Direct Preference Optimization With Unobserved Preference Heterogeneity

Keertana Chidambaram Stanford University vck@stanford.edu Karthik Vinay Seetharaman Stanford University kvseet17@stanford.edu Vasilis Syrgkanis Stanford University vsyrgk@stanford.edu

Abstract

RLHF has emerged as a pivotal step in aligning language models with human objectives and values. It typically involves learning a reward model from human preference data and then using reinforcement learning to update the generative model accordingly. Conversely, Direct Preference Optimization (DPO) directly optimizes the generative model with preference data, skipping reinforcement learning. However, both RLHF and DPO assume uniform preferences, overlooking the reality of diverse human annotators. This paper presents a new method to align generative models with varied human preferences. We propose an Expectation-Maximization adaptation to DPO, generating a mixture of models based on latent preference types of the annotators. We then introduce a min-max regret ensemble learning model to produce a single generative method to minimize worst-case regret among annotator subgroups with similar latent factors. Our algorithms leverage the simplicity of DPO while accommodating diverse preferences. Experimental results validate the effectiveness of our approach in producing equitable generative policies.

1 Introduction

Reinforcement Learning from Human Feedback (RLHF) has emerged as one of the leading methods to align Language Models (LMs) to human preferences [34, 45, 51]. RLHF focuses on learning a single reward model from human preference data and uses that to fine-tune and align the LM. To sidestep potentially expensive reinforcement learning, Direct Preference Optimization (DPO) [37] is an alignment method that optimizes the LM policy directly using the preference data. However, DPO implicitly uses the same reward model as RLHF to train the LM. This reward model reflects the majority opinion of the preference data annotators

and caters to that majority. If the annotator population is not representative of the general population, then this comes at the cost of neglecting groups underrepresented in the annotators, leading to misrepresentation of preferences. On the other hand, if the annotator population is representative, then opinions of minority groups in the general population are shunned, causing bias and discrimination.

Most papers that try to deal with this issue learn a reward model and then use a standard RL framework such as PPO to align the LM. However, DPO has several advantages over RLHF, eliminating the need for a reward model and leading to a more stable pipeline. [62] utilizes these benefits by developing an algorithm that directly optimizes policy by implicitly learning a multi-objective reward model. However, methods that rely on a multidimensional reward model [50, 62] implicitly or explicitly have two main drawbacks. First, these methods typically require annotators to rate data on a multi-dimensional scale, with each dimension corresponding to a different objective like safety or accuracy. This data is both more costly and harder to obtain compared to binary preference data [10]. Second, the different rating objectives must be determined ahead of the data collection stage. This can be a difficult task as there are many latent factors that might affect the preferences of annotators [43], which can be difficult to discern. For example, if we collect ratings based on helpfulness and harmfulness similar to [5], these rankings might not fully explain some preference decisions made because of cultural, political, or geographical inclinations.

We propose a pipeline of two algorithms to sidestep the need for RLHF for a heterogeneous population, allowing us to cater to diverse preferences without the need for reinforcement learning, letting us reap the added benefits of DPO. In particular, we propose Expectation Maximization Direct Preference Optimization (EM-DPO) and MinMax Direct Preference Optimization (MinMax-DPO). EM-DPO uses an EM algorithm [15] to simultaneously learn the distribution of user preference types as well as policies for each type. Note that, if we already knew the group each user belonged to, we could simply train an optimal DPO policy on each group separately. Since we do not, we think of our data as being generated by latent mixture model, where for each user we first draw a latent preference type and then draw a set of annotation data based on the preference type. We show that one can combine ideas from DPO with the EM algorithm for learning mixture models and directly learn a distribution of latent types, as well as a regularized optimal policy for each type. MinMax-DPO then takes these optimal policies and learns one model to best serve the needs of the population. Figure 1 shows the proposed pipeline.

2 Related Literature

RLHF With Diverse Preferences: One of the chief issues in RLHF is that of diverse populations; different annotators could have very different preferences [17]. Several studies have tried to solve the diverse population problem by learning more expressive reward functions and then using them to perform RLHF. For example, [39, 23, 11] maintains and learns several reward models at once. Similarly, [50] learns a multi-dimensional reward model where each dimension provides rewards based on a different objective such as safety or usefulness. [55] proposes a policy-agnostic method to perform multi-objective LLM alighment. Alternatively, [43, 27] learns a distribution over fixed reward models. Finally, these reward models are combined using various strategies [6, 23, 39] to get a final reward model which is then used to perform RLHF. [11] also learns multiple reward models, but performs RL by maximizing the minimum reward thereby ensuring that the final model is fair. The paper draws on elements of social choice theory, which [13] argues is an effective path forward for RLHF research in general, specifically regarding issues with aggregating preferences. [14] outlines a correspondence between the key principles and desiderata of social choice into the RLHF context.

In an orthogonal approach, [59] utilizes metalearning to learn diverse preferences. In general, trying to do RLHF with many reward models becomes expensive, making extending DPO [37] an attractive alternative. [47] proposes SPO to sidestep reinforcement learning using the concept of a minimax winner from social choice theory, but only in the case of homogeneous preferences. In concurrent work, [35] proposes a personalized RLHF algorithm which learns clustered policies via a hard Expectation Maximization algorithm using DPO. We instead propose a soft-clustering algorithm, which enjoys stronger theoretical guarantees [15]. [35] also proposes an algorithm to aggregate estimated reward functions for a heterogeneous population. [40] also deals with the idea of aggregating reward models to increase robustness. We instead propose a complete pipeline to learn one equitable policy for a heterogeneous population without appealing to reward model estimation at all.

DPO Generalizations: Since DPO's inception [37], there has been a growing line of literature on its generalizations, some of which we highlight here. [25] generalizes DPO to the case of multiple SFT models, while [61] generalizes to multiple objectives. [57, 36] work on extending DPO to work at the token level. [49] extends DPO to work with other types of divergence terms, while [54] relates DPO to DRO in order to robustify it. [4] augments DPO with a computable advantage function to create a hybrid between DPO and RLHF.

Additional work that more generally relates to the fields of reward modeling and preference-based reinforcement learning can be found in Appendix A.

3 Background

In this section, we discuss traditional alignment methods that assume uniform preference among the whole population, namely RLHF [63, 45, 34] and DPO [37].

Reinforcement Learning from Human Feedback (RLHF) The RLHF pipeline has two inputs. The first is a language model, denoted as π_{SFT} , which is pre-trained on internet-scale data and then fine-tuned using supervised learning. The second input is a static annotator preference dataset, $\mathcal{D} = \{x, y_w, y_l, h\}$. To collect this dataset, for a given prompt x, pairs of responses (y_1, y_2) are generated from $\pi_{\text{SFT}}(\cdot|x)$. Then, a human annotator $h \in \mathcal{H}$ (where \mathcal{H} represents the population of all annotators) selects the preferred response between



Figure 1: Proposed pipeline to find the optimal policy. Step 1: We gather binary preference data from heterogeneous annotators. Step 2: We run an expectation-maximization algorithm EM-DPO to soft assign annotators to clusters and to find an ensemble of optimal policies. Step 3: We run a regret-based algorithm Max-Min DPO to learn a linear combination of the optimal policies that is equitable.

 y_1 and y_2 . We use the notation y_w and y_l to represent the winning and losing responses, respectively.

To model the ground truth for how annotators from the previous step choose between pairs of responses, a common assumption is that the preference data is linked to a "ground truth" reward model via the Bradley-Terry-Luce model [9, 37, 34]. Let $r^*(x, y)$ represent this true reward function for all annotators. Then, according to the Bradley-Terry-Luce model, the probability that an annotator prefers one response over the other is given by:

$$p_{r^*}(y_1 \succ y_2 | x)$$
(1)
=
$$\frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))}$$

=
$$\sigma(r^*(x, y_1) - r^*(x, y_2))$$

In the first step of RLHF, a reward model $r_{\phi}(x, y)$ (parametrized by ϕ) is fit using the preference data \mathcal{D} to approximate the true reward function r^* . This is done by minimizing the following log-likelihood loss:

$$\mathcal{L}(r_{\phi}; \mathcal{D}) = -\mathbb{E}_{(x, y_1, y_2) \sim \mathcal{D}}[p_{r_{\phi}}(y_1 \succ y_2 | x)]$$
$$= -\mathbb{E}_{\mathcal{D}}[\sigma(r_{\phi}(x, y_1) - r_{\phi}(x, y_2))]$$

The second and final step is fine-tuning with reinforcement learning (RL) using the learned reward model $r_{\phi}(x, y)$. More specifically, the Proximal Policy Optimization (PPO) [41] is used in training the LM. The PPO algorithm optimizes the following objective:

$$\pi_{\theta}^{*} = \operatorname*{arg\,max}_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(y, x)] \quad (2)$$
$$-\beta \mathcal{D}_{\mathrm{KL}} [\pi_{\theta}(y|x)] ||\pi_{\mathrm{SFT}}(y|x)]$$

Direct Preference Optimization (DPO) DPO optimizes the same objective as PPO as given above in 2 but bypasses learning the reward model by directly optimizing with the preference data by combining 1 and 2, resulting in a pipeline that is significantly simpler and also exhibits greater stability [37]:

$$\mathcal{L}(\pi_{\theta}; \pi_{\text{SFT}}, \mathcal{D}, \beta) = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{SFT}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{SFT}}(y_l | x)} \right) \right] \\ \pi_{\theta}^* = \arg \min_{\pi_{\theta}} \mathcal{L}(\pi_{\theta}; \pi_{\text{SFT}}, \mathcal{D}, \beta)$$

4 DPO Extension for Diverse Annotators

Both Reinforcement Learning from Human Feedback (RLHF) and Direct Preference Optimization (DPO) assume uniform preferences across the population and learn a single reward model, either implicitly or explicitly. However, human preferences and values are inherently diverse. Consequently, RLHF and DPO tend to align with the majority opinion among annotators, introducing bias and potentially marginalizing minority perspectives. To mitigate this issue, we propose a pipeline consisting of two algorithms: Expectation-Maximization DPO (EM-DPO) for clustering diverse preference distributions and learning the optimally aligned policy for each cluster and Min-Max Regret DPO, which fairly aggregates the learned policies to minimize worst-case regret for any sub-group of annotators.

4.1 Learning an emsemble of LLMs using EM-DPO

The Expectation-Maximization Algorithm [15, 32] deals with settings with mixture data. Data are produced by first drawing a set of latent factors Z and then drawing a set of observed variables V | Z. The parameters of the likelihood determine both the distribution of the latent factors $p(Z; \theta)$ as well as the conditional likelihood $p(V | Z; \theta)$. At step t of the algorithm, we have a current candidate parameter vector θ_t and calculate θ_{t+1} as follows:

$$\begin{aligned} \theta_{t+1} &= \operatorname*{arg\,max}_{\theta} Q(\theta \mid \theta_t) \\ &:= \mathbb{E}_{Z \sim p(\cdot \mid V, \theta_t)} \left[\log(p(V, Z \mid \theta)) \right] \end{aligned}$$

In our setting, the latent factors $Z = (Z_1, \ldots, Z_n)$ correspond to the unobserved heterogeneity types Z_i of an annotator $i \in [n]$ and $V = (V_1, \ldots, V_m)$ correspond to the chosen preferences $y_1^{ij} \succeq y_2^{ij}$ for each of the prompts X_{ij} assigned to the annotator. We assume for simplicity that each annotator is assigned m prompts and we let $V_{ij} = (X_{ij}, y_1^{ij} \succeq y_2^{ij})$, where X_{ij} is the prompt and $y_1^{ij} \succeq y_2^{ij}$ is the preference for that prompt. Our parameters θ are (ϕ, ρ, η) , where ϕ are the parameters for the groupwise policies, η the latent distribution of user types and ρ are parameters that determine the distribution of prompts X_{ij} .

With some calculation, we find that a parameterization of the policy $\pi_{\phi,z}$ implies a parameterization of the likelihood (see Appendix B):

$$p(V_i \mid Z_i; \theta) = \prod_{j=1}^m \sigma_\phi(Z_i, V_{ij}) \, p(X_{ij} \mid Z_i; \rho)$$

where the function σ_{ϕ} is similar to the parameterization introduced in DPO:

$$\sigma_{\phi}(z, x, y_1, y_2)$$

:= $\sigma \left(\beta \log \frac{\pi_{\phi, z}(y_1|x)}{\pi_{\text{SFT}}(y_1|x)} - \beta \log \frac{\pi_{\phi, z}(y_2|x)}{\pi_{\text{SFT}}(y_2|x)}\right)$

Note that the latent factors take values in a set of K discrete values $\{z_1, \ldots, z_K\}$. In this case, we can assume a fully non-parametric likelihood $p(Z; \theta)$, where $\eta = p(z_k; \theta) \in \Delta(K)$, the K-dimensional simplex. Subsequently, we can decompose the criterion as:

$$Q(\theta \mid \theta_t)$$

= $\mathbb{E}_{Z \sim p(\cdot \mid V, \theta_t)} \left[\sum_{i=1}^n \log(p(V_i, Z_i \mid \theta)) \right]$
= $\mathbb{E}_{Z \sim p(\cdot \mid V, \theta_t)} \left[\sum_{i=1}^n \log(p(V_i \mid Z_i; \theta)) + \log(p(Z_i; \theta)) \right]$

For further simplification, we note that $p(Z_i; \theta) = \sum_{k=1}^{K} \eta_k \mathbf{1}\{Z_i = z_k\}$. Further, assuming that $p(V_i \mid Z_i; \theta)$ does not depend on the vector η , so that $p(V_i \mid Z_i; \theta) = p(V_i \mid Z_i; \phi, \rho)$, the original criterion decomposes into two separate optimization problems:

$$\eta_{t+1} = \arg \max_{\eta} \mathbb{E}_{Z \sim p(\cdot | V, \theta_t)} \left[\sum_{i=1}^n \log \left(\sum_{k=1}^K \eta_k 1\{Z_i = z_k\} \right) \right]$$

$$\phi_{t+1} = \arg \max_{\phi, \rho} \mathbb{E}_{Z \sim p(\cdot | V, \theta_t)} \left[\sum_{i=1}^n \log(p(V_i | Z_i; \phi, \rho)) \right]$$

For the *E*-step, we must characterize the posterior distribution of the latent factors. Under the assumption that the contexts are un-correlated with the unobserved preference types, which is natural in the context of LLM fine-tuning, since contexts are randomly assigned to annotators, we can derive that (see Appendix C):

$$p(z_k \mid V_i; \theta) = \frac{\eta_k \prod_{j=1}^m \sigma_\phi(z_k, V_{ij})}{\sum_{\ell=1}^K \eta_\ell \prod_{j=1}^m \sigma_\phi(z_\ell, V_{ij})}$$

For the *M*-step, we must solve the two optimization problems given above. The solution for η can be derived in closed form, while the solution for ϕ is independent of the term $p(X_{ij} | Z_i; \rho)$:

$$\eta_{k,t+1} = \frac{1}{n} \sum_{i=1}^{n} p(z_k \mid V_i; \theta_t)$$

$$\phi_{t+1} = \arg\max_{\phi} \sum_{i=1}^{n} \mathbb{E}_{Z_i \sim p(\cdot \mid V_i; \theta_t)} \left[\sum_{j=1}^{m} \log(\sigma_{\phi}(Z_i, V_{ij})) \right]$$

A full derivation is in Appendix D. This gives rise to the following EM algorithm:

Algorithm 1 EM-DPO: Expectation-Maximization Direct Preference Optimization

- 1: **Input:** Preference data \mathcal{D} indexed for all human annotators \mathcal{I} and containing m_i demonstrations for each human annotator *i*.
- 2: **Input:** pre-trained group-wise models $\pi_{\phi_0,z}$; $\forall z \in \{z_1, \dots, z_k\}$.
- 3: Initialize $\eta_0 = (1/K, ..., 1/K)$
- 4: for t in $\{0, ..., T\}$ do
- 5: **E.** Calculate posterior $p(z_k | V_i; \theta_t)$ for each annotator $i \in \mathcal{I}$:

$$\gamma_{i,k} = \frac{\eta_{k,t} \prod_{j=1}^{m_i} \sigma_{\phi_t}(z_k, V_{ij})}{\sum_{\ell=1}^{K} \eta_{\ell,t} \prod_{j=1}^{m_i} \sigma_{\phi_t}(z_\ell, V_{ij})}$$

6: **M.** Update parameters ϕ , η :

$$\eta_{k,t+1} = \frac{\sum_{i \in \mathcal{I}} \gamma_{i,k}}{|\mathcal{H}|}$$

$$\phi_{t+1} = \arg\max_{\phi} \sum_{i \in \mathcal{I}} \sum_{k=1}^{K} \gamma_{i,k} \sum_{j=1}^{m_i} \log(\sigma_{\phi}(z_k, V_{ij}$$

7: end for

8: **Return:** Policies $\{\pi_{\phi_t,z} : z \in \{z_1, \ldots, z_k\}\}$ and posterior preference weights $\{\gamma_{i,k} : i \in \mathcal{I}\}$.

Note that if we do not share parameters across the policies for each preference type z, i.e. we have separate parameters ϕ_z for each $z \in \{z_1, \ldots, z_K\}$, then the optimization in the final step of EM-DPO also decomposes into separate policy optimization problems for each preference type:

$$\phi_{z_k,t+1} = \arg\min_{\phi_{z_k}} \sum_{i \in \mathcal{I}} \sum_{j=1}^{m_i} \gamma_{i,k} \log(\sigma_{\phi}(z_k, V_{ij}))$$

Note that the latter is simply a weighted DPO problem, where each demonstration (h, j), which corresponds to the *j*-th demonstrations from annotator *i*, is assigned weight $\gamma_{i,k}$ when optimizing the policy parameters for preference type z_k . Alternatively, for multi-tasking purposes, some parameters can be shared parameters across policies for each preference type, in which case the final optimization problem should be solved simultaneously via stochastic gradient descent over the joint parameters ϕ .

4.2 Fair aggregation of LLMs using Min-Max DPO

So far, we have shown how to calculate a separate policy that optimizes for each preference population z. Our ultimate goal is to output a single policy. Hence, we need to trade-off optimizing for the preferences of different groups and find a policy that strikes a good balance.

In that respect, to equitably cater to all K subpopulations, we focus on identifying a policy that minimizes the worst-case regret among the subpopulations. To avoid having to retrain a new policy, we will restrict ourselves to selecting an ensemble among the already trained policies. As such, we define the ensemble space of policies as:

$$\Pi = \left\{ \sum_{k=1}^{K} w_k \pi_{\phi, z_k} : w \in \Delta(K) \right\}$$

)) If we had access to the reward functions $r_z^*(y, x)$, then for any policy π , the expected reward that population z receives would be:

$$R_z(\pi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(\cdot|x)} \left[r_z^*(y, x) \right]$$

Note that if we were to solely focus on population z, we would be optimizing the expected reward objective above, regularized so as not to deviate from the reference policy. This would yield policy $\pi_z^* = \pi_{\phi^*,z}$, where ϕ_* are the policy parameters we calculated based on the EM-DPO algorithm.

Our goal is to find an ensemble policy π such that no population z has very large regret towards choosing their population-preferred policy π_z^* . Our minimax regret optimization problem can be simply stated as:

$$\pi_* = \arg\min_{\pi \in \Pi} \max_{k=1}^{K} \left[R_{z_k}(\pi_{z_k}^*) - R_z z_k \pi) \right]^+$$

where $[x]^+ = \max\{x, 0\}$. Note that we only consider the positive part of the regret.

Why min-max regret? Max-min reward is another fairness criterion that can be applied to the RLHF problem to ensure equity, as discussed in [11]. However, this criterion has two major drawbacks. Firstly, the reward model is not uniquely identifiable from preference data. Two reward models r(x, y) and r'(x, y) are equivalent if r(x, y) - r'(x, y) = f(x) [37]. Therefore, directly maximizing the minimum reward is ineffective due to this scaling. We could fix this by standardizing the reward model to set the minimum reward to zero - if r(x, y) is the recovered reward function, we can use $r'(x, y) = r(x, y) - \min_{y} r(x, y)$, which is an equivalent reward model. Even then, there is another issue with the max-min reward criterion: The max-min reward focuses on improving rewards for users with the lowest reward, while the min-max regret function targets users with the highest regrets. These groups differ when users with low rewards also have low regrets. As an example, consider a setting with fixed context and three responses. If two users have reward vectors [0, 0.01,(0.02) and [0, 10, 1] respectively, then the max-min reward objective will choose response 3 to maximize user 2's reward. However, user 1 is nearly indifferent between the three choices 1, whereas user 2 strongly prefers option 2. Therefore, it is more ideal to choose option 2, which the min-max regret criteria chooses.

We now show that the min-max regret objective can also be optimized over, without access to the explicit reward functions, but solely based on the policies we have already trained. We can rewrite our objective as (see Appendix E):

$$\min_{w \in \Delta(K)} \max_{z \in \{z_0, z_1, \dots, z_K\}} \sum_{k=1}^K w_k \cdot \left(\mathcal{L}_{z, z} - \mathcal{L}_{z, z_k}\right),$$

where

$$\mathcal{L}_{z,z'} := \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{z'}^*(\cdot|x)} \left[\log \left(\frac{\pi_z^*(y|x)}{\pi_{\text{SFT}}(y|x)} \right) \right].$$

Letting \mathcal{R} denote the $(K + 1) \times K$ matrix whose (k, k') entry (for $0 \leq k \leq K, 1 \leq k' \leq K$) corresponds to $\mathcal{R}_{k,k'} := \mathcal{L}_{z_k, z_k} - \mathcal{L}_{z_k, z_{k'}}$, we can rewrite the above objective as:

$$\min_{w \in \Delta(K)} \max_{p \in \Delta(K+1)} p^{\top} \mathcal{R} w$$

This is simply a finite action zero-sum game, where the minimizing player has K actions and the maximizing player has K + 1 actions. A large variety of methods can be utilized to calculate an equilibrium of this zero-sum game and hence identify the minimax regret optimal mixture weights w_* . For instance, we can employ optimistic Hedge vs. optimistic Hedge dynamics, which are known to achieve fast convergence rates in such finite action zero-sum games [38] and then use the average of the solutions over the iterates of training, as described in Algorithm 2. Algorithm 2 MinMax-DPO: Direct Optimization for Min-Max Regret Ensemble

- 1: **Input:** Distribution \mathcal{D} of contexts x.
- 2: **Input:** Population-specific optimal policies $\pi_z^* \equiv \pi_{\phi_*,z}$ returned from EM-DPO
- Input: Number of iterations T and a sufficiently small, albeit constant, independent of T, step-size η
- 4: Calculate discrepancies for $z, z' \in \{z_1, \ldots, z_k\}$:

$$\mathcal{L}_{z,z'} := \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{z'}^*(\cdot|x)} \left[\log \left(\frac{\pi_z^*(y|x)}{\pi_{\text{SFT}}(y|x)} \right) \right]$$

with the convention that $\mathcal{L}_{z_0,z_0} = \mathcal{L}_{z_0,z_k} = 0$

5: Calculate $(K+1) \times K$ regret matrix \mathcal{R} , whose $k \in \{0, \dots, K\}$ and $k' \in \{1, \dots, K\}$ entry is:

$$\mathcal{R}_{k,k'} := w'_k \left(\mathcal{L}_{z_k, z_k} - \mathcal{L}_{z_k, z'_k} \right)$$

6: Initialize $w_0 = (1/K, \dots, 1/K)$ and $p_0 = (1/(K+1), \dots, 1/(K+1))$ 7: for t in $\{0, \dots, T\}$ do

$$w_{t} \propto w_{t-1} \exp\left\{-\eta \cdot \left(2\mathcal{R}^{\top} p_{t-1} - \mathcal{R}^{\top} p_{t-2}\right)\right\}$$
$$p_{t} \propto p_{t-1} \exp\left\{\eta \cdot \left(2\mathcal{R} w_{t-1} - \mathcal{R} w_{t-2}\right)\right\}$$

8: end for

9: **Return:** Policy $\pi_* = \sum_{k=1}^K w_k^* \pi_{\phi_*, z_k}$, where $w^* = \frac{1}{T} \sum_{t=1}^T w_{k,t}$

The solution π_* returned by Algorithm 2 constitutes a $O(\log(K) \log(T) T^{-1})$ -approximate solution to the min-max regret problem (a direct consequence of the results in [38]). This completes our overall direct preference optimization procedure with unobserved heterogeneous preferences.

One can also optimize a new policy π that does not correspond to an ensemble of the base policies $\pi_{\phi_*,z}$ by solving the saddle point problem:

$$\min_{\pi} \max_{z} L_z(\pi) \triangleq R_z(\pi) - R_z(\pi_z^*)$$

which has already been shown, can be expressed as a function of π_z^* and π_{SFT} . This saddle point can be solved by policy-gradient vs multiplicative weight dynamics, or for faster convergence via optimistic policy gradient descent vs optimistic multiplicative weight dynamics:

$$\phi_{t+1} = \phi_t - 2\nabla_{\phi} \sum_{z} p_{t,z} L_z(\pi_{\phi_t}) + \nabla_{\phi} \sum_{z} p_{t-1,z} L_z(\pi_{\phi_t-1}) p_{t+1,z} \propto p_{t,z} \exp\{\eta \cdot (2L_z(\pi_{\phi_t}) - L_z(\pi_{\phi_{t-1}}))\}$$

5 Experiments

5.1 Multi-Armed Bandit Experiment

5.1.1 Settings

As a warm-up exercise, we consider the multiarmed bandit setting [3]. In this setting, the context represents the prompt and the bandit arms represent the possible responses for the given context. For our experiment, we consider a simplified case with three arms and we model the sub-population reward model using a linear function, similarly to linear contextual bandits [16]:

$$r_z^*(y,x) = x^T \theta_z(y) + \mathcal{N}(0,\sigma)$$

where $\theta_z(y)$ is the model parameters corresponding to the latent variable z for arm y, x is the context, and $\mathcal{N}(0,\sigma)$ represents noise with 0 mean and standard deviation $\sigma = 0.01$. We generate 200 annotators drawn randomly from three subpopulations with probabilities 0.6, 0.3 and 0.1 respectively. The preferences of annotators within each sub-population is homogeneous and therefore, each sub-population is associated with a single reward model. We fix $\theta_z(y)$ for any given subpopulation with values $\theta_i(i) = [10, 10, 10]$ and $\theta_i(j) = [0, 0, 0], j \neq i$ for sub-population i. We generate 10 preference data pairs per annotator. For each data point, first we draw a context vector x uniformly randomly from the hypercube $[0, 1]^3$. Then, a pair of responses is generated from a uniformly random reference policy $\pi_{\text{SFT}} = (1/3, 1/3, 1/3).$ The annotator then chooses a response $y_w \succ y_l$ based on their reward model $r_z^*(y, x)$. We implement EM-DPO and MinMax-DPO for this data. Appendix F shows hyperparameters for the experiment.

5.1.2 Results

We run standard DPO and MinMax-DPO on this experimental setup and calculate the average regret



Figure 2: DPO vs. MinMax-DPO Regret Plot



Figure 3: Convergence of learned weights in MinMax-DPO.

per user group. The results of this are shown in Figure 2.

We can see that training DPO over the whole population leads to the policy completely optimizing for the first user group's preference (i.e., majority opinion), leading to maximal regret for the other two groups. However, MinMax-DPO achieves the social optimum and respects the preferences of all three groups, as shown by the equal regret among all three groups. Figure 3 shows the convergence of the learned weights in the MinMax-DPO algorithm; we see relatively quick convergence to the optimal weights, which are close to uniform. This is expected as all three sub-groups have perfectly contradicting opinions because they each prefer a different response.

5.2 IMDb Movie Reviews

5.2.1 Settings

We conduct our experiment on Mistral 7B v3.0 [24]. To construct our preference dataset, we use a

publicly available synthetic preference data generated from IMDb reviews [30] (see Appendix F.1 for details). Our dataset consists of 60,000 preference pairs over IMDb movie review completions, randomly and evenly distributed among 1,200 users (50 pairs per user). Of these users, 66.66% (group 0) prefer reviews with greater positive sentiment, while the remaining 33.33% (group 1) favor grammatically correct reviews.

5.2.2 Algorithms

As before, we compare the performance of our pipeline (EM-DPO plus MinMax-DPO) with performing regular DPO on the full dataset. We implement an additional benchmark: Cluster DPO. In this algorithm we perform clustering based on the chosen response (more details in Appendix [?]). Once the clusters are determined we perform DPO on each of the individual clusters and combine them using our aggregation algorithm MinMax DPO. We investigate the quality of both algorithm both after the clustering stage and after the aggregation stage.

5.2.3 Metrics

Accuracy: Here, we define accuracy as the percentage of data points (x, y_1, y_2) , where x is the prompt, y_1 is the chosen response, and y_2 is the rejected response, such that

$$\log \frac{\pi_{\phi,z}(y_1|x)}{\pi_{\text{SFT}}(y_1|x)} > \log \frac{\pi_{\phi,z}(y_2|x)}{\pi_{\text{SFT}}(y_2|x)},$$

or equivalently, the percentage of data points where the chosen response is given a higher "reward" than the rejected one.

Reward Margins: Reward margins as defined as:

$$r_m(x, y_1, y_2) = \beta \log \frac{\pi_{\phi, z}(y_1|x)}{\pi_{\text{SFT}}(y_1|x)} - \beta \log \frac{\pi_{\phi, z}(y_2|x)}{\pi_{\text{SFT}}(y_2|x)}$$

on the evaluation dataset for each of the five policies, i.e. the average difference between the rewards of the chosen and the rejected responses.

5.3 Results

We report the margins in the main text and report the accuracy metrics in Appendix [?]. Figure 4 shows the policies learned after the clustering step (i.e. EM-DPO and k means clustering then DPO), we also report the regular DPO metrics for comparison. EM-DPO clearly learns clusters of higher



Figure 4: Reward margins after the clustering step



Figure 5: Reward margins after the aggregation step

quality than k means clustering, both of which are of higher quality than regular DPO. In figure Figure 5, we report the margins after aggregation. In addition to just reporting the results of MinMax-DPO we also report a weighted version of the results. The tuple in the x-axis, for e.g. (1,4), gives the final clustering results of the policies by weighing the regret of each sub-group, i.e. in this case, the grammar groups's regret stays the same while the sentiment group's regrets are multiplied by 4. This allows us to define a "weighted" version of fairness, where we have the flexibility to control the priorities of each group differently. Naturally, for different possible weights, we get better results if the clusters are of higher quality. We see that EM-DPO out performs both the regular DPO as well as k means clustering followed by DPO in all cases.

5.4 Discussion & Limitations

We provide a robust framework to train equitable policies for a heterogeneous population with diverse preferences. By extending the DPO algorithm, we are able to sidestep reinforcement learning entirely, enjoying the added stability that DPO provides while making it more applicable to realworld situations and datasets. We demonstrate our findings on a contextual bandit experiment as well as a larger-scale LLM experiment, showing how our algorithm, MinMax-DPO, generates a far more socially equitable policy than standard DPO in diverse populations where some groups may be underrepresented.

Based on our results, we raise some limitations and directions for future work. Our derivations operate off of the assumption that contexts are uncorrelated given the preference type of the annotator; this may not be the case in the real world as, to increase accuracy of data collection, annotators may be given prompts more tuned to their skill sets. We also assume annotators report their preferences honestly, which may not be the case - this raises important questions regarding incentive compatibility.

References

- [1] Youssef Abdelkareem, Shady Shehata, and Fakhri Karray. Advances in preference-based reinforcement learning: A review. In 2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC), pages 2527–2532. IEEE, 2022.
- [2] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021.
- [3] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR, 2024.
- [4] Anirudhan Badrinath, Prabhat Agarwal, and Jiajing Xu. Hybrid preference optimization: Augmenting direct preference optimization with auxiliary objectives. *arXiv* preprint arXiv:2405.17956, 2024.
- [5] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

- [6] Michiel Bakker, Martin Chadwick, Hannah Sheahan, Michael Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matt Botvinick, et al. Fine-tuning language models to find agreement among humans with diverse preferences. Advances in Neural Information Processing Systems, 35:38176–38189, 2022.
- [7] Andreea Bobu, Andi Peng, Pulkit Agrawal, Julie Shah, and Anca D Dragan. Aligning robot and human representations. *arXiv preprint arXiv:2302.01928*, 2023.
- [8] Michael Bowling, John D Martin, David Abel, and Will Dabney. Settling the reward hypothesis. In *International Conference on Machine Learning*, pages 3003–3020. PMLR, 2023.
- [9] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [10] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- [11] Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Amrit Singh Bedi, and Mengdi Wang. Maxmin-rlhf: Towards equitable alignment of large language models with diverse human preferences. *arXiv preprint arXiv:2402.08925*, 2024.
- [12] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [13] Vincent Conitzer, Rachel Freedman, Jobst Heitzig, Wesley H Holliday, Bob M Jacobs, Nathan Lambert, Milan Mossé, Eric Pacuit, Stuart Russell, Hailey Schoelkopf, et al. Social choice for ai alignment: Dealing with diverse human feedback. *arXiv preprint arXiv:2404.10271*, 2024.
- [14] Jessica Dai and Eve Fleisig. Mapping social choice theory to rlhf. *arXiv preprint arXiv:2404.13038*, 2024.
- [15] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977.
- [16] Maria Dimakopoulou, Zhengyuan Zhou, Susan Athey, and Guido Imbens. Balanced linear contextual bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3445–3453, 2019.
- [17] Vincent Dumoulin, Daniel D Johnson, Pablo Samuel Castro, Hugo Larochelle, and Yann Dauphin. A density estimation perspective on learning from pairwise human preferences. *arXiv preprint arXiv:2311.14115*, 2023.

- [18] Owain Evans, Andreas Stuhlmüller, and Noah Goodman. Learning the preferences of ignorant, inconsistent agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- [19] Justin Fu, Anoop Korattikara, Sergey Levine, and Sergio Guadarrama. From language to goals: Inverse reinforcement learning for vision-based instruction following. *arXiv preprint arXiv:1902.07742*, 2019.
- [20] Joey Hong, Kush Bhatia, and Anca Dragan. On the sensitivity of reward inference to misspecified human models. *arXiv preprint arXiv:2212.04717*, 2022.
- [21] Minyoung Hwang, Gunmin Lee, Hogun Kee, Chan Woo Kim, Kyungjae Lee, and Songhwai Oh. Sequential preference ranking for efficient reinforcement learning from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- [22] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning from human preferences and demonstrations in atari. *Advances in neural information processing systems*, 31, 2018.
- [23] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *arXiv preprint arXiv:2310.11564*, 2023.
- [24] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [25] Hung Le, Quan Tran, Dung Nguyen, Kien Do, Saloni Mittal, Kelechi Ogueji, and Svetha Venkatesh. Multi-reference preference optimization for large language models. *arXiv preprint arXiv:2405.16388*, 2024.
- [26] Kimin Lee, Laura Smith, Anca Dragan, and Pieter Abbeel. B-pref: Benchmarking preference-based reinforcement learning. *arXiv preprint arXiv:2111.03026*, 2021.
- [27] Dexun Li, Cong Zhang, Kuicai Dong, Derrick Goh Xin Deik, Ruiming Tang, and Yong Liu. Aligning crowd feedback via distributional preference reward modeling. *arXiv preprint arXiv:2402.09764*, 2024.
- [28] David Lindner and Mennatallah El-Assady. Humans are not boltzmann distributions: Challenges and opportunities for modelling human feedback and interaction in reinforcement learning. *arXiv preprint arXiv:2206.13316*, 2022.
- [29] Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

- [30] Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150, 2011.
- [31] Anirudha Majumdar, Sumeet Singh, Ajay Mandlekar, and Marco Pavone. Risk-sensitive inverse reinforcement learning via coherent risk models. In *Robotics: science and systems*, volume 16, page 117, 2017.
- [32] Todd K Moon. The expectation-maximization algorithm. *IEEE Signal processing magazine*, 13(6):47–60, 1996.
- [33] Khanh Nguyen, Hal Daumé III, and Jordan Boyd-Graber. Reinforcement learning for bandit neural machine translation with simulated human feedback. *arXiv preprint arXiv:1707.07402*, 2017.
- [34] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [35] Chanwoo Park, Mingyang Liu, Kaiqing Zhang, and Asuman Ozdaglar. Principled rlhf from heterogeneous feedback via personalization and preference aggregation. *arXiv preprint arXiv:2405.00254*, 2024.
- [36] Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to q*: Your language model is secretly a q-function. *arXiv preprint arXiv:2404.12358*, 2024.
- [37] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [38] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. *Advances in Neural Information Processing Systems*, 26, 2013.
- [39] Alexandre Rame, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpolating weights finetuned on diverse rewards. *Advances in Neural Information Processing Systems*, 36, 2024.
- [40] Alexandre Ramé, Nino Vieillard, Léonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier Bachem, and Johan Ferret. Warm: On the benefits of weight averaged reward models. *arXiv preprint arXiv:2401.12187*, 2024.
- [41] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- [42] Rohin Shah, Noah Gundotra, Pieter Abbeel, and Anca Dragan. On the feasibility of learning, rather than assuming, human biases for reward inference. In *International Conference on Machine Learning*, pages 5670–5679. PMLR, 2019.
- [43] Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learning: Understanding and accounting for hidden context in rlhf. *arXiv preprint arXiv:2312.08358*, 2023.
- [44] Joar Max Viktor Skalse and Alessandro Abate. The reward hypothesis is false. 2022.
- [45] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- [46] Haoran Sun, Yurong Chen, Siwei Wang, Wei Chen, and Xiaotie Deng. Mechanism design for llm finetuning with multiple reward models. *arXiv preprint arXiv:2405.16276*, 2024.
- [47] Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaximalist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*, 2024.
- [48] Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, et al. Secrets of rlhf in large language models part ii: Reward modeling. *arXiv preprint arXiv:2401.06080*, 2024.
- [49] Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints. *arXiv preprint arXiv:2309.16240*, 2023.
- [50] Haoxiang Wang, Yong Lin, Wei Xiong, Rui Yang, Shizhe Diao, Shuang Qiu, Han Zhao, and Tong Zhang. Arithmetic control of llms for diverse user preferences: Directional preference alignment with multi-objective rewards. *arXiv preprint arXiv:2402.18571*, 2024.
- [51] Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*, 2023.
- [52] Christian Wirth, Riad Akrour, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-based reinforcement learning methods. *Journal of Machine Learning Research*, 18(136):1–46, 2017.
- [53] Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*, 2021.

[54] Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jiawei Chen, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. Towards robust alignment of language models: Distributionally robustifying direct preference optimization. *arXiv preprint arXiv:2407.07880*, 2024.

- [55] Kailai Yang, Zhiwei Liu, Qianqian Xie, Tianlin Zhang, Nirui Song, Jimin Huang, Ziyan Kuang, and Sophia Ananiadou. Metaaligner: Conditional weak-to-strong correction for generalizable multiobjective alignment of language models. *arXiv preprint arXiv:2403.17141*, 2024.
- [56] Yueqin Yin, Zhendong Wang, Yi Gu, Hai Huang, Weizhu Chen, and Mingyuan Zhou. Relative preference optimization: Enhancing llm alignment through contrasting responses across identical and diverse prompts. *arXiv preprint arXiv:2402.10958*, 2024.
- [57] Yongcheng Zeng, Guoqing Liu, Weiyu Ma, Ning Yang, Haifeng Zhang, and Jun Wang. Tokenlevel direct preference optimization. *arXiv preprint arXiv:2404.11999*, 2024.
- [58] Jiangchuan Zheng, Siyuan Liu, and Lionel M Ni. Robust bayesian inverse reinforcement learning with sparse behavior noise. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.
- [59] Huiying Zhong, Zhun Deng, Weijie J Su, Zhiwei Steven Wu, and Linjun Zhang. Provable multi-party reinforcement learning with diverse human feedback. *arXiv preprint arXiv:2403.05006*, 2024.
- [60] Li Zhou and Kevin Small. Inverse reinforcement learning with natural language goals. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 35, pages 11116–11124, 2021.
- [61] Zhanhui Zhou, Jie Liu, Jing Shao, Xiangyu Yue, Chao Yang, Wanli Ouyang, and Yu Qiao. Beyond one-preference-fits-all alignment: Multi-objective direct preference optimization. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10586–10613, 2024.
- [62] Zhanhui Zhou, Jie Liu, Chao Yang, Jing Shao, Yu Liu, Xiangyu Yue, Wanli Ouyang, and Yu Qiao. Beyond one-preference-fits-all alignment: Multi-objective direct preference optimization. *arXiv preprint arXiv:2310.03708*, 2023.
- [63] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

Appendix A Additional Related Work

Preference-Based Reinforcement Learning: Reinforcement learning from preferences has been an active research area for some time, providing a way to train on tasks for which explicitly defining rewards is hard [52, 26, 1]. In particular, [12, 22] show that using human preferences to guide reinforcement learning (RLHF) is particularly effective on a variety of tasks, such as training robots. More recently, RLHF has become a very popular technique to fine-tune language models to do a variety of tasks such as summarization [34, 63, 45, 53]. RLHF has also been used to align language models [5, 2]. [10] details several open problems in the field of RLHF, including those related to the feedback itself, particularly the inverse relation between richness and efficiency. Some work has been done on this problem with regards to language-based feedback in particular [19, 60] as well as in more general settings [21], but specific applications to LLMs have not been fully explored.

Challenges with Reward Modeling: In general, human preferences can be difficult to represent using reward models [20], and the validity of reward modeling itself is still somewhat debated [8, 7, 44]. Some work has also been done to take personality into account when reward modeling [28, 26], but this area remains open. In general, taking human irrationality into account when reward modeling (to optimize a more accurate reward function) leads to a trade-off between efficiency and accuracy [42, 33]. Work has been done on inverse RL with particular models of suboptimality such as myopia [18], noise [58], and risk-sensitivity [31], but dealing with general irrationalities remains open.

The proper use and collection of data remains an issue with RLHF. [46] analyzes LLM fine-tuning as a mechanism design problem where agents may have the incentive to misreport their preferences. Data can also often have issues or certain data points may not be as effective as others; [48] proposes methods to deal with incorrect or ambiguous preference pairs, while [56] proposes an extension to DPO which uses contrastive learning to discern between more and less preferred responses to prompts.

Appendix B Likelihood Parameterization

Note that, in our situation, the latent factors and observed variables (Z_i, V_i) are independent across

annotators and therefore, the likelihood and the prior factorizes across the annotators. Moreover, conditional on the latent factor, the V_{ij} are independently distributed across j and for each j the conditional likelihood takes a logistic form, as follows:

$$p(V_i \mid Z_i; \theta)$$

$$= \prod_{j=1}^m p(V_{ij} \mid Z_i; \theta)$$

$$= \prod_{j=1}^m p(y_1^{ij} \succeq y_2^{ij}, X_{ij} \mid Z_i; \theta)$$

$$= \prod_{j=1}^m p(y_1^j \succeq y_2^{ij} \mid X_{ij}, Z_i; \theta) p(X_{ij} \mid Z_i; \theta)$$

$$= \prod_{j=1}^m \sigma \left(r^* \left(Z, X_{ij}, y_1^{ij} \right) - r^* \left(Z, X_{ij}, y_2^{ij} \right) \right) p(X_{ij} \mid Z_i; \theta)$$

where r^* denotes the true reward for the annotator, as in Section 3.

The first part $\sigma\left(r^*(Z, X_j, y_1^j) - r^*(Z, X_j, y_2^j)\right)$ can also be written in closed form in terms of the policy parameters $\pi_{\phi^*, z}$ for each preference type as designated by the same observation as in [37]:

$$\sigma(r^*(z, x, y_1) - r^*(z, x, y_2)) = \sigma\left(\beta \log \frac{\pi_{\phi^*, z}(y_1|x)}{\pi_{\text{SFT}}(y_1|x)} - \beta \log \frac{\pi_{\phi^*, z}(y_2|x)}{\pi_{\text{SFT}}(y_2|x)}\right)$$

where $\pi_{\phi^*,z}$ optimizes the type specific regularized objective:

$$\pi_{\phi^*,z} = \underset{\pi}{\arg\max} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [r^*(z, x, y)] - \beta \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [\mathcal{D}_{\mathrm{KL}}[\pi(y|x)| |\pi_{\mathrm{SFT}}(y|x)]]$$

We will introduce the shorthand notation:

$$\sigma_{\phi}(z, x, y_1, y_2)$$

:= $\sigma \left(\beta \log \frac{\pi_{\phi, z}(y_1|x)}{\pi_{\text{SFT}}(y_1|x)} - \beta \log \frac{\pi_{\phi, z}(y_2|x)}{\pi_{\text{SFT}}(y_2|x)}\right)$

Thus a parameterization of the policy space $\pi_{\phi,z}$, implies a parameterization of the likelihood:

$$p(V_i \mid Z_i; \theta) = \prod_{j=1}^m \sigma_\phi(Z_i, V_{ij}) \, p(X_{ij} \mid Z_i; \theta),$$

as desired.

Appendix C *E*-**Step Derivation**

Here, we derive the posterior distribution $p(Z \mid V; \theta) = \prod_{i=1}^{n} p(Z_i \mid V_i; \theta)$ for any given parameter θ . We apply Bayes rule:

$$p(z_k | V_i; \theta)$$

$$= \frac{p(V_i, z_k; \theta)}{p(V_i; \theta)}$$

$$= \frac{p(V_i | z_k; \theta) p(z_k; \theta)}{\sum_{\ell=1}^{K} p(V_i | z_\ell; \theta) p(z_\ell; \theta)}$$

$$= \frac{p(V_i | z_k; \phi) \eta_k}{\sum_{\ell=1}^{K} p(V_i | z_\ell; \phi) \eta_\ell}$$

$$= \frac{\prod_{j=1}^{m} \sigma_{\phi}(z_k, V_{ij}) p(X_{ij} | z_\ell; \theta) \eta_k}{\sum_{\ell=1}^{K} \prod_{j=1}^{m} \sigma_{\phi}(z_\ell, V_{ij}) p(X_{ij} | z_\ell; \theta) \eta_\ell}.$$

In the context of LLMs, the quantity X_{ij} is the prompt and the prompts are randomly assigned to annotators, so we would expect no correlation between the preference type of the annotator and the prompt assigned to them. Thus, all prompts are equally likely given the preference type of the annotator. Hence, we make the following assumption: **ASSUMPTION 1** (Un-correlated Contexts and La-

tent Preference Types). For all $k, \ell \in [K]$:

$$p(X_{ij} \mid Z_i = z_k; \theta) = p(X_{ij} \mid Z_i = z_\ell; \theta)$$
$$:= \rho(X_{ij})$$

Based on this assumption, we can then write:

$$p(z_k \mid V_i; \theta) = \frac{\prod_{j=1}^m \sigma_\phi(z_k, V_{ij})\rho(X_{ij})\eta_k}{\sum_{\ell=1}^K \prod_{j=1}^m \sigma_\phi(z_\ell, V_{ij})\rho(X_{ij})\eta_\ell}$$
(3)

Note that we can write:

$$\sum_{\ell=1}^{K} \prod_{j=1}^{m} \sigma_{\phi}(z_{\ell}, V_{ij}) \rho(X_{ij}) \eta_{\ell}$$
$$= \sum_{\ell=1}^{K} \prod_{j=1}^{m} \rho(X_{ij}) \cdot \prod_{j=1}^{m} \sigma_{\phi}(z_{\ell}, V_{ij}) \eta_{\ell}$$
$$= \prod_{j=1}^{m} \rho(X_{ij}) \cdot \sum_{\ell=1}^{K} \prod_{j=1}^{m} \sigma_{\phi}(z_{\ell}, V_{ij}) \eta_{\ell}$$

Thus, the terms $\prod_{j=1}^{m} \rho(X_j)$ cancel from the numerator and denominator in Equation (3), leading to the simplified formula that is independent of π :

$$p(z_k \mid V_i; \theta) = \frac{\eta_k \prod_{j=1}^m \sigma_\phi(z_k, V_j)}{\sum_{\ell=1}^K \eta_\ell \prod_{j=1}^m \sigma_\phi(z_\ell, V_j)}$$

Appendix D *M*-**Step Derivation**

We aim to solve the following two optimization problems:

$$\eta_{t+1} = \arg \max_{\eta} \mathbb{E}_{Z \sim p(\cdot | V, \theta_t)} \left[\sum_{i=1}^n \log \left(\sum_{k=1}^K \eta_k 1\{Z_i = z_k\} \right) \right]$$

$$\phi_{t+1} = \arg \max_{\phi, \rho} \mathbb{E}_{Z \sim p(\cdot | V, \theta_t)} \left[\sum_{i=1}^n \log(p(V_i | Z_i; \phi, \rho)) \right]$$

(4)

The first optimization problem in Equation (4) admits a closed-form solution. Letting $w_{k,t} = \sum_{i=1}^{n} p(z_k \mid V_i; \theta_t)$

$$\mathbb{E}_{Z \sim p(\cdot | V, \theta_t)} \left[\sum_{i=1}^n \log \left(\sum_{k=1}^K \eta_k 1\{Z_i = z_k\} \right) \right]$$
$$= \sum_{i=1}^n \sum_{k=1}^K p(z_k | V_i; \theta_t) \log(\eta_k)$$
$$= \sum_{k=1}^K w_{k,t} \log(\eta_k)$$

Thus the optimization problem that determines η_{t+1} takes the simple form $\max_{\eta \in \Delta(K)} \sum_{k=1}^{K} w_{k,t} \log(\eta_k)$. The Lagrangian of this problem is $L(\eta, w_t, \lambda) = \sum_{k=1}^{K} w_{k,t} \log(\eta_k) + \lambda^T (\eta - 1)$. The KKT condition is:

$$\frac{w_{k,t}}{\eta_{k,t+1}} = \lambda$$

$$\implies \eta_{k,t+1} \propto w_{k,t}$$

$$\implies \eta_{k,t+1} = \frac{w_{k,t}}{\sum_k w_{k,t}}$$

Moreover, since $\sum_k p(z_k \mid V_i; \theta_t) = 1$, we have $\sum_k w_{k,t} = n$. Thus, the above simplifies to:

$$\eta_{k,t+1} = \frac{1}{n} w_{k,t}$$
$$= \frac{1}{n} \sum_{i=1}^{n} p(z_k \mid V_i; \theta_t)$$

For the second optimization problem in Equation (4), we further decompose the objective:

$$\log(p(V_i \mid Z_i; \phi, \rho))$$

=
$$\sum_{j=1}^{m} \log(p(V_{ij} \mid Z_i; \phi, \rho))$$

=
$$\sum_{j=1}^{m} \log(\sigma_{\phi}(Z_i, V_{ij})) + p(X_{ij} \mid Z_i; \rho)$$

Assuming that the parameter ρ that determines that $p(X \mid Z; \rho)$, according to Assumption 1 is not subject to joint constraints with the parameter ϕ , we can drop the second part in the objective, when optimizing for ϕ :

$$\varphi_{t+1} = \arg \max_{\phi} \mathbb{E}_{Z \sim p(\cdot|V,\theta_t)} \left[\sum_{i=1}^n \sum_{j=1}^m \log(\sigma_{\phi}(Z_i, V_{ij})) \right]$$
$$= \sum_{i=1}^n \mathbb{E}_{Z_i \sim p(\cdot|V_i;\theta_t)} \left[\sum_{j=1}^m \log(\sigma_{\phi}(Z_i, V_{ij})) \right]$$

1

Moreover, since ρ does not enter in the update rules for η , ϕ , nor in the calculation of the posterior, we can ignore it in our EM-DPO algorithm.

Appendix E Min-Max Regret Objective Derivation

We can write, by linearity of expectation:

$$R_{z}(\pi) - R_{z}(\pi_{z}^{*})$$

$$= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{z}^{*}(\cdot|x), y' \sim \pi(\cdot|x)} \left[r_{z}^{*}(y, x) - r_{z}^{*}(y', x) \right]$$

$$= \beta \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{z}^{*}(\cdot|x)} \left[\log \left(\frac{\pi_{z}^{*}(y|x)}{\pi_{\text{SFT}}(y|x)} \right) \right]$$

$$- \beta \mathbb{E}_{x \sim \mathcal{D}, y' \sim \pi(\cdot|x)} \left[\log \left(\frac{\pi_{z}^{*}(y'|x)}{\pi_{\text{SFT}}(y'|x)} \right) \right]$$

For any $z, z' \in \{z_1, \ldots, z_k\}$, we will let:

$$\mathcal{L}_{z,z'} := \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{z'}^*(\cdot|x)} \left[\log \left(\frac{\pi_z^*(y|x)}{\pi_{\mathsf{SFT}}(y|x)} \right) \right]$$

Given the policy parameters we estimated in the EM-DPO section, these quantities can be calculated as simple empirical averages over the annotated data. Moreover, note that since our policy $\pi \in \Pi$ is a mixture policy over the policies $\pi_{z'}^*$ for $z' \in \{z_1, \ldots, z_k\}$ with weights $w \in \Delta(K)$, we can write:

$$R_{z}(\pi) - R_{z}(\pi_{z}^{*}) = \beta \left(\mathcal{L}_{z,z} - \sum_{k=1}^{K} w_{k} \cdot \mathcal{L}_{z,z_{k}} \right)$$
$$= \beta \sum_{k=1}^{K} w_{k} \cdot \left(\mathcal{L}_{z,z} - \mathcal{L}_{z,z_{k}} \right),$$

Thus, our minimax regret objective can be simply written as:

$$\min_{w \in \Delta(K)} \max_{z \in \{z_1, \dots, z_k\}} \left[\sum_{k=1}^K w_k \cdot (\mathcal{L}_{z,z} - \mathcal{L}_{z,z_k}) \right]^+$$
$$= \min_{w \in \Delta(K)} \max_{z \in \{z_1, \dots, z_k\}} \max \left\{ 0, \sum_{k=1}^K w_k \cdot (\mathcal{L}_{z,z} - \mathcal{L}_{z,z_k}) \right\}$$

Introducing a fake preference population z_0 that always has 0 regret, i.e. $\mathcal{L}_{z_0,z_0} = \mathcal{L}_{z_0,z_k} = 0$, we can re-write the above objective simply as:

$$\min_{w \in \Delta(K)} \max_{z \in \{z_0, z_1, \dots, z_k\}} \sum_{k=1}^K w_k \cdot (\mathcal{L}_{z, z} - \mathcal{L}_{z, z_k})$$

Appendix F Additional Experiment Details

F.1 IMDb Data Generation

We use the IMDb dataset [30] to generate a synthetic preference dataset. More specifically, we use a publicly available adaptation of the IMDb dataset¹. This dataset uses the first 20 tokens from the original IMDb dataset [30] as the prompt and then two responses are generated for each prompt using a GPT-2 Large model that is fine-tuned on the IMDb dataset. We synthetically generate preference data for two user groups using this dataset. We split the 50,000 available data points into 256 test examples and the remaining are training examples. To contruct a user, we first assign the population-type (i.e. that they prefer grammatically correct response or that the prefer the positive response) and then sample 50 data points randomly from the train set for each user. We use LanguageTool² to automatically to find the number of grammatical errors in a given text and divide this number by the length of the text to get a correctness score. The grammar-type user prefer the response with a higher correctness score. For the positive-sentiment preferring user, we choose the most positive response by prompting GPT-4.

F.2 More IMDb Results

Figure 6 and 7 show the accuracy metrics for both the clustering and the aggregation steps.

¹Modified IMDb dataset

²LanguageTool



Figure 6: Accuracy after the clustering step



Figure 7: Accuracy after the aggregation step

F.3 Cluster-DPO

The Cluster-DPO policy is generated as follows. We naively cluster the users into 2 user sub-groups using k-means clustering on the average embedding of all the preferred texts of that user. Embeddings are generated using the RoBERTa-Large model [29]. Then, we train a DPO policy on each cluster separately and combine them using Algorithm 2; we are essentially replacing the EM-DPO step with a k means clustering step in the proposed algorithm pipeline.

F.4 Hyperparameters

Table 1 shows the hyperparameters for the bandit experiment and Table 2 for the LLM experiment. We ran the bandit experiment on one A100 GPU. On average, the code took approximately 1 hour to run. The LLM experiment was run on 5x NVIDIA A6000 GPUs. Every step of the EM algorithm took about 40 minutes to run for a grand total of 13 hours.

Hyperparameter	Value
Neural Network Layers	3
Neural Network Hidden Dimension	10
Learning Rate	0.01
Optimizer	Adam
DPO Regularization Constant Beta	1
Max Epochs for Optimization	1000
Max Steps for EM-DPO	100
Max Steps for MINMAX-DPO	1000
Seed (numpy and torch)	123

Table 1: Hyper-parameters for the contextual bandit experiment

Parameter	Value
Learning Rate	5e-7
Beta	0.1
Max Text Length (Prompt + Response)	512
No. of Training Epochs	1
No. of Evaluation Examples	256
Optimizer	RMSprop
No. of Warmup Steps for Learning Rate	150
No. of Iterations of the EM Algorithm	10
No. of Prompts to Estimate the Regret Matrix	512
Eta for Algorithm 2	0.05
Total Steps for Algorithm 2	10000
No. of Examples for Evaluation	256
Seed (DPO), Seed1 (Evaluation), Seed2 (Evaluation)	0, 42, 62

Table 2: Hyper-parameters for the IMDb LLM experiment