:306–336, 2022.
- Gal Yona and Guy Rothblum. Probably approximately metric-fair learning. In *International Conference on Machine Learning*, pages 5680–5688. PMLR, 2018.
- Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International conference on machine learning*, pages 325–333. PMLR, 2013.

Appendices

We derive every result stated in the main body in appendix A, and provide supplementary experiments in appendix B.

A Proof Compendium

We now derive all lemmas, theorems, corollaries, and other results stated in the paper body. This appendix is broken into two subappendices, the first (appendix A.1) shows the properties and relationships between the various axioms of section 3, and the second (appendix A.2) shows all results in section 4 related to sample complexity and FPAC learning.

A.1 Properties of Axiomatic Aggregator Functions

Before delving into the proofs of all results given in the paper body, we first state a standard result that shows the power-mean $M_p(\mathbf{u}; \mathbf{w})$ is monotonically increasing in p , usually referred to as the *power-mean inequality*.

Lemma A.1 (Power-Mean Inequality). Suppose $-\infty \leq p \leq q \leq \infty$. Then for any $\mathbf{u} \in \mathbb{R}_{0+}^g$, $\mathbf{w} \in \Delta_g$, it holds

$$M_p(\mathbf{u}; \mathbf{w}) \leq M_q(\mathbf{u}; \mathbf{w}) .$$

We now show lemma 3.2.

Lemma 3.2 (Equivalence of Weighted Axioms). Consider some aggregator function $M(\cdot; \mathbf{w})$. It always holds that WS (axiom 2) \wedge WD (axiom 3) \Leftrightarrow WA (axiom 9).

Proof. Recall the axioms in question:

Weighted Symmetry (WS): For all permutations π over \mathcal{G} , it holds that $M(\mathbf{u}; \mathbf{w}) = M(\pi(\mathbf{u}); \pi(\mathbf{w}))$.

Weighted Decomposability (WD): Suppose $\alpha \in (0, 1)$. Then $M(\mathbf{u}; \mathbf{w}) = M(\langle \mathbf{u}_1, \mathbf{u}_1, \mathbf{u}_2, \dots \rangle; \langle \alpha \mathbf{w}_1, (1 - \alpha) \mathbf{w}_1, \mathbf{w}_2, \dots \rangle)$.

Weighted Additivity (WA): Suppose $g' \in \mathbb{Z}_+ \cup \{\infty\}$, $\mathbf{u}' \in \mathbb{R}_{0+}^{g'}$, and weights vector $\mathbf{w}' \in \Delta_{g'}$ such that for all $u \in \mathbb{R}_{0+}$, it holds that $\sum_{i \in \mathcal{G}} \mathbf{w}_i \mathbb{1}_u(\mathbf{u}_i) = \sum_{i \in \mathcal{G}'} \mathbf{w}'_i \mathbb{1}_u(\mathbf{u}'_i)$. Then $M(\mathbf{u}; \mathbf{w}) = M(\mathbf{u}'; \mathbf{w}')$.

We first show the reverse direction, i.e., WS (axiom 2) \wedge WD (axiom 3) \Leftarrow WA (axiom 9). Note that the WD holds by definition, as the two weight terms $\alpha \mathbf{w}_1$ and $(1 - \alpha) \mathbf{w}_1$ that share sentiment \mathbf{u}_1 in WD are combined within the summation of WA. Now, observe that WS holds by commutativity of summation over countable sets, thus the LHS and RHS summations in the WA definitions both remain invariant under arbitrary permutation.

We now show the forward direction, i.e., WS (axiom 2) \wedge WD (axiom 3) \implies WA (axiom 9). This result is less direct, but observe that together, (WS) and (WD) can be used to consolidate the weights of all $\mathbf{u}_i, \mathbf{u}_j$ s.t. $\mathbf{u}_i = \mathbf{u}_j$. In particular for each unique \mathbf{u}_i , we can produce some *unique minimal reduction* \mathbf{u}^* and \mathbf{w}^* over population \mathcal{G}^* such that

- (1) $\mathbf{u}_1^* < \mathbf{u}_2^* < \mathbf{u}_3^* < \dots$;
- (2) for all group indices $i \in \mathcal{G}^*$, there exists some $j \in \mathcal{G}$ such that $\mathbf{u}_i^* = \mathbf{u}_j$; &
- (3) for all group indices $i \in \mathcal{G}^*$, it holds $\mathbf{w}_i^* = \sum_{j \in \mathcal{G}} \mathbf{w}_j \mathbb{1}_u(\mathbf{u}_j)$ for some $u \in \mathbb{R}_{0+}$.

Now, observe that *exactly the same* \mathbf{u}^* and \mathbf{w}^* are produced by repeating this process for \mathbf{u}' and \mathbf{w}' over population \mathcal{G}' , thus we may conclude that for all $u \in \mathbb{R}_{0+}$, it holds that

$$\underbrace{\sum_{i \in \mathcal{G}} \mathbf{w}_i \mathbb{1}_u(\mathbf{u}_i)}_{\text{UNIQUE MINIMAL REDUCTION for } \mathbf{u}, \mathbf{w}} = \underbrace{\sum_{i \in \mathcal{G}^*} \mathbf{w}_i^* \mathbb{1}_u(\mathbf{u}_i^*)}_{\text{UNIQUE MINIMAL REDUCTION for } \mathbf{u}', \mathbf{w}'}} = \sum_{i \in \mathcal{G}'} \mathbf{w}'_i \mathbb{1}_u(\mathbf{u}'_i) .$$

We may thus conclude WA. □

We now show lemma 3.3.

Lemma 3.3 (Transfer Principle Equivalencies). Consider some aggregator function $M(\cdot; \mathbf{w})$. The following relate properties (axioms) that $M(\cdot; \mathbf{w})$ obeys.

- 1) PDTP (11) \implies WTP (8); &

2) Suppose axioms 1–7. Then WTP (8) \implies PDTP (11).

Proof. We first show item 1. Observe that the PDTP (11) \implies WTP (8) follows directly, as the PDTP requires that a broad class of equitable (dis)utility transfers are favorable, whereas the WTP requires only that *there exist* some (dis)utility transfer between two particular groups that is favorable.

The conditional reverse implication of item 2 is a bit more subtle. Suppose axioms 1–7. Then, by theorem 3.5 item 1,⁶ we may conclude $M(\mathbf{u}; \mathbf{w}) = M_p(\mathbf{u}; \mathbf{w})$ for some $p \in \mathbb{R}$. Now, suppose for the sake of argument (shown below) that axiom 8 does not hold for welfare functions with $p > 1$, or for malfare functions with $p < 1$. We may thus conclude $p \leq 1$ in the welfare case, and $p \geq 1$ in the malfare case, in either case for which axiom 11 is known to hold, which completes the result.

The remaining step is to show that axiom 8 *does not hold* if $p > 1$ for welfare functions, or $p < 1$ for malfare functions. First, observe that for $p \neq 0$, the monotonic transformation $pM_p^p(\mathbf{u}; \mathbf{w}) = p \sum_{i=1}^g w_i u_i^p$ of $M_p(\mathbf{u}; \mathbf{w})$ is *convex* for $p > 1$, and *concave* for $0 \neq p < 1$. Now, for any $\mathbf{u} \in \mathbb{R}_{0+}^g$, let $i \doteq \operatorname{argmin}_i u_i$ and $j \doteq \operatorname{argmax}_j u_j$, and suppose some $\varepsilon > 0$ s.t. $u_i + w_j \varepsilon < u_j - w_i \varepsilon$. Any “equitable transfer” of the form $W(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w})$ obeys $W(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w}) > W(\mathbf{u}; \mathbf{w})$ for welfare if $p > 1$. Similarly, an “equitable transfer” of the form $\Lambda(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w})$ obeys $\Lambda(\mathbf{u} + \varepsilon w_j \mathbb{1}_i - \varepsilon w_i \mathbb{1}_j; \mathbf{w}) < \Lambda(\mathbf{u}; \mathbf{w})$ for malfare if $p < 1$. Both cases are apparent from the monotonic transform, as the $\frac{w_i}{w_i}$ and $\frac{w_j}{w_j}$ weighting terms cancel, leaving only transfers along the curvature of the $(\cdot)^p$ power function. In either case, the WTP is violated, thus we may conclude $p \leq 1$ for welfare functions, and $p \geq 1$ for malfare functions. Finally, note that similar logic applies for the case of $p = 0$, instead using a logarithmic monotonic transform. \square

We now show theorem 3.5.

Theorem 3.5 (Aggregator Function Properties). Suppose aggregator function $M(\mathbf{u}; \mathbf{w})$, and assume arbitrary sentiment vector $\mathbf{u} \in \mathbb{R}_{0+}^g$ and weights vector $\mathbf{w} \in \Delta_g$. The following then hold.

- 1) *Power-Mean Factorization*: Axioms 1–7 imply there exists some $p \in \mathbb{R}$ such that $M(\mathbf{u}; \mathbf{w}) = M_p(\mathbf{u}; \mathbf{w})$.
- 2) *Fair Welfare and Malfare*: Axioms 1–8 imply $p \in (-\infty, 1]$ for welfare and $p \in [1, \infty)$ for malfare.

Proof. The key to showing this result is to note that, as mentioned in the text, theorem 2.4 of Cousins [2021] draws the same *conclusion*, but under different *assumptions*. The proof strategy is thus to show that our seemingly weaker assumptions actually imply (in fact, are equivalent to) the assumptions of the aforementioned result. In particular, for item 1, it suffices to conclude axioms 1, 4–7 & 9, and for item 2, we need only additionally conclude axiom 11.

We now show item 1. Observe that we assume axioms 1–7 directly, leaving only axiom 9 (WA), which by lemma 3.2, is implied by axioms 2 & 3. This concludes item 1.

We now show item 2. Observe that after assuming our axioms, we need only show axiom 11 (PDTP), which by lemma 3.3 item 2, is implied by the assumed axioms 1–7, in conjunction with axiom 8, which is also assumed. This concludes item 2. \square

We now show lemma 3.7.

Lemma 3.7 (Consequences of Strong Axioms). Suppose power-mean aggregator function $M_p(\cdot; \mathbf{w})$. The following then hold:

- 1) Strengthening the SM axiom (i.e., 1 \rightarrow 12) implies $p > 0$.
- 2) Strengthening the WTP axiom (i.e., 8 \rightarrow 13) implies $p \neq 1$, thus $p < 1$ for welfare and $p > 1$ for malfare.
- 3) Strict PDTP (i.e., 11 \rightarrow 14) implies $p \neq 1$ and $p \neq \pm\infty$, thus $p \in (-\infty, 1)$ for welfare and $p \in (1, \infty)$ for malfare.

Proof. We first show item 1. First note that the desideratum follows directly from the following claim: “If $\mathbf{u} \neq \mathbf{0}$ and $\min_{i \in G} u_i = 0$, then $(M_p(\mathbf{u}; \mathbf{w}) = 0) \Leftrightarrow (p \leq 0)$.” In particular, here $M(\mathbf{u}; \mathbf{w}) = M(\mathbf{0}; \mathbf{w}) = 0$ for $\mathbf{u} \neq \mathbf{0}$ would violate axiom 12, thus by contraposition, axiom 12 implies $p > 0$. We thus need only show this claim, which follows via analysis of the power mean.

It is straightforward to see that since $\mathbf{u} \neq \mathbf{0}$, it holds $p > 0 \implies M_p(\mathbf{u}; \mathbf{w}) > 0$, thus by the contrapositive, $M_p(\mathbf{u}; \mathbf{w}) = 0 \implies p \leq 0$. To see the converse, first observe that since $\min_i u_i = 0$, it holds $p = 0 \implies M_p(\mathbf{u}; \mathbf{w}) = 0$. The case of

⁶Note that theorem 3.5 makes use of this result to show item 2, but we use only theorem 3.5 item 1 here, thus there is no cyclic dependency.

$p < 0$ appears to be a bit more subtle, however observe that the power mean is monotonically increasing in p , and thus in this case, we observe the sandwich inequality

$$0 = M_{-\infty}(\mathbf{u}; \mathbf{w}) \leq M_p(\mathbf{u}; \mathbf{w}) \leq M_0(\mathbf{u}; \mathbf{w}) = 0 .$$

It thus holds that, for $\mathbf{u} \neq \mathbf{0}$, $p \leq 0 \implies M_p(\mathbf{u}; \mathbf{w}) = 0$. Both directions of the bijection of the claim have now been shown, thus item 1 is complete.

We now show item 2. From theorem 3.5 item 2, we already have $p \leq 1$ for welfare and $p \geq 1$ for malfare, so in both cases, all we need do is show $p \neq 1$, which we now show via the contrapositive. Now observe that for all ε , it holds that $M_1(\mathbf{u} + \mathbb{1}_i \mathbf{w}_j \varepsilon; \mathbf{w}) = M_1(\mathbf{u} - \mathbb{1}_j \mathbf{w}_i \varepsilon; \mathbf{w})$, hence for no choice of ε does the transfer result in a strict improvement to welfare or malfare. We thus conclude $p = 1 \implies \neg(\text{axiom 13})$, hence $(\text{axiom 13}) \implies p \neq 1$, which concludes item 2.

We now show item 3. Similar logic to item 2 precludes the case of $p = 1$ (indeed, observe that this must be so, as SPDTP implies PDTP, by similar reasoning to that found in the proof of lemma 3.3 item 1), this portion of item 3 is a direct corollary of item 2. Now, observe that the egalitarian cases $p \in \pm\infty$ are also inadmissible, essentially because they are only sensitive to the extreme values of \mathbf{u} and thus transfer between any two non-extreme u_i, u_j , i.e., transfer between i, j s.t. $\inf_k u_k < u_i < u_j < \sup_k u_k$, has no impact on the egalitarian power means. We thus conclude that under axiom 14 (SPDTP), it holds that $p \neq \pm\infty$. \square

We now show lemma 3.9.

Lemma 3.9 (Consequences of Extreme Axioms). Suppose as in lemma 3.7. The following then hold.

1) 0B (axiom 15) $\Leftrightarrow p \leq 0$. 2) FC (axiom 16) $\Leftrightarrow p < 0$.

Proof. We first show item 1. We first show that 0B (axiom 15) $\implies p \leq 0$. This is clear by contrapositive, as for any $p > 0$, it holds $M_p((0, 1); \langle \frac{1}{2}, \frac{1}{2} \rangle) = \frac{1}{2^{1/p}} > 0$, thus 0B does not hold.

We now show the converse, i.e., $p \leq 0 \implies$ 0B (axiom 15). In particular, we seek to show that $\lim_{\mathbf{u}_i \rightarrow 0^+} M_0(\mathbf{u}; \mathbf{w}) = 0$. We first address the case of $p < 0$. Observe that

$$\begin{aligned} \lim_{\mathbf{u}_i \rightarrow 0^+} M_p(\mathbf{u}; \mathbf{w}) &= \lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\varepsilon \rightarrow 0^+} \left(\frac{1}{\sum_{i=1}^g \frac{w_i}{(\mathbf{u}_i + \varepsilon)^{-p}}} \right)^{-\frac{1}{p}} && \text{DEFINITION 3.4 (POWER MEAN)} \\ &= \left(\frac{1}{\lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\varepsilon \rightarrow 0^+} \sum_{i=1}^g \frac{w_i}{(\mathbf{u}_i + \varepsilon)^{-p}}} \right)^{-\frac{1}{p}} && \text{LIMIT LAWS} \\ &= \left(\frac{1}{\infty} \right)^{-\frac{1}{p}} = 0 . && \text{LIMIT LAWS} \end{aligned}$$

We now address the case of $p = 0$; in particular, observe that

$$\begin{aligned} \lim_{\mathbf{u}_i \rightarrow 0^+} M_0(\mathbf{u}; \mathbf{w}) &= \lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\rho \rightarrow 0} \lim_{\varepsilon \rightarrow 0^+} M_\rho(\mathbf{u} + \varepsilon \mathbf{1}; \mathbf{w}) && \text{DEFINITION 3.4 (POWER MEAN)} \\ &= \lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\varepsilon \rightarrow 0^+} \prod_{i=1}^g (\mathbf{u}_i + \varepsilon)^{w_i} && \text{GEOMETRIC MEAN LIMIT} \\ &= \lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\varepsilon \rightarrow 0^+} \exp \left(\sum_{i=1}^g w_i \ln(\mathbf{u}_i + \varepsilon) \right) && \text{LOGARITHMIC IDENTITIES} \\ &= \exp \left(\lim_{\mathbf{u}_i \rightarrow 0^+} \lim_{\varepsilon \rightarrow 0^+} \sum_{i=1}^g w_i \ln(\mathbf{u}_i + \varepsilon) \right) && \text{LIMIT LAWS} \\ &= \exp(-\infty) = 0 . && \text{LIMIT LAWS} \end{aligned}$$

We thus have that, for any $p \leq 0$, it holds $\lim_{\mathbf{u}_i \rightarrow 0^+} M_0(\mathbf{u}; \mathbf{w}) = 0$. This completes the converse statement, and thus concludes item 1.

We now show item 2. To see this result, first observe that for any $\mathbf{u} \succ \mathbf{0}$, say, $\mathbf{u} \doteq \mathbf{1}$, it holds that $\lim_{c \rightarrow \infty} M_p(\mathbf{u} + c\mathbf{1}_i; \mathbf{w}) = \infty$ if and only if $p \geq 0$. In particular, this can be seen by observing that, for the forward direction, that $p < 0 \implies M_p(\mathbf{u}; \mathbf{w}) \leq \sqrt[p]{\mathbf{w}_j} \mathbf{u}_j \leq \infty$ for any $j \in \mathcal{G}$, thus $p \geq 0$, and for the reverse direction,

$$M_0(\mathbf{u} + c\mathbf{1}_i; \mathbf{w}) = \lim_{c \rightarrow \infty} (1 + c)^{\mathbf{w}_i} = \infty ,$$

and then that $M_0(\mathbf{u}; \mathbf{w}) \leq M_p(\mathbf{u}; \mathbf{w})$ for any $p \geq 0$, i.e., lemma A.1. Note that the same would hold for all \mathbf{u} if we assumed $p > 0$, but not for \mathbf{u} s.t. $\mathbf{u}_j = 0$ for some $j \neq i$ if $p = 0$, however, we only need the *existence* of a single \mathbf{u} for which the statement holds.

From here, logically, we have that

$$\left(\exists \mathbf{u} \in \mathbb{R}_{0+}^g \quad \text{s.t.} \quad \lim_{c \rightarrow \infty} M_p(\mathbf{u} + c\mathbf{1}_i; \mathbf{w}) = \infty \right) \Leftrightarrow (p \geq 0) .$$

Now, observe that, by contraposition of the bijection, it holds that

$$(p < 0) \Leftrightarrow \underbrace{\left(\forall \mathbf{u} \in \mathbb{R}_{0+}^g : \lim_{c \rightarrow \infty} M_p(\mathbf{u} + c\mathbf{1}_i; \mathbf{w}) < \infty \right)}_{\text{FC AXIOM}} ,$$

and observe that the RHS is, by definition, the FC axiom. □

We now show lemma 3.10.

Lemma 3.10 (Power-Mean Differentiation). Suppose $\mathbf{u}_{\setminus i} \succ \mathbf{0}$, some weights vector $\mathbf{w} \in \Delta_g$, and $p \in \mathbb{R}$. The power mean then differentiates in \mathbf{u}_i as follows.

- 1) If $\mathbf{u}_i > 0$, then $\frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})}$.
- 2) If $p < 0$, then $\lim_{\mathbf{u}_i \rightarrow 0+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = -\sqrt[p]{\frac{1}{\mathbf{w}_i}}$.
- 3) If $p \in [0, 1)$, then $\lim_{\mathbf{u}_i \rightarrow 0+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = \infty$.

Proof. Observe first that item 1 is an elementary application of the chain, power, and summation rules; the only subtlety to this result arises in the remaining cases.

We now show item 2. This case is difficult, as naïve application of item 1 results in an indeterminate $\frac{0}{0}$ form. An experienced practitioner of the calculus of infinitesimals may expect results via L'Hôpital's rule, however in this case, said approach is unwieldy, and a simple limit calculus argument yields the desideratum much more concisely. Observe now that the result follows as

$$\begin{aligned} \lim_{\mathbf{u}_i \rightarrow 0+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) &= \lim_{\mathbf{u}_i \rightarrow 0+} \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})} && \text{ITEM 1} \\ &= \lim_{\mathbf{u}_i \rightarrow 0+} \mathbf{w}_i M_p^{1-p} \left(j \mapsto \frac{\mathbf{u}_j}{\mathbf{u}_i}; \mathbf{w} \right) && \text{MULTIPLICATIVE LINEARITY} \\ &= \mathbf{w}_i \left(\sum_{j=1}^g \mathbf{w}_j \lim_{\mathbf{u}_i \rightarrow 0+} \left(\frac{\mathbf{u}_j}{\mathbf{u}_i} \right)^p \right)^{\frac{1-p}{p}} && \text{LIMIT LAWS} \\ &= \mathbf{w}_i \left(\mathbf{w}_i 1^p + \sum_{j=1, j \neq i}^g \mathbf{w}_j 0^p \right)^{\frac{1-p}{p}} && \text{LIMIT LAWS} \\ &= \mathbf{w}_i^{1 + \frac{1-p}{p}} = \mathbf{w}_i^{\frac{1}{p}} = -\sqrt[p]{\frac{1}{\mathbf{w}_i}} . && \text{ALGEBRA} \end{aligned}$$

We now show item 3. We split this case into two subcases; namely the $p = 0$ and $p > 0$ subcases, essentially because whether $M_p(\mathbf{u}; \mathbf{w}) = 0$ in the limit is of material significance to the proof technique.

We begin with the $p \in (0, 1)$ case. This case is simpler than the case of $p = 0$, as here $M_p(\mathbf{u}; \mathbf{w}) \neq 0$ in the limit, thus the difficulty of resolving the $\frac{0}{0}$ indeterminate form in the limit vanishes entirely. In fact, in this case we have a finite nonzero denominator, and an infinite numerator. In particular, we may observe the result as

$$\begin{aligned}
 \lim_{\mathbf{u}_i \rightarrow 0^+} \frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) &= \lim_{\mathbf{u}_i \rightarrow 0^+} \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})} && \text{ITEM 1} \\
 &= \mathbf{w}_i \frac{\lim_{\mathbf{u}_i \rightarrow 0^+} \mathbf{u}_i^{p-1}}{\lim_{\mathbf{u}_i \rightarrow 0^+} M_p^{p-1}(\mathbf{u}; \mathbf{w})} && \text{LIMIT LAWS} \\
 &= \underbrace{\mathbf{w}_i M_p^{1-p}(\mathbf{u}; \mathbf{w})}_{\text{POSITIVE FINITE}} \left(\lim_{\mathbf{u}_i \rightarrow 0^+} \frac{1}{\mathbf{u}_i} \right)^{1-p} = \infty . && \text{ALGEBRA}
 \end{aligned}$$

We now show the case of $p = 0$. Direct proof is more subtle than for $p \in (0, 1)$, but can be derived via reasoning akin to that of item 2. However, it is much easier to observe that the power mean exhibits continuity in p , and therefore taking $\lim_{p \rightarrow 0}$ via either the $p < 0$ or $p \in (0, 1)$ case yields the desideratum. \square

We now show lemma 3.12.

Lemma 3.12 (Power-Mean Lipschitz Continuity). Suppose $p \in \mathbb{R}$, sentiment vectors $\mathbf{u}, \mathbf{u}' \in \mathbb{R}_{0^+}^g$, and weights vector $\mathbf{w} \in \Delta_g$. The following then hold.

- 1) Suppose $p \geq 1$. Then $M_p(\cdot; \mathbf{w})$ is $\sqrt[p]{\mathbf{w}_{\max^-}} \|\cdot\|_1$, $1 - M_p(\cdot; \mathbf{w})$, and $1 - \|\cdot\|_\infty$ Lipschitz.
- 2) Suppose $p < 0$. Then $M_p(\cdot; \mathbf{w})$ is $\frac{1}{\sqrt[p]{\mathbf{w}_{\min^-}}} \|\cdot\|_\infty$ Lipschitz.

Proof. Items 1 & 2 follow directly from lemma 3.10, and consideration of the curvature and monotonicity of these functions. Briefly put, observe that malfare functions (i.e., $p \geq 1$ weighted power-means) exhibit monotonically-increasing convexity, thus derivatives increase as $\mathbf{u}_i \rightarrow \infty$, whereas welfare functions (i.e., $p \leq 1$ weighted power-means) exhibit monotonically-increasing concavity, thus derivatives increase as $\mathbf{u}_i \rightarrow 0^+$. We now show each result in detail.

We first show item 1. We begin with the $\sqrt[p]{\mathbf{w}_{\max^-}} \|\cdot\|_1$ Lipschitz property. Observe that for any group index $i \in \mathcal{G}$, and any sentiment value $\mathbf{u}_i > 0$, it holds that

$$\frac{\partial}{\partial \mathbf{u}_i} M_p(\mathbf{u}; \mathbf{w}) = \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})} \leq \lim_{\mathbf{u}_i \rightarrow \infty} \frac{\mathbf{w}_i \mathbf{u}_i^{p-1}}{M_p^{p-1}(\mathbf{u}; \mathbf{w})} = \frac{\mathbf{w}_i}{\mathbf{w}_i^{\frac{p-1}{p}}} = \mathbf{w}_i^{1 - \frac{p-1}{p}} = \sqrt[p]{\mathbf{w}_i} .$$

From here, maximizing over group indices yields the $\sqrt[p]{\mathbf{w}_{\max^-}} \|\cdot\|_1$ Lipschitz characterization.

We now show the $1 - M_p(\cdot; \mathbf{w})$ Lipschitz property. Observe that $|M_p(\mathbf{u}; \mathbf{w}) - M_p(\mathbf{u}'; \mathbf{w})| \leq M_p(|\mathbf{u} - \mathbf{u}'|; \mathbf{w})$ follows via the subadditivity of $p \geq 1$ power-mean functions, i.e., they are convex and have the unique zero of $M_p(\mathbf{0}; \mathbf{w}) = 0$.

Finally, to see the $1 - \|\cdot\|_\infty$ Lipschitz property, observe that $M_p(|\mathbf{u} - \mathbf{u}'|; \mathbf{w}) \leq M_\infty(|\mathbf{u} - \mathbf{u}'|; \mathbf{w}) = \|\mathbf{u} - \mathbf{u}'\|_\infty$, which follows from monotonicity of the power mean in p , i.e., lemma A.1.

We now show item 2. A $\frac{1}{\sqrt[p]{\mathbf{w}_{\min^-}}} \|\cdot\|_\infty$ Lipschitz continuity guarantee can easily be seen by maximizing derivatives, via lemma 3.10 item 2 (note that this limit maximizes the derivative, since the welfare function is concave and increasing). It may seem surprising that we could get the same Lipschitz constant for $\|\cdot\|_\infty$, however observe that even taking two values of $\mathbf{u}_i, \mathbf{u}_j$ to 0 simultaneously actually results in *smaller change*, as it is effectively the same as increasing the weight \mathbf{w}_i of a single group, thus the same analysis yields an $\|\cdot\|_\infty$ Lipschitz constant. \square

We now show lemma 3.13.

Lemma 3.13 (Power-Mean Hölder Continuity). Suppose $\mathbf{u} \in [0, r]^g$, group index $i \in \mathcal{G}$, weights vector $\mathbf{w} \in \Delta_g$, and assume where appropriate that $\mathbf{u}_i + \varepsilon \leq r$. The power mean then obeys the following Hölder continuity criteria.

- 1) *Generic Welfare Hölder Condition:* Suppose $p \leq 1$. Then $|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})| \leq r^{1-\mathbf{w}_i} \varepsilon^{\mathbf{w}_i}$, and $M_p(\cdot; \mathbf{w})$ is $r^{1-\mathbf{w}_{\min^-}} \mathbf{w}_{\min^-} \|\cdot\|_\infty$ Hölder continuous.
- 2) *Positive Welfare Hölder Condition:* Suppose $p \in (0, 1]$. Then $|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})| \leq r^{1-p} \frac{\mathbf{w}_i}{p} \varepsilon^p$. Furthermore, $M_p(\cdot; \mathbf{w})$ meets the following Hölder conditions:

- A) $r^{1-p} \frac{w_{\max}}{p} - p - \|\cdot\|_1$;
 B) $r^{1-p} \frac{1}{p} - p - M_1(\|\cdot\|; \mathbf{w})$; &
 C) $r^{1-p} \frac{1}{p} - p - \|\cdot\|_\infty$.

Proof. Due to the complicated and multifaceted nature of this result, we break the proof into several parts. Before showing the first item, we begin with two auxiliary results that will prove useful throughout.

We first note that a generic way to show $\lambda - \alpha - \|\cdot\|$ Hölder continuity w.r.t. \mathbf{u}_i is to show that

$$\sup_{\mathbf{u}, \varepsilon} \frac{|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})|}{\varepsilon^\alpha} \leq \lambda ,$$

and similar techniques can be used to analyze Hölder continuity w.r.t. norms over per-group differences.

We now note that due to its scale-dependence, it is often convenient to show local Hölder continuity, i.e., Hölder continuity over a bounded region. For simplicity, we assume that utility values have range $[0, 1]$, and then extend this analysis to a larger region through the multiplicative linearity axiom. The remainder of the proof assumes WLOG this range, and the below analysis is applied to produce the final (range-dependent) result.

Observe that if $f(x)$ exhibits multiplicative linearity (as do all power means, by the multiplicative linearity axiom), and is $\lambda - \alpha$ Hölder continuous, then for any $r > 0$, it holds that $x \mapsto r f(\frac{x}{r})$ is $r^{1-\alpha} \lambda - \alpha$ Hölder continuous. To see this, first suppose $f(x)$ exhibits multiplicative linearity and is $\lambda - \alpha$ Hölder continuous. Then $g(x) \doteq r f(\frac{x}{r})$ obeys

$$\frac{|g(x) - g(y)|}{|x - y|^\alpha} = \frac{|r f(\frac{x}{r}) - r f(\frac{y}{r})|}{|x - y|^\alpha} = \frac{r |f(\frac{x}{r}) - f(\frac{y}{r})|}{r^\alpha \left| \frac{x}{r} - \frac{y}{r} \right|^\alpha} = r^{1-\alpha} \frac{|f(\frac{x}{r}) - f(\frac{y}{r})|}{\left| \frac{x}{r} - \frac{y}{r} \right|^\alpha} \leq r^{1-\alpha} \lambda .$$

We now show item 1. We first consider the case of $p = 0$, i.e., we analyze the Nash social welfare $M_0(\cdot; \mathbf{w})$. In this case, observe that the most rapid change to $M_0(\mathbf{u}; \mathbf{w})$ occurs as some \mathbf{u}_i approaches zero, and furthermore, the degree of change is maximized when each remaining $\mathbf{u}_j = 1$, i.e., is maximized (this much is clear from concavity). In particular, for each $i \in \mathcal{G}$, taking $\alpha = \mathbf{w}_i$, here we have

$$\begin{aligned} \sup_{\mathbf{u}, \varepsilon} \frac{|M_0(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_0(\mathbf{u}; \mathbf{w})|}{\varepsilon^{\mathbf{w}_i}} &\leq \frac{1}{\varepsilon^{\mathbf{w}_i}} |M_0(\mathbf{0} + \varepsilon \mathbf{1}_i + \mathbf{1}_{\mathcal{G} \setminus \{i\}}; \mathbf{w}) - M_0(\mathbf{0} + \mathbf{1}_{\mathcal{G} \setminus \{i\}}; \mathbf{w})| && \begin{array}{l} \text{CONCAVITY} \\ \text{MONOTONICITY} \end{array} \\ &= \frac{1}{\varepsilon^{\mathbf{w}_i}} \exp(\mathbf{w}_i \ln(\varepsilon) + (1 - \mathbf{w}_i) \ln(1)) && \text{DEFINITION OF } M_0(\cdot; \mathbf{w}) \\ &= \frac{1}{\varepsilon^{\mathbf{w}_i}} \exp(\mathbf{w}_i \ln(\varepsilon)) = 1 , && \text{ALGEBRA} \end{aligned}$$

from which we may conclude $\lambda = 1$. This is enough to bound the Hölder constants for the $\|\cdot\|_1$ norm, however observe that even taking two values of $\mathbf{u}_i, \mathbf{u}_j$ to 0 simultaneously actually results in *slower growth*, as it is effectively the same as increasing the weight \mathbf{w}_i , and the Hölder constants are actually higher for smaller weights values \mathbf{w}_i . We thus conclude that the same bounds hold for the $\|\cdot\|_\infty$ case.

The above completes item 1 for $p = 0$, so we now show that the result holds for all $p \leq 1$. In other words, we show that $p = 0$ is in some sense the “worst case” for small-scale local deviations. To see this, observe that, for any $\varepsilon > 0$, $i \in \mathcal{G}$, it holds that $M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})$ is decreasing as $p \rightarrow 0$, from both the positive and negative sides. We thus conclude that $\lambda - \alpha - \|\cdot\|$ Hölder continuity for $M_0(\mathbf{u}; \mathbf{w})$ implies the same for $M_p(\mathbf{u}; \mathbf{w})$.

We now show item 2. Assume $p \in (0, 1]$. Observe then that

$$\begin{aligned} \sup_{\mathbf{u}, \varepsilon} \frac{|M_p(\mathbf{u} + \varepsilon \mathbf{1}_i; \mathbf{w}) - M_p(\mathbf{u}; \mathbf{w})|}{\varepsilon^p} &= \sup_{\varepsilon \in (0, 1)} \frac{|M_p(\langle \varepsilon, 1 \rangle; \langle \mathbf{w}_i, 1 - \mathbf{w}_i \rangle) - M_p(\langle 0, 1 \rangle; \langle \mathbf{w}_i, 1 - \mathbf{w}_i \rangle)|}{\varepsilon^p} && \begin{array}{l} \text{CONCAVITY} \\ \text{MONOTONICITY} \end{array} \\ &= \sup_{\varepsilon \in (0, 1)} \frac{(\mathbf{w}_i \varepsilon^p + (1 - \mathbf{w}_i))^{\frac{1}{p}} - (1 - \mathbf{w}_i)^{\frac{1}{p}}}{\varepsilon^p} && \text{DEFINITION OF } M_p(\cdot; \mathbf{w}) \\ &\leq \sup_{\varepsilon \in (0, 1)} \frac{\frac{1}{p} \mathbf{w}_i \varepsilon^{p-1} + (1 - \mathbf{w}_i)^{\frac{1}{p}-1} - (1 - \mathbf{w}_i)^{\frac{1}{p}-1}}{\varepsilon^p} && \text{SEE BELOW} \\ &= \frac{\mathbf{w}_i}{p} . && \text{ALGEBRA} \end{aligned}$$

For the step marked SEE BELOW, suppose $a, b \geq 0$ s.t. $a + b \leq 1$. Then for all $c \geq 1$, it holds that $(a + b)^c \leq ca + b^c$. This algebraic manipulation yields the result.

From here, item A follows immediately, and item B follows via a similar argument (i.e., the total weighted deviation, as measured by $M_p(\cdot; \mathbf{w})$, plays the role of ε). Finally, item C follows from item B by noting that, for any $\mathbf{u} \in \mathbb{R}_{0+}^g$ and weights vector $\mathbf{w} \in \Delta_g$, it holds that $\|\mathbf{u}\|_{1, \mathbf{w}} \leq \|\mathbf{u}\|_\infty \leq \|\mathbf{u}\|_1$. \square

A.2 Analysis of Fair-PAC Learning

We now show theorem 4.1.

Theorem 4.1 (Hölder Continuity and Welfare Optimality). Suppose $W(\cdot; \mathbf{w})$ is λ - α - $\|\cdot\|_{\mathbb{W}}$ Hölder continuous w.r.t. some norm $\|\cdot\|_{\mathbb{W}}$, and additive error bounds ε that obey (3). Then

$$\sup_{h \in \mathcal{H}} \left| W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ h]; \mathbf{w}\right) - W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ h]; \mathbf{w}\right) \right| \leq \lambda \|\varepsilon\|_{\mathbb{W}}^\alpha.$$

Consequently, the *empirical welfare maximizer*

$$\hat{h} \doteq \sup_{h \in \mathcal{H}} W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ h]; \mathbf{w}\right)$$

approximates the *true welfare maximizer*

$$h^* \doteq \sup_{h^* \in \mathcal{H}} W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ h^*]; \mathbf{w}\right),$$

in terms of welfare-optimality, as it holds that

$$W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ \hat{h}]; \mathbf{w}\right) \geq W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ h^*]; \mathbf{w}\right) - 2\lambda \|\varepsilon\|_{\mathbb{W}}^\alpha.$$

Proof. The first portion of the result follows directly from the assumption, and the definition of Hölder continuity (definition 3.11).

The next applies a standard technique in learning theory, wherein the first bound is applied twice: once for h^* and once more for \hat{h} , alongside the fact that, by definition \hat{h} realizes the supremum over the empirical welfare. In particular, we have

$$\begin{aligned} W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ \hat{h}]; \mathbf{w}\right) &\geq W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ \hat{h}]; \mathbf{w}\right) - \lambda \|\varepsilon\|_{\mathbb{W}}^\alpha && \text{FIRST PORTION ON } \hat{h} \\ &\geq W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ h^*]; \mathbf{w}\right) - \lambda \|\varepsilon\|_{\mathbb{W}}^\alpha && W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ \hat{h}]; \mathbf{w}\right) \geq W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ h^*]; \mathbf{w}\right) \\ &\geq W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ h^*]; \mathbf{w}\right) - 2\lambda \|\varepsilon\|_{\mathbb{W}}^\alpha. && \text{FIRST PORTION ON } h^* \end{aligned}$$

We now show theorem 4.2. \square

Theorem 4.2 (Welfare Sample Complexity). Suppose sample complexity function $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d})$ for hypothesis class \mathcal{H} , and some welfare function $W(\cdot; \mathbf{w})$ that is λ - α - $\|\cdot\|_\infty$ Hölder continuous. Then the sample complexity function

$$m_{\mathbb{W}, \mathcal{H}}(\varepsilon, \delta, g, r, \mathbf{d}) \leq m_{\mathcal{H}}\left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}, \frac{\delta}{g}, r, \mathbf{d}\right)$$

is sufficient, i.e., for at least this many samples from each of the g groups, (5) holds. Moreover, for this sample size, with probability at least $1 - \delta$, the empirical welfare maximizer is 2ε -optimal.

Proof. This result essentially follows from theorem 4.1 and the definitions of sample complexity and Hölder continuity. By definition, a sample of size at least $m_{\mathcal{H}}\left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}, \delta, r, \mathbf{d}\right)$ ensures a probability $1 - \delta$ bound on the supremum deviation for a single group, and thus by union bound, a sample of size at least $m_{\mathcal{H}}\left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}, \frac{\delta}{g}, r, \mathbf{d}\right)$ ensures a probability $1 - \delta$ on the $\|\cdot\|_\infty$ norm of per-group supremum deviations over all groups, i.e., it shall hold with the above probability that $\|\varepsilon\|_\infty \leq \sqrt[\alpha]{\frac{\varepsilon}{\lambda}}$. Then, applying theorem 4.1 yields

$$\sup_{h \in \mathcal{H}_d} \left| W\left(i \mapsto \mathbb{E}_{\mathcal{D}_i}[\mathbf{u} \circ h]; \mathbf{w}\right) - W\left(i \mapsto \hat{\mathbb{E}}_{\mathcal{D}_i}[\mathbf{u} \circ h]; \mathbf{w}\right) \right| \leq \lambda \|\varepsilon\|_\infty^\alpha \leq \lambda \left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}\right)^\alpha = \varepsilon.$$

As we have ε -estimated $W(\cdot; \mathbf{w})$ with this sample, we may conclude that $m_{\mathbb{W}, \mathcal{H}}(\varepsilon, \delta, g, r, \mathbf{d}) \leq m_{\mathcal{H}}\left(\sqrt[\alpha]{\frac{\varepsilon}{\lambda}}, \delta, r, \mathbf{d}\right)$. Finally, the statement about approximate optimality of the empirical welfare maximizer follows from the second portion of theorem 4.1. \square

We now show theorem 4.4.

Theorem 4.4 (Characterizing FPAC Learnability). Suppose some weighted power-mean welfare function $W_p(\cdot; \mathbf{w})$, utility function u with range r , and hypothesis class \mathcal{H} with sample complexity function $m_{\mathcal{H}}(\varepsilon, \delta, r, \mathbf{d}) \in \text{Poly}(\frac{1}{\varepsilon}, \log \frac{1}{\delta}, r, \mathbf{d})$. We then bound the sample complexity $m \doteq m_{\mathcal{W}, \mathcal{H}}(\varepsilon, \delta, W, g, \mathbf{d})$ of FPAC learning \mathcal{H} w.r.t. welfare class $\mathcal{W} \doteq \{W(\cdot; \mathbf{w})\}$ as

- 1) $m \leq m_{\mathcal{H}}(\sqrt[p]{\frac{\varepsilon}{2\lambda}}, \frac{\delta}{g}, r, \mathbf{d}) \in \text{Poly}\left(\sqrt[p]{\lambda}, \frac{1}{\sqrt[p]{\varepsilon}}, \log \frac{1}{\delta}, \log g, r, \mathbf{d}\right)$;
- 2) $p \in (0, 1] \implies m \in \text{Poly}\left(\sqrt[p]{r}, \frac{1}{\sqrt[p]{p}}, \frac{1}{\sqrt[p]{\varepsilon}}, \log \frac{1}{\delta}, \log g, \mathbf{d}\right)$;
- 3) $p = 0 \implies m \in \text{Poly}\left(\frac{w_{\min} \sqrt{r}}{w_{\min} \sqrt{\varepsilon}}, \log \frac{1}{\delta}, \log g, \mathbf{d}\right)$;
- 4) $p < 0 \implies m \in \text{Poly}\left(\frac{1}{\varepsilon}, \frac{1}{|r| \sqrt[p]{w_{\min}}}, \log \frac{1}{\delta}, \log g, r, \mathbf{d}\right)$; &
- 5) for any $c \in (0, 1)$, if $|p| \geq c$ and group weights obey the *nonnegligibility condition* $w_{\min} \geq \frac{c}{g}$, then $m \in \text{Poly}^{\frac{1}{c}}\left(\frac{1}{\varepsilon}, \frac{1}{g}, \log \frac{1}{\delta}, r, \mathbf{d}\right)$.

Proof. In each case, the FPAC learning algorithm $\mathcal{A}(\mathcal{D}_{1:g}, W, \varepsilon, \delta, \mathbf{d})$ is simply empirical welfare maximization on a sufficiently large sample, thus we need only bound the size of such a sufficient sample. Each item of this result is essentially a direct consequence of theorem 4.2, with lemmata 3.12 & 3.13 to bound Lipschitz and Hölder constants. It thus suffices to bound the constants λ for $\lambda \|\cdot\|_{\infty}$ Lipschitz continuity, or λ and α for $\lambda \|\cdot\|_{\infty}^{\alpha}$ Hölder continuity, for each of the classes under consideration.

In particular, item 1 follows from theorem 4.2 applied to any $W(\cdot; \mathbf{w})$ in the class under consideration. Then, items 2 & 3 follow from item 1, using lemma 3.13 items 2C & 1, respectively, to bound λ and α , and item 4 follows similarly, except using lemma 3.12 item 2 to bound the Lipschitz constant λ (thus $\alpha = 1$).

Finally, item 5 is slightly more involved, but again essentially reduces to item 1. In particular, observe that $|p| \geq c$ means we need not consider $p \approx 0$, and since c is constant, any exponential dependence on c remains polynomial in the remaining variables. Along with the *nonnegligibility condition* $w_{\min} \geq \frac{c}{g}$, this allows us to control the dependence of the Lipschitz constant λ for $p \leq -c$ as $\lambda \leq \frac{1}{|r| \sqrt[p]{w_{\min}}} \leq (\frac{g}{c})^{\frac{1}{c}}$, thus $\lambda \in \text{Poly}^{\frac{1}{c}}(g, \frac{1}{c})$ for $p \leq -c$. Similarly, for $p \geq c$, note that for welfare functions we need only consider $p \leq 1$, and observe that for $\alpha = c$, we have $\lambda = \frac{r^{1-c}}{c}$ by lemma 3.13 item 2C, which yields only $\text{Poly}^{\frac{1}{c}}(\frac{1}{\varepsilon}, \frac{1}{g})$ sample complexity terms. In either case, the desideratum is shown. \square

B Supplementary Experiments

We now present two additional one-armed bandit experiments using beta and Bernoulli noise models. Here utility samples are range $[0, 1]$ i.i.d. random variables with mean \mathbf{u}_i for each group i . For the Bernoulli model, we use $\text{BERNOULLI}(\mathbf{u}_i)$ random variables, and for the beta model, we use $\text{BETA}(\mathbf{u}_i, 1 - \mathbf{u}_i)$, which acts as continuous approximation of a $\text{BERNOULLI}(\mathbf{u}_i)$ coin, avoiding issues of discreteness, with exactly half the variance. The main difference here is that the variance of each estimator is now dependent on \mathbf{u}_i , being either $\frac{\mathbf{u}_i(1-\mathbf{u}_i)}{m_i} \leq \frac{1}{4m_i}$ in the Bernoulli case, or $\frac{\mathbf{u}_i(1-\mathbf{u}_i)}{2m_i} \leq \frac{1}{8m_i}$ in the beta case, as opposed to $\frac{1}{12m_i}$ in the uniform case. For all values of \mathbf{u}_i sufficiently far from $\frac{1}{2}$, these variance values are much smaller than under the uniform noise model, so we use only $m = 50$ samples unless otherwise noted.

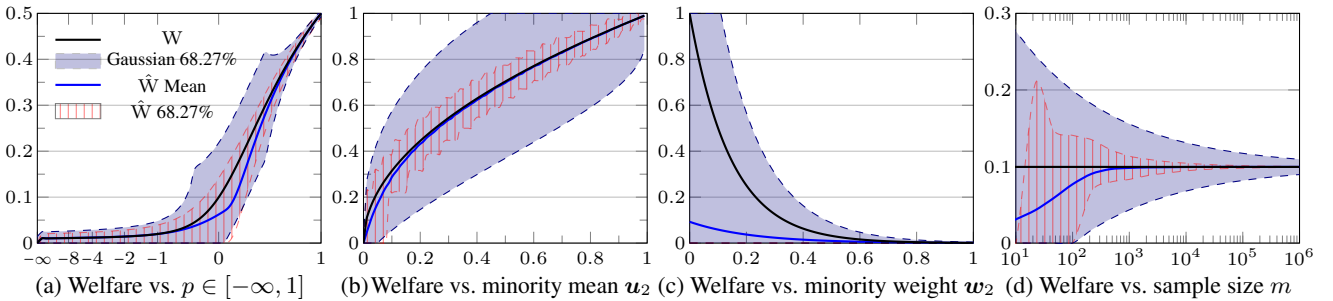


Figure A1: Estimating the Welfare of a 1-Armed Bandit under Bernoulli Noise. Each plot studies the response of welfare to one parameter, and the remaining parameters are selected from $p = 0$, $\mathbf{u} = (0.99, 0.01)$, $\mathbf{w} = (\frac{1}{2}, \frac{1}{2})$, and $m = 50$. All axes are linear, except A1a, which plots $p \in [-\infty, 1]$ by transforming $x = \frac{1}{\pi} \arctan(1 - p)$, and A1d, which is logarithmic in x .

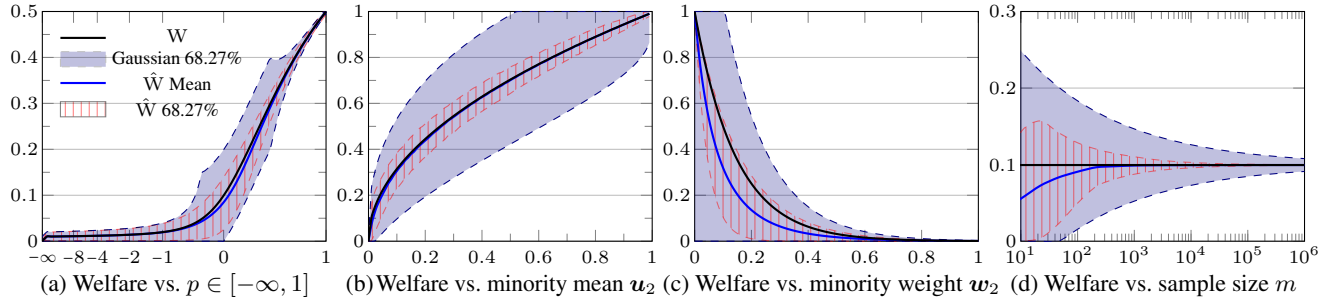


Figure A2: Estimating the Welfare of a 1-Armed Bandit under Beta Noise. Each plot studies the response of welfare to one parameter, and the remaining parameters are selected from $p = 0$, $\mathbf{u} = \langle 0.99, 0.01 \rangle$, $\mathbf{w} = \langle \frac{1}{2}, \frac{1}{2} \rangle$, and $m = 50$. All axes are linear, except A2a, which plots $p \in [-\infty, 1]$ by transforming $x = \frac{1}{\pi} \arctan(1 - p)$, and A2d, which is logarithmic in x .

The remainder of the experimental setup is identical to that under the uniform noise model, as described in section 5. In particular, we vary the parameters p , minority utility u_2 , minority group weight w_2 , and sample size m in order to study performance around the particularly challenging $p \approx 0$ and $w_{\min} \approx 0$ domains, and present the results in figures A1 & A2.

The beta and Bernoulli experiments are largely similar to the uniform noise experiment of section 5. In figures A1b & A2b which adjust the minority group utility u_2 , we observe that as the minority utility tends to 0, the empirical confidence intervals remain surprisingly wide, especially when considering that the variance of this coin is extremely small (also tending to 0). In contrast, as the coin bias tends to 1, confidence intervals get much smaller, as here again variance goes to 0, and here the welfare function is not sensitive to small changes. Note also that, generally speaking, the lower variances under the beta noise model (figure A2) result in tighter confidence bounds than the Bernoulli noise model (figure A1).