

Optimized Distortion in Linear Social Choice

Luise Ge, Gregory Kehne, Yevgeniy Vorobeychik

Washington University in St. Louis

Abstract

Social choice theory offers a wealth of approaches for selecting a candidate on behalf of voters based on their reported preference rankings over options. When voters have underlying utilities for these options, however, using preference rankings may lead to suboptimal outcomes vis-à-vis utilitarian social welfare. Distortion is a measure of this suboptimality, and provides a worst-case approach for developing and analyzing voting rules when utilities have minimal structure. However in many settings, such as common paradigms for value alignment, alternatives admit a vector representation, and it is natural to suppose that utilities are parametric functions thereof. We undertake the first study of distortion for linear utility functions. Specifically, we investigate the distortion of linear social choice for deterministic and randomized voting rules. We obtain bounds that depend only on the dimension of the candidate embedding, and are independent of the numbers of candidates or voters. Additionally, we introduce poly-time instance-optimal algorithms for minimizing distortion given a collection of candidates and votes. We empirically evaluate these in two real-world domains: recommendation systems using collaborative filtering embeddings, and opinion surveys utilizing language model embeddings, benchmarking several standard rules against our instance-optimal algorithms.

Code — <https://github.com/luisegehajjing/Optimized-Distortion-in-Linear-Social-Choice>

Extended version — <https://arxiv.org/pdf/2510.20020>

1 Introduction

Conventional voting rules map a profile of voter preference rankings over a set of candidates (options, outcomes, choices) to a winner. If, instead of rankings, we were instead to associate candidates with utility values for each voter, in many settings it would be natural to choose a social-welfare-maximizing candidate. But when voters’ preference rankings are generated by latent and unknown utilities, ranking-based voting rules can select candidates that are not welfare-maximizing. Thus, the limited information about utilities obtained from rankings can be viewed as an imperfect means of approximating a welfare-maximizing choice, and the quality of this approximation—the worst-case ratio of optimal welfare to welfare obtained by a voting rule—is known

as *distortion*. Since its introduction by Procaccia and Rosenschein (2006), an extensive literature has used distortion as means to compare voting rules, identify new rules, and understand the inherent difficulty in different settings of identifying desirable outcomes from ordinal and otherwise incomplete information (Anshelevich et al. 2021)

If we make no assumptions, no deterministic rules have bounded distortion, and (when there are many voters) the best randomized rule chooses a winning candidate uniformly at random. Starting with Procaccia and Rosenschein (2006), a prevalent approach is to assume each voter’s utilities are nonnegative and sum to 1 (also see Aziz (2020) for discussion of this assumption). Another mild but incomparable choice is to assume each voter’s utilities span the range $[0, 1]$, which generalizes approval preferences when each voter has at least one approval and one disapproval. In all such settings, even optimal voting rules have distortion polynomial in the number of candidates (Ebadian et al. 2024b).

What if we wish to make an approximately welfare-maximizing choice from among many candidates, and voters’ utilities are more structured than unit-sum or unit-range? Our motivating setting is reinforcement learning from human feedback (RLHF), where options are naturally represented as vector embeddings and voting data consists of rankings. In this context, a voter’s utility for an option reflects how well the option aligns with what the voter cares about—mathematically captured as the inner product between the voter’s preference vector and the option’s embedding. This linear utility model is both a common assumption in RLHF (Zhu, Jordan, and Jiao 2023; Ge et al. 2024) and arises naturally in recommendation systems (Pennock et al. 2000) and multi-objective decision making (Ngatchou, Zarei, and El-Sharkawi 2005), where candidates have feature representations that voters weight according to their preferences. More broadly, structured representation spaces for both voters and candidates are increasingly relevant as AI clones and artificial agents acting on behalf of humans are studied and deployed (Tan 2024; Lin 2024; Liang 2025), motivating the study of parametric utility structures over such spaces. Linear utilities provide a natural starting point for this investigation.

Unlike in many traditional applications of social choice, such as political elections, in RLHF and similar settings the space of options is extremely large. For example, the space

of all possible conversational responses to prompts is so vast as to defy enumeration. Consequently, distortion bounds that depend on the number of candidates m in these settings have little bite, since m can be exponential—or larger—in the number of features d .

We investigate distortion when utility functions are linear, and both the candidates’ and voters’ embeddings into common feature space \mathbb{R}^d are non-negative and normalized to 1. We focus on ℓ^1 normalization, and defer discussion of ℓ^2 normalization to Appendix C. Notably, under ℓ^1 normalization, our model recovers the unit-sum setting when $m = d$ and candidates embed to the standard basis of \mathbb{R}^d .

A significant benefit of measuring distortion over linear utilities in RLHF and alignment settings is its robustness to the introduction of duplicate candidates, i.e. clones. This is a criterion of particular importance in the application of social choice methods to alignment problems (Procaccia, Schiffer, and Zhang 2025; Berker et al. 2025). By way of illustration, consider a profile $\vec{\sigma}$ of voters’ rankings of C , and construct $\vec{\sigma}_{\cup c'}$ by adding a single duplicate c' of some $c \in C$ which all voters rank (weakly) directly below c . First, can the space of profile-consistent voter utilities over the original candidates C change from $\vec{\sigma}$ to $\vec{\sigma}_{\cup c'}$? For unit-sum and unit-range utilities, yes; for d -dimensional linear utilities, no. Second, m -dependent distortion bounds degrade from m to $m' = m + 1$, while d -dependent bounds do not.

Summary of contributions. On the positive side, we introduce three novel voting rules with distortion bound a function of *only the dimension* d . We summarize these results in Table 1. The first is the deterministic *max coordinate plurality* (MCP) rule. Given the candidate embeddings in $\mathbb{R}_{\geq 0}^d$, it first chooses one candidate that maximizes the value of each coordinate, and then selects a plurality winner among this subset. We show that MCP attains distortion $O(d^3)$ (Theorem 3). The second and third are randomized rules that are adaptations of the stable lottery rule of Ebadian et al. (2024b) to both the settings where candidate embeddings into \mathbb{R}^d are known and unknown. We consider the case when rules can access candidate embeddings to be the standard setting. Here we construct a *linear stable lotteries* rule, which attains the asymptotically optimal $\Theta(\sqrt{d})$ distortion for linear utilities (Theorem 6). When, as in more traditional social choice settings, access to candidate vectors is unavailable, a version of the stable lotteries rule which only samples candidates from suitably-sized (and randomly chosen) committees attains distortion $O(d)$; we dub this the *pure stable lotteries* rule (Theorem 9). We also show that random dictatorship attains $O(d^3)$ distortion in this setting (Theorem 7).

Alongside rules with provable worst-case guarantees, we develop a linear programming (LP)-based approach for computing instance-optimal distortion-minimizing candidates for a given set of options and voter preferences. While the resulting LP has an infinite set of constraints, we identify an efficient separation oracle that enables efficient optimization over it. We then empirically evaluate the extent to which our instance-optimal approach outperforms other rules.

C Embedding	Randomized	Deterministic
Known	$O(\sqrt{d})$ (Thm. 6)	$O(d^3)$ (Thm. 3)
	$\Omega(\sqrt{d})$ (Thm. 4)	$\Omega(d^2)$ (Thm. 2)
Unknown	$O(d)$ (Thm. 9)	-

Table 1: Distortion in d -dimensional linear social choice when candidates and voters are ℓ^1 -normalized. Note the unknown setting inherits lower bounds from the known setting.

Further related work. Our model diverges from two major strands of prior work in computational social choice: (i) classical frameworks with minimal utility assumptions, and (ii) metric distortion approaches that adopt latent representations, but differ fundamentally in objective and structure.

In classical frameworks, various assumptions are imposed directly on the utility values assigned by voters to candidates—such as unit-sum (each voter’s utilities sum to one), unit-range (utilities span a fixed range, e.g., $[0, 1]$), and approval (where each utility is either 0 or 1). Our approach strictly generalizes this first model by imposing an ℓ^1 normalization constraint on voter vectors and treating candidates as standard basis vectors, i.e., with $m = d$ (ℓ^∞ normalization would similarly generalize the second and third).

Metric distortion approaches are more closely aligned with our setting in that they also assume latent (metric) structure. However, their motivation typically centers on distance-based objectives in utility space—such as minimizing disutility in facility location problems—rather than maximizing social welfare. Although the relation $\text{disutil}_v(c) := 1 - u_v(c) = 1 - v^T c$ satisfies the triangle inequality and thus fits into the metric distortion framework, minimizing social cost in this sense is not equivalent to maximizing social welfare—our primary objective. As a result, our approach departs from metric distortion in both its goal and its fundamental structure.

Finally, we remark that while our work builds on the linear social choice framework introduced by Ge et al. (2024), we impose additional structure on the voter and candidate vectors for analysis purpose. In addition, although our work and that of Gözl, Haghtalab, and Yang (2025) share an interest in alignment, their model of the *distortion of alignment* is substantively different from ours. First, they bound utilities within $[0, 1]$ but otherwise make no assumptions about latent preference structure. Second, they assume distributions over both voters and candidates, with pairwise comparisons drawn i.i.d. from these distributions; distortion is then evaluated in expectation over this distribution. In contrast, we make no distributional assumptions and evaluate distortion in the worst case over all possible voter-candidate profiles. Finally, we consider truthful ordinal rankings, while Gözl, Haghtalab, and Yang (2025) assume a Bradley–Terry noise model for reported preferences, which enables some access to voter preference intensities. This reflects different modeling goals: we aim for robustness under accurate ordinal preference reporting, while they seek insights under probabilistic reporting with minimal utility structure assumptions.

2 Model and Preliminaries

We consider a setting with n voters V and m candidates C , where each voter $v \in V$ and each candidate $c \in C$ corresponds to a vector in $\mathbb{R}_{\geq 0}^d$. The d dimensions correspond to positive attributes that determine both voters' preferences and candidates' characteristics. We use superscripts to denote coordinates of a vector.

We assume that the utility of a voter v for a candidate c is given by the inner product $u_v(c) = v^\top c$. This utility function $u_v : C \rightarrow \mathbb{R}_{\geq 0}$, in turn, induces a preference ranking σ_v where $c \succ_v c'$ if $u_v(c) \geq u_v(c')$, with ties broken arbitrarily. Below we will introduce some constraints for v and c , and use \mathcal{U} to denote the set of all feasible utility functions induced by such v and c .

Let $\vec{u} = \{u_v\}_{v \in V}$ be the utility profile of all voters, and let $\vec{\sigma} = \{\sigma_v\}_{v \in V}$ be the preference profile of all voters. We use the notation $\vec{u} \triangleright \vec{\sigma}$ to indicate that u_v is consistent with σ_v for each voter $v \in V$. Voting rules have access to the rankings $\vec{\sigma}$, but *not* to the utilities \vec{u} .

If we know the utility functions u_v of all voters, a natural candidate selection criterion is *utilitarian welfare*, which is the sum of voter utilities: $\text{UW}(\vec{u}, c) := \sum_{v \in V} u_v(c)$. In this case the winning candidate $c^*(\vec{u})$ is the one with highest welfare, i.e., $c^*(\vec{u}) \in \arg \max_{c \in C} \text{UW}(\vec{u}, c)$.

Let f be a (possibly randomized) voting rule that maps a preference profile $\vec{\sigma}$ to a winner $c \in C$. For a fixed profile $\vec{\sigma}$, the *instance distortion* of f on $\vec{\sigma}$ is the worst-case ratio between the optimal utilitarian welfare and the utilitarian welfare of f ; that is,

$$D(f, \vec{\sigma}) := \max_{\vec{u} \triangleright \vec{\sigma}} \frac{\text{UW}(\vec{u}, c^*(\vec{u}))}{\mathbb{E}_{c \sim f(\vec{\sigma})} [\text{UW}(\vec{u}, c)]}.$$

For theoretical results, we are primarily interested in the overall *distortion* of f , which is the worst case over profiles:

$$D(f) := \max_{\vec{u} \triangleright \vec{\sigma}, \vec{u} \in \mathcal{U}^n} D(f, \vec{\sigma}).$$

We impose the following structural assumptions.

- **Non-negativity.** All voter and candidate vectors lie in the positive orthant, i.e., $v, c \in \mathbb{R}_{\geq 0}^d$. This ensures all utilities are non-negative, and avoids a mixed-sign objective.
- **Normalization.** We will assume $\|v\|_p = 1$ and $\|c\|_p = 1$ for all v, c . For $v \in V$ this can be viewed as constraining all voters to have equal influence on the mechanism (c.f. Aziz (2020)), and for both candidates and voters as identifying the *relative* magnitude of embedding components. We focus on ℓ^p norms for $p = 1$ in the main body,¹ but also provide results for ℓ^2 normalization in the supplemental material.
- **Expressiveness.** For theoretical analysis, we further assume $V \subset \text{Cone}(C)$, meaning every voter can be expressed as a non-negative linear combination of candidates. Intuitively, this can be understood as requiring that the candidate set is rich enough to describe voter preferences. This excludes the case where voters have 0 utility

¹Note that this together with non-negativity is equivalent to assuming $v, c \in \Delta_d$, where Δ_k denotes the k -dimensional simplex.

for all alternatives, ensuring each voter has meaningful preferences over the candidate set.

Normalization of both v and c is necessary to keep the model well-posed and the distortion finite, as the following example illustrates:

Example 1. Consider $n - 1$ voters who prefer c_1 to c_2 and one voter who strongly prefers c_2 due to a large utility spike in one coordinate. If candidate c_2 has an unbounded entry, it can yield unbounded utilitarian welfare even though $n - 1$ voters rank c_1 above c_2 .

We study several established voting rules. We define them here for convenience.

Definition 1. (Randomized Scoring Rules (RSRs)) Let $\vec{s} = (s^1, \dots, s^m)$ be a scoring vector with $s^1 \geq s^2 \geq \dots \geq s^m \geq 0$. For a candidate $c \in C$, let $\text{rank}_v(c)$ denote the position of c in voter v 's ranking (with $\text{rank}_v(c) = 1$ meaning a is ranked first). Then the score assigned to c by agent v is:

$$\text{score}_v(c, \vec{s}) := s^{\text{rank}_v(c)}.$$

The total score of c across all voters is:

$$\text{score}_V(c, \vec{s}) := \sum_{i \in n} \text{score}_v(c, \vec{s}).$$

The randomized scoring rule $f_{\vec{s}}^{\text{rand}}$ selects each alternative $c \in C$ with probability proportional to its total score:

$$\Pr[f_{\vec{s}}^{\text{rand}}(\vec{\sigma}) = c] := \frac{\text{score}_V(c, \vec{s})}{n \cdot \|\vec{s}\|_1}.$$

Notable examples include *Random Dictatorship*, which corresponds to the plurality scoring rule $\vec{s} = (1, 0, \dots, 0)$, and the *Randomized Harmonic* rule (Boutilier et al. 2015), which uses $\vec{s} = (1 + \frac{H_m}{m}, \frac{1}{2} + \frac{H_m}{m}, \dots, \frac{1}{m} + \frac{H_m}{m})$, where H_m is the m^{th} harmonic number.

We will also make use of stable lotteries, which were shown to exist for any preference profile by Cheng et al. (2020, Lemma 4).

Definition 2 (Stable Lotteries). Given a preference profile $\vec{\sigma}$, a committee W , and a candidate c , let $S_c(W) := \{v \in V : c \succ_v W\}$ be the set of voters who prefer c to all alternatives in W . We say a distribution \mathcal{W} over committees of size k is a stable lottery if for all $c \in C$,

$$\mathbb{E}_{W \sim \mathcal{W}} [|S_c(W)|] \leq \frac{n}{k}.$$

This has seen much recent use in computational social choice; our direct inspiration is its application by Ebadian et al. (2024b) to the design of distortion-optimal randomized rules in the unit-sum and related settings.

Special case: unit-sum utilities. The special case in which each candidate is a standard basis vector, i.e. $C = \{e_1, \dots, e_d\}$, recovers the well-studied unit-sum model, in which distortion was first introduced (Procaccia and Rosenschein 2006), and for which the worst-case distortion-optimal randomized rule is has distortion $\Theta(\sqrt{m})$ (Boutilier et al. 2015; Ebadian et al. 2024b), while the worst-case optimal deterministic rule has distortion $\Theta(m^2)$ (Caragiannis et al. 2017). By generalizing this setting, we inherit its lower bounds for both deterministic and randomized rules (Theorems 2 and 4).

3 Deterministic Rules

What utility guarantees can deterministic rules provide in this setting? How well do prominent deterministic rules fare? We begin with a useful lower bound on every individual voter's utility for their favorite candidate. Throughout, we use $\text{CH}(C)$ to denote the *convex hull* of a set C .

Lemma 1. *For any v and any candidates $C = \{c_j\}_{j \in [m]}$, the maximum utility of v is at least $\max_{c \in C} u_v(c) \geq \frac{1}{d}$.*

Proof. Let $\tilde{c} := \arg \max_{c \in C} v^T c$. For any $\alpha \in \Delta_m$, $v^T \tilde{c} \geq v^T (\sum_{i \in [m]} \alpha_i c_i) = \sum_{i \in [m]} \alpha_i (v^T c_i)$ since the maximum upper bounds any convex combination. Moreover, for ℓ^1 normalization, $v \in \text{Cone}(C)$ implies $v \in \text{CH}(C)$, i.e. $v = \sum_{i \in [m]} \beta_i c_i$ for some $\beta \in \Delta_m$. Substituting α with this specific β , we have $v^T \tilde{c} \geq v^T v \geq \frac{1}{d}$ by Cauchy-Schwarz (in particular, since $\|v\|_2 \cdot d \geq \|v\|_1 = 1$). \square

An analogous claim (Lemma 2 in Appendix C) holds for ℓ^2 normalization.

A natural deterministic rule is plurality (f_{Plur}). In the unit-sum setting, f_{Plur} is known to have distortion $\Theta(\min(n, m) \cdot m)$, which is asymptotically optimal (Caragiannis et al. 2017). We find that f_{Plur} attains worst-case distortion $\Theta(\min(n, m) \cdot d)$ for linear utilities. This matches the unit-sum guarantee when $d = m$, but scales poorly for $m, n \gg d$. The proof is deferred to Appendix A.

Theorem 1. $D(f_{\text{Plur}}) = \Theta(\min(m, n) \cdot d)$.

Is this dependence on m and n unavoidable for all deterministic rules, or can it be circumvented? For deterministic rules, we inherit lower bounds from the work of Caragiannis et al. (2017) in the unit-sum setting.

Theorem 2. [Theorem 1 of Caragiannis et al. (2017)] *Suppose $m, n \geq d$. Then for any deterministic voting rule f , we have $D(f) = \Omega(d^2)$.*

This lower bound depends only on d , and raises the question of whether it is possible to improve upon plurality by avoiding a dependence on n and m in the distortion bound. Indeed it is: we introduce a rule we dub *maximum coordinate plurality (MCP)* and denote by f_{MCP} .

Definition 3 (Maximum Coordinate Plurality). *Let*

$$\hat{C} := \left\{ c_i \in C \mid c_i = \arg \max_{c \in C} c^i \text{ for each } i \in [d] \right\}$$

be a set of at most d candidates such that for each coordinate i , a candidate maximal in that coordinate is included. The maximum coordinate plurality rule f_{MCP} then restricts the profile to \hat{C} and selects the plurality winner.

Theorem 3. $D(f_{\text{MCP}}) = O(d^3)$.

Proof. For any fixed voter v , it has its maximum coordinate at least $\frac{1}{d}$; call this coordinate i . Since $v \in \text{CH}(C)$, the candidate $c_i \in \hat{C}$ therefore has i -th coordinate at least $\frac{1}{d}$, and so the welfare conferred to v by their favorite choice in \hat{C} is at least $\max_{c \in \hat{C}} v^T c \geq \frac{1}{d^2}$. Since $|\hat{C}| \leq d$, this plurality winner $\hat{c} \in \hat{C}$ receives at least $\frac{n}{d}$ votes, and so $\text{UW}(\hat{c}) \geq \frac{n}{d} \frac{1}{d^2}$. As $\max_{c \in C} \text{UW}(c) \leq n$, the claim follows. \square

Though this guarantee still exceeds the lower bound by a factor of d , the $\Omega(d^2)$ lower bound is already quite large. Can this be improved by randomizing over candidates?

4 Randomized Rules

As previously mentioned, we also inherit a distortion lower bound for any randomized rule from the unit-sum setting.

Theorem 4. (Boutilier et al. 2015) *Suppose $m \geq d$ and $n \geq \sqrt{d}$. Then there exists a preference profile $\vec{\sigma}$ such that for any randomized rule f , we have $D(f, \vec{\sigma}) = \Omega(\sqrt{d})$.*

We include a proof for completeness in Appendix A.

4.1 Known Candidate Embeddings

Consider the center of the simplex Δ_d , denoted by $\mu := (\frac{1}{d}, \dots, \frac{1}{d})$. Observe that $\text{UW}(\mu) = \frac{n}{d}$, which implies a distortion at most d is possible when $\mu \in C$. However in general, $\mu \notin \text{CH}(C)$. This motivates the goal of *approximating* the uniform candidate μ by a point within the convex hull of candidates $\text{CH}(C)$, assuming the candidate embeddings are known. We therefore introduce a randomized rule whose output distribution matches the *reverse information projection* of μ onto $\text{CH}(C)$.

Definition 4 (Uniform Projection Rule (f_{UProj})). *Given candidate locations $C \subset \mathbb{R}_{\geq 0}^d$, the uniform projection rule f_{UProj} defines a distribution $\{p_c\}_{c \in C}$ over candidates such that the expected candidate vector $\hat{c} := \sum_{c \in C} p_c \cdot c$ minimizes the Kullback–Leibler (KL) divergence from the uniform vector, i.e., $\hat{c} := \arg \min_{x \in \text{CH}(C)} \text{KL}(\mu \| x)$.*

Theorem 5. *The expected welfare of f_{UProj} is at least $\text{UW}(f_{\text{UProj}}) \geq \frac{n}{d}$. As a consequence, $D(f_{\text{UProj}}) = O(d)$.*

Proof of Theorem 5. First observe that

$$\text{KL}(\mu \| x) = \sum_{i=1}^d \mu^i \ln \frac{\mu^i}{x^i} = \sum_i \mu^i \ln \mu^i - \sum_{i=1}^d \mu^i \ln x^i,$$

so minimizing $\text{KL}(\mu \| x)$ over the convex set $\text{CH}(C)$ is equivalent to minimizing the smooth, convex function $f(x) = -\sum_{i=1}^d \mu^i \ln x^i$ subject to $x \in \text{CH}(C)$. The first-order optimality condition gives $\nabla f(x^*)^\top (v - x^*) \geq 0$ for every voter $v \in \text{CH}(C)$. Since $\frac{\partial f}{\partial x^i}(x) = -\frac{\mu^i}{x^i}$, we have

$$\nabla f(x^*)^\top (v - x^*) = -\sum_{i=1}^d \frac{\mu^i}{x^{*i}} (v^i - x^{*i}) \geq 0$$

$$\implies \sum_{i=1}^d \frac{\mu^i v^i}{x^{*i}} \leq \sum_{i=1}^d \mu^i = 1.$$

Since $\mu^i = \frac{1}{d}$, this yields $\sum_{i=1}^d \frac{v^i}{x^{*i}} \leq d$. Now define the auxiliary sequences $a_i := \sqrt{\frac{v^i}{x^{*i}}}$ and $b_i := \sqrt{v^i x^{*i}}$. Then $\sum_{i=1}^d a_i b_i = \sum_{i=1}^d v^i = 1$ and, by, Cauchy–Schwarz,

$$1 = \left(\sum_i a_i b_i \right)^2 \leq \left(\sum_i a_i^2 \right) \left(\sum_i b_i^2 \right)$$

$$= \left(\sum_i \frac{v^i}{x^{*i}} \right) \left(\sum_i v^i x^{*i} \right).$$

Combining with $\sum_i v^i/x^{*i} \leq d$ gives $\sum_{i=1}^d x^{*i} v^i \geq 1/d$. Now as we have n voters and the above inequality holds for arbitrary v , the total utilitarian welfare is at least $\frac{n}{d}$.

The bound on the distortion of f_{UProj} then follows because the maximum utility for any voter is at most 1. \square

This randomized $O(d)$ -distortion rule can be seen as generalizing uniform candidate selection in the unit-sum setting. A natural candidate for improving upon this is the harmonic rule f_{HR} , which obtains near-optimal distortion $\Theta(\sqrt{m \log m})$ (Boutillier et al. 2015; Bhaskar, Dani, and Ghosh 2018) for unit-sum utilities. However this performance does not generalize to linear utilities, wherein the introduction of duplicate candidates does not constrain the underlying utility profile. Instead, it turns out that here f_{HR} has distortion unbounded in d ; we discuss this and other randomized positional scoring rules in Section 4.2.

Fortunately, Theorem 5 does indirectly lead to distortion sublinear in d . Using it, we may adapt the stable lottery rule of (Ebadian et al. 2024b) to the linear utilities setting to achieve a better distortion bound.

Definition 5 (Linear Stable Lottery Rule (f_{LSLR})). *Given a stable lottery \mathcal{W} over committees of size $k = \sqrt{d}$, the linear stable lottery rule f_{LSLR} on profile $\vec{\sigma}$ chooses each $c \in C$ with probability $\frac{1}{2\sqrt{d}} \Pr_{W \sim \mathcal{W}(\vec{\sigma})}[c \in W] + \frac{1}{2} \Pr_{c' \sim f_{\text{UProj}}(\vec{\sigma})}[c' = c]$.*

Here f_{UProj} takes the place of uniform selection.

Theorem 6. $D(f_{\text{LSLR}}) = O(\sqrt{d})$.

Proof. We follow the proof of Ebadian et al. (2024b) closely, albeit in our notation. The principal difference is our use of f_{UProj} instead of uniform selection over candidates.

To that end, we combine Definition 2 and Theorem 5 in order to relate the expected welfare conferred from each part of f_{LSLR} . For any committee W of size k and any $a \in W$, let $S_a(W) \subseteq V$ denote the voters for which $a \succ_v W$, and $\bar{S}_a(W)$ its complement. Then

$$\begin{aligned} \sum_{v \in V} \sum_{c \in W} u_v(c) &= \sum_{v \in S_a(W)} \sum_{c \in W} u_v(c) + \sum_{v \in \bar{S}_a(W)} \sum_{c \in W} u_v(c) \\ &\geq \sum_{v \in S_a(W)} \sum_{c \in W} u_v(c) + \sum_{v \in \bar{S}_a(W)} u_v(c^*) \\ &\geq \sum_{v \in S_a(W)} (u_v(c^*) - 1) + \sum_{v \in \bar{S}_a(W)} u_v(c^*) \\ &= \sum_{v \in V} u_v(c^*) - |\{S_a(W)\}|. \end{aligned} \quad (1)$$

Taking the expectation over a stable lottery \mathcal{W} and applying Definition 2, we obtain

$$\mathbb{E}_{W \sim \mathcal{W}} \left[\sum_{c \in W} \text{UW}(c) \right] \geq \text{UW}(c^*) - \mathbb{E}_{W \sim \mathcal{W}} [|\{S_a(W)\}|]$$

$$\begin{aligned} &\geq \text{UW}(c^*) - \frac{n}{k} \\ &\geq \text{UW}(c^*) - \frac{d}{k} \mathbb{E}_{c \sim f_{\text{UProj}}} [\text{UW}(c)], \end{aligned} \quad (2)$$

where the last step follows from Theorem 5.

We now let $x = \frac{1}{2}x_1 + \frac{1}{2}x_2$ denote the distribution over C of f_{LSLR} , where x_1 is the stable lottery part and x_2 is the f_{UProj} part. Then applying (2) and letting $k = \sqrt{d}$,

$$\begin{aligned} k \cdot \text{UW}(x_1) &= \mathbb{E}_{W \sim \mathcal{W}} \left[\sum_{c \in W} \text{UW}(c) \right] \\ &\geq \text{UW}(c^*) - \frac{d}{k} \mathbb{E}_{c \sim f_{\text{UProj}}} [\text{UW}(c)] \\ k \cdot \text{UW}(x_1) + \frac{d}{k} \text{UW}(x_2) &\geq \text{UW}(c^*) \\ \mathbb{E}_{c \sim f_{\text{LSLR}}} [\text{UW}(c)] &\geq \frac{1}{2\sqrt{d}} \text{UW}(c^*). \\ \frac{\text{UW}(c^*)}{\mathbb{E}_{c \sim f_{\text{LSLR}}} [\text{UW}(c)]} &\leq 2\sqrt{d} \end{aligned}$$

This proves the claim. \square

Computing a stable lottery \mathcal{W} requires only ordinal information; however in order to identify the distribution from which f_{UProj} samples, the embedding of C into $\mathbb{R}_{\geq 0}^d$ must be known by the rule. This raises the question of what can be done with ordinal preferences when both voter *and* candidate embeddings are unknown.

4.2 Unknown Candidate Embeddings

Even when the locations of the candidate vectors $c \in \mathbb{R}_{\geq 0}^d$ are unknown, some established rules exhibit distortion that is bounded as a function of d . Our primary example is random dictatorship, for which proof is deferred to Appendix A.

Theorem 7. $D(f_{\text{RD}}) = \Omega(d)$, and also $D(f_{\text{RD}}) = O(d^3)$.

One feature of this model is that the worst-case distortion of many randomized positional scoring rules (Appendix 1) is quite similar. In particular, consider the RSRs where, for given m , the score vector is given by $\vec{s} = \frac{1}{S_m}(s_1, s_2, \dots, s_m)$ for some fixed sequence s_1, s_2, \dots , and where $S_i = \sum_{j \leq i} s_j$. (Note that this contains f_{RD} and f_{HR} and the uniform distribution over candidates, but not Borda.) Then there are two cases. If $S_m \rightarrow S$ converges, then the distortion of f_s on *any* instance is within S/s_1 of f_{RD} , and we can make the behavior of f_s approach that of f_{RD} by cloning each candidate sufficiently many times. And if S_m diverges, then cloning bad candidates can lead to poor performance on instances which are easy for f_{RD} . We illustrate this by relating the performance of f_{HR} to that of f_{RD} , the proof of which also appears in Appendix A.

Theorem 8. *The distortion of f_{HR} is unbounded as a function of d . However, for any profile $\vec{\sigma}$,*

$$D(f_{\text{HR}}, \vec{\sigma}) \leq D(f_{\text{RD}}, \vec{\sigma}) \cdot (\log m + 1).$$

Thus, the distortion of f_{HR} is never much better than that of f_{RD} , and if—as we expect—the worst-case profiles for f_{RD} have $\text{poly}(d)$ candidates, then its worst-case distortion cannot be much better than that of f_{RD} , even for small m .

Can we improve upon this bound when candidate embeddings are unknown? The role that f_{UProj} plays in the design and analysis of f_{LSLR} —and the role uniform selection plays in the stable lotteries for Ebadian et al. (2024b)—is in some sense both an absolute lower bound on welfare when the maximum candidate welfare is not large, and sample access to it. We might then let f_{RD} take the role of f_{UProj} in f_{LSLR} to attain distortion $O(d^{3/2})$ for unknown embeddings.

However even in the absence of a better rule, it turns out stable lotteries can still be used by furnishing a direct lower bound on the maximum welfare. Consider the following:

Definition 6 (Pure Stable Lottery Rule (f_{PSLR})). *Given a stable lottery \mathcal{W} over committees of size $2d$, the pure stable lottery rule f_{PSLR} chooses $c \in C$ w.p. $\frac{1}{2d} \Pr_{W \sim \mathcal{W}}[c \in W]$.*

Using larger committees and that $\max_{c \in C} \text{UW}(c) \geq \frac{n}{d}$, we avoid the need for a second sampling component.

Theorem 9. $D(f_{\text{PSLR}}) = O(d)$.

Proof (sketch). The difference between this and the proof of Theorem 6 is that the lottery is over committees of size $k = 2d$, and that we do not use f_{UProj} as in (2). We instead look to Theorem 5, which demonstrates there always exists a distribution $\{p_c\}$ over candidates such that $\mathbb{E}_{c \sim f_{\text{UProj}}}[\text{UW}(c)] \geq \frac{n}{d}$. Then by averaging, there always exists a $c \in C$ such that $\text{UW}(c) \geq \frac{n}{d}$, and in particular $\text{UW}(c^*) \geq \frac{n}{d}$. Picking things up just before (2) with $k = 2d$,

$$\begin{aligned} \mathbb{E}_{W \sim \mathcal{W}} \left[\sum_{c \in W} \text{UW}(c) \right] &\geq \text{UW}(c^*) - \mathbb{E}_{W \sim \mathcal{W}} [\{S_a(W)\}] \\ &\geq \text{UW}(c^*) - \frac{n}{2d} \geq \frac{1}{2} \cdot \text{UW}(c^*). \end{aligned}$$

The rest of the proof proceeds as before, though without the need for x_2 . It appears in full in Appendix A. \square

5 Optimizing Distortion

In the preceding sections we analyzed the asymptotic behavior of voting rules, but the lower bounds are often driven by pathological cases. In practice, we are more interested in this welfare approximation guarantee for a given profile. That is: *Given a preference profile, can we compute the best possible distortion—and design a rule that achieves it?*

The computational tractability of the distortion-instance-optimal rule in the unit-sum setting was established by Boutilier et al. (2015) and clarified by Ebadian et al. (2024a). For linear utilities we also answer both questions affirmatively, showing that instance-optimal rules—both deterministic and randomized—can be computed efficiently.

Definition 7 (Feasible Region for \bar{v}). *Let $\mathcal{F}_j = \{v \in \Delta_d : \langle c_a - c_b, v \rangle \geq 0 \text{ whenever } c_a \succ_j c_b\}$ be the set of voter vectors consistent with σ_{v_j} . Then the feasible region for the average voter $\bar{v} = \frac{1}{n} \sum_{j=1}^n v_j$ is $\mathcal{F} = \frac{1}{n} \sum_{j=1}^n \mathcal{F}_j$.*

Theorem 10 (Deterministic Instance-Optimal Rule). *Given a preference profile $\vec{\sigma}$ and any $\varepsilon > 0$, one can compute a deterministic voting rule f such that $D(f, \vec{\sigma}) \leq \min_{c \in C} D(c, \vec{\sigma}) + \varepsilon$ in time $\text{poly}(n, m, \log \frac{1}{\varepsilon})$.*

Proof. Let \mathcal{F} be the feasible region for the average voter vector \bar{v} induced by $\vec{\sigma}$ (Definition 7). For each pair $c_1, c_2 \in C$, we aim to compute the best possible upper bound on $\frac{\text{UW}(c_1)}{\text{UW}(c_2)} = \frac{\bar{v}^\top c_1}{\bar{v}^\top c_2}$ over all $\bar{v} \in \mathcal{F}$. Since we cannot compute this ratio exactly without knowing \bar{v} , we upper-bound it by the largest value $\frac{1}{\beta}$ such that $\bar{v}^\top c_1 \leq \frac{1}{\beta} \bar{v}^\top c_2 \forall \bar{v} \in \mathcal{F}$.

Equivalently, for fixed c_1, c_2 , we find the largest $\beta_{c_1, c_2} \in [0, 1]$ such that the following LP has a non-negative optimal value: $\min_{\bar{v} \in \mathcal{F}} \langle \beta c_1 - c_2, \bar{v} \rangle$. This is doable via binary search over β to precision ε , requiring $O(\log \frac{1}{\varepsilon})$ iterations per pair.

Once the distortion bounds are computed for all $O(m^2)$ candidate pairs, we select the candidate $\hat{c} \in C$ minimizing the worst-case distortion $\hat{c} := \arg \min_{c \in C} \max_{c' \in C} \beta_{c', c}$. Returning this candidate defines the deterministic rule f , and by construction we have $D(f, \vec{\sigma}) \leq \min_{c \in C} D(c, \vec{\sigma}) + \varepsilon$. Since each linear program is solvable in polynomial time in n, m, d , and the number of binary search iterations is $O(\log \frac{1}{\varepsilon})$, the total runtime is $\text{poly}(n, m, d, \log \frac{1}{\varepsilon})$. \square

Note that the procedure described above can be used to compute the distortion of any rule on a given preference profile. We also apply it in our experiments.

Theorem 11 (Randomized Instance-Optimal Rule). *Given a preference profile $\vec{\sigma}$ and any $\varepsilon > 0$, a randomized voting rule f with $D(f, \vec{\sigma}) \leq \min_{f' \in \Delta(C)} D(f', \vec{\sigma}) + \varepsilon$ can be computed in time $\text{poly}(n, m, \frac{1}{\varepsilon})$, where the minimum is taken over all distributions over candidates.*

Proof. Let $p \in \Delta_C$ be a distribution over candidates, and let $\hat{c} := \sum_{i=1}^m p_i c_i$ be the expected candidate vector under p . We aim to find a distribution p that minimizes the distortion:

$$D(p, \vec{\sigma}) = \sup_{\bar{v} \in \mathcal{F}} \frac{\max_{c \in C} \langle c, \bar{v} \rangle}{\langle \hat{c}, \bar{v} \rangle}.$$

Since this expression may be unbounded, we instead maximize its reciprocal $\beta \in [0, 1]$, subject to the constraint:

$$\langle \hat{c}, \bar{v} \rangle \geq \beta \cdot \max_{c \in C} \langle c, \bar{v} \rangle \quad \forall \bar{v} \in \mathcal{F}.$$

This yields the following optimization problem:

$$\max_{\substack{p \in \Delta(C) \\ \beta \in [0, 1]}} \left\{ \beta : \sum_{i=1}^m p_i \langle c_i, \bar{v} \rangle \geq \beta \cdot \max_{c \in C} \langle c, \bar{v} \rangle \quad \forall \bar{v} \in \mathcal{F} \right\}.$$

Although \mathcal{F} contains infinitely many constraints, we can apply the ellipsoid method using a separation oracle. For a candidate solution $(\hat{\beta}, \hat{p})$, feasibility reduces to verifying:

$$\sum_{i=1}^m \hat{p}_i \langle c_i, \bar{v} \rangle \geq \hat{\beta} \cdot \langle \hat{c}, \bar{v} \rangle \quad \forall c \in C, \bar{v} \in \mathcal{F}.$$

For each $c \in C$, this can be checked by solving an LP:

$$\min_{\bar{v} \in \mathcal{F}} \left\langle \sum_{i=1}^m \hat{p}_i c_i - \hat{\beta} \cdot c, \bar{v} \right\rangle.$$

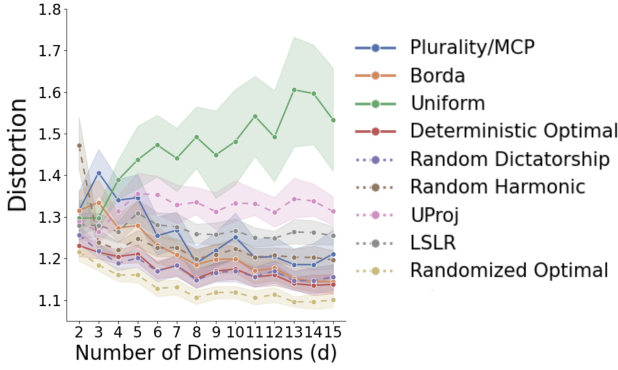


Figure 1: Instance distortion $D(f, \bar{\sigma})$ for MovieLens

If the minimum is nonnegative for all $c \in C$, the solution is feasible; otherwise, the auxiliary LP provides a separating hyperplane. Since each $\bar{v} \in \mathcal{F}$ is defined via $O(m^2n)$ linear constraints, each LP is of polynomial size. Running the ellipsoid method with this oracle to precision ε yields a solution in time $\text{poly}(n, m, d, \frac{1}{\varepsilon})$. \square

Practical implementation. Although the ellipsoid method is polynomial-time, it is often impractical. We instead use the more efficient *column generation* (Dantzig and Wolfe 1960) (see Appendix B). Because the algorithm uses only pairwise comparisons, it also supports partial rankings.

6 Experiments

We evaluate our instance-optimal rules alongside f_{UProj} , f_{MCP} , f_{LSLR} , and several other common rules, on two real-world datasets. Specifically, we measure both the *instance distortion* (relative to the underlying utility profile \bar{v}) and the *worst-case distortion* (over all utility profiles consistent with the observed preferences $\bar{\sigma}$). Due to space constraints, we report results for the latter, and only as a function of the dimension d in the main text. Results varying n, m , running times, and instance distortion, are deferred to Appendix D.

Datasets. We consider the MovieLens 100K (Harper and Konstan 2015) and abortion opinion survey (Fish et al. 2024) datasets. MovieLens contains 100K movie ratings from 1000 users (voters) over 1700 movies. We translate it into our setting via approximate matrix decomposition for each dimension d , adding the nonnegativity and normalization constraints. We then randomly subsample 20 movie preference votes for each dimension d with $n = 100$ and $m = 25$. The abortion opinions dataset includes ratings from 100 individuals on 5 different opinion statements. For this setting, we embed the candidates using Matryoshka embeddings (Kusupati et al. 2022), followed by dimensionality reduction via principal component analysis (PCA). The embedding choice for either dataset is not essential to our method and is intended only as a representative example.

Results. Figure 2 presents instance-based distortion bounds for MovieLens (left) and abortion opinion survey (right). Our first observation is that, as we expect, the instance-optimal approaches we propose outperform all alternatives,

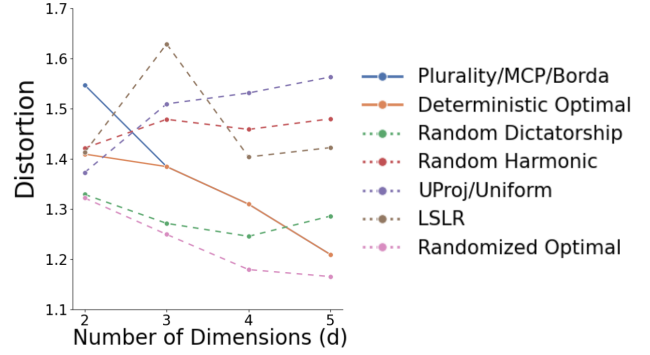


Figure 2: Instance distortion $D(f, \bar{\sigma})$ for Abortion Survey

with improvement appearing to increase with d . This is true for both deterministic and randomized rules. Second, and somewhat surprisingly, we find that worst-case distortion bounds *decrease with dimension d* , despite the fact that distortion lower bounds increase in d . Thus, as the dimension of candidate representations increases, using rankings in place of utilities becomes empirically less consequential.

7 Conclusion

We initiate the study of distortion in the setting of linear social choice, which we anticipate will play a role in the integration of social choice into neural and representation-based models. For deterministic rules we introduce the maximum coordinate plurality (MCP) rule, which we prove obtains worst-case distortion $\mathcal{O}(d^3)$ and is within $\Theta(d)$ of optimal. Since MCP requires access to candidate positions, we leave open whether d -dependent distortion is even *possible* for deterministic rules without such access.

For randomized rules, we construct the *linear stable lotteries rule*, and show it attains asymptotically optimal worst-case distortion $\Theta(\sqrt{d})$ when candidate locations are known. We show stable lotteries yield $O(d)$ distortion even when candidate locations are unknown, leaving unresolved the optimal distortion in this setting. We also establish preliminary results for ℓ^2 voter and candidate normalization.

Beyond the utilitarian welfare objective, future work might also explore *Nash welfare* or *proportional fairness*, as studied by Ebadian et al. (2024b), or worst-case *regret* minimization (Caragiannis et al. 2017; Kahng and Kehne 2022). Finally, given the connection between linear social choice and reward learning, we plan to investigate the robustness of various rules under both learning errors and strategic voters.

Acknowledgments

We thank Dominik Peters for helpful suggestions regarding the implementation of stable lotteries. This work was supported in part by the NSF (IIS-2214141, CCF-2403758), ARO (W911NF-25-1-0059), ONR (N00014-24-1-2663), the Foresight Institute, and Amazon.

References

- Anshelevich, E.; Filos-Ratsikas, A.; Shah, N.; and Voudouris, A. A. 2021. Distortion in Social Choice Problems: The First 15 Years and Beyond. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-21)*.
- Aziz, H. 2020. Justifications of Welfare Guarantees Under Normalized Utilities. *ACM SIGecom Exchanges*.
- Berker, R. E.; Casacuberta, S.; Robinson, I.; Ong, C.; Conitzer, V.; and Elkind, E. 2025. From Independence of Clones to Composition Consistency: A Hierarchy of Barriers to Strategic Nomination. In *Proceedings of the 26th ACM Conference on Economics and Computation*.
- Bhaskar, U.; Dani, V.; and Ghosh, A. 2018. Truthful and Near-Optimal Mechanisms for Welfare Maximization in Multi-winner Elections. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Boutilier, C.; Caragiannis, I.; Haber, S.; Lu, T.; Procaccia, A. D.; and Sheffet, O. 2015. Optimal Social Choice Functions: A Utilitarian View. *Artificial Intelligence*.
- Caragiannis, I.; Nath, S.; Procaccia, A. D.; and Shah, N. 2017. Subset Selection via Implicit Utilitarian Voting. *Journal of Artificial Intelligence Research*.
- Cheng, Y.; Jiang, Z.; Munagala, K.; and Wang, K. 2020. Group Fairness in Committee Selection. *ACM Transactions on Economics and Computation (TEAC)*.
- Dantzig, G. B.; and Wolfe, P. 1960. Decomposition Principle for Linear Programs. *Operations Research*.
- Ebadian, S.; Filos-Ratsikas, A.; Latifian, M.; and Shah, N. 2024a. Computational Aspects of Distortion. In *The 23rd International Conference on Autonomous Agents and Multi-agent Systems*.
- Ebadian, S.; Kahng, A.; Peters, D.; and Shah, N. 2024b. Optimized Distortion and Proportional Fairness in Voting. *ACM Transactions on Economics and Computation*.
- Fish, S.; Gözl, P.; Parkes, D. C.; Procaccia, A. D.; Rusak, G.; Shapira, I.; and Wüthrich, M. 2024. Generative Social Choice. In *Proceedings of the 25th ACM Conference on Economics and Computation*.
- Ge, L.; Halpern, D.; Micha, E.; Procaccia, A. D.; Shapira, I.; Vorobeychik, Y.; and Wu, J. 2024. Axioms for AI Alignment from Human Feedback. In *Advances in Neural Information Processing Systems*.
- Gözl, P.; Haghtalab, N.; and Yang, K. 2025. Distortion of AI Alignment: Does Preference Optimization Optimize for Preferences? In *Advances in Neural Information Processing Systems*.
- Harper, F. M.; and Konstan, J. A. 2015. The MovieLens Datasets: History and Context. *ACM Transactions on Interactive Intelligent Systems (TiS)*.
- Kahng, A.; and Kehne, G. 2022. Worst-Case Voting When the Stakes Are High. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Kusupati, A.; Bhatt, G.; Rege, A.; Wallingford, M.; Sinha, A.; Ramanujan, V.; Howard-Snyder, W.; Chen, K.; Kakade, S.; Jain, P.; and Others. 2022. Matryoshka Representation Learning. In *Advances in Neural Information Processing Systems*.
- Liang, A. 2025. Artificial Intelligence Clones. In *Proceedings of the 26th ACM Conference on Economics and Computation, EC 2025*.
- Lin, B. 2024. Grindr Aims to Build the Dating World’s First AI ‘Wingman’. *The Wall Street Journal*. October 5.
- Ngatchou, P.; Zarei, A.; and El-Sharkawi, A. 2005. Pareto Multi Objective Optimization. In *Proceedings of the 13th International Conference on Intelligent Systems Application to Power Systems*. IEEE.
- Pennock, D. M.; Horvitz, E.; Giles, C. L.; and Others. 2000. Social Choice Theory and Recommender Systems: Analysis of the Axiomatic Foundations of Collaborative Filtering. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Procaccia, A. D.; and Rosenschein, J. S. 2006. The Distortion of Cardinal Preferences in Voting. In *Proceedings of the International Workshop on Cooperative Information Agents*. Springer.
- Procaccia, A. D.; Schiffer, B.; and Zhang, S. 2025. Clone-Robust AI Alignment. In *Proceedings of the International Conference on Machine Learning*. PMLR.
- Tan, E. 2024. Dating Apps Suck. A.I. Clones Are Making Them Even Weirder. *The New York Times*. December 11.
- Zhu, B.; Jordan, M.; and Jiao, J. 2023. Principled Reinforcement Learning with Human Feedback from Pairwise or K-wise Comparisons. In *Proceedings of the International Conference on Machine Learning*. PMLR.