

From Efficiency to Equity: Measuring Fairness in Preference Learning

Shreeyash Gowaikar

Department of Computer Science
Birla Institute of Technology and Science, Pilani
f20201719@goa.bits-pilani.ac.in

Hugo Berard

UNESCO Chair in Urban Landscape
University of Montréal
hugo.berard@umontreal.ca

Rashid Mushkani

UNESCO Chair in Urban Landscape
University of Montréal
rashid.ahmad.mushkani@umontreal.ca

Shin Koseki

UNESCO Chair in Urban Landscape
University of Montréal
shin.koseki@umontreal.ca

Abstract

As AI systems, particularly generative models, increasingly influence decision-making, ensuring that they are able to fairly represent diverse human preferences becomes crucial. This paper introduces a novel framework for evaluating epistemic fairness in preference learning models inspired by economic theories of inequality and Rawlsian justice. We propose metrics adapted from the Gini Coefficient, Atkinson Index, and Kuznets Ratio to quantify fairness in these models. We validate our approach using two datasets: a custom visual preference dataset (AI-EDI-Space) and the Jester Jokes dataset. Our analysis reveals variations in model performance across users, highlighting potential epistemic injustices. We explore pre-processing and in-processing techniques to mitigate these inequalities, demonstrating a complex relationship between model efficiency and fairness. This work contributes to AI ethics by providing a framework for evaluating and improving epistemic fairness in preference learning models, offering insights for developing more inclusive AI systems in contexts where diverse human preferences are crucial.

1 Introduction

The rapid advancement of generative artificial intelligence (AI) has brought unprecedented capabilities in natural language processing and content generation. However, these developments have also raised significant concerns about the potential for these systems to perpetuate or amplify epistemic injustice [18]. Epistemic injustice, a concept introduced by Miranda Fricker [13], refers to wrongs done to individuals in their capacity as knowers. In the context of generative AI, this manifests as the misrepresentation or misunderstanding of views from minorities and marginalized groups, potentially subjecting them to epistemic violence by denying their own subjective experience.

Generative AI systems, such as large language models, risk producing epistemic injustices by embedding biases through their training data and amplifying narratives from the Global North while silencing voices from the Global South. This imposition of a singular framing onto diverse global perspectives fails to recognize the multitude of views and experiences that exist worldwide. The challenge, therefore, is to develop generative AI systems that can fairly represent all existing views and perspectives, acknowledging and respecting the diversity of human experiences.

To address these alignment challenges, researchers have turned to techniques like Reinforcement Learning with Human Feedback (RLHF) [39]. RLHF typically involves a multi-step process: first,

gathering a dataset where human annotators indicate their preferences among a set of AI-generated options; second, training a reward model on this dataset to predict which options were preferred; and finally, using this reward model to fine-tune the generative model to align it with human preference. While this approach can improve alignment between AI outputs and human values, it raises new questions about epistemic justice.

The concept of the "tyranny of the majority," long recognized in political theory [23], becomes relevant in this context. If the reward model is biased towards certain groups, it may capture the preferences of dominant groups at the expense of marginalized voices. This scenario could lead to a digital manifestation of majority rule, where the opinions and values of the numerical majority consistently overshadow those of minority groups in AI-generated content.

To mitigate this risk, we argue that reward models should be trained on a diverse set of preferences from annotators who are as representative of the global population as possible. However, even with a diverse dataset, the underrepresentation of minorities may still lead to their perspectives being overlooked in the final model. Therefore, it is crucial to develop methods for measuring and addressing this form of epistemic injustice within reward models.

In this paper, we propose novel metrics, inspired by economic literature on inequality and fair allocation, to quantify the extent to which reward models equally capture the preferences of all users. We demonstrate the application of these metrics on two preference learning tasks, revealing that epistemic injustice can persist even in models with high overall accuracy. Our findings underscore the importance of looking beyond aggregate performance metrics to ensure equitable representation of diverse perspectives.

Furthermore, our research highlights a critical gap in the field: the scarcity of datasets containing individual annotations that would allow for more nuanced analysis of preference distribution across different groups. We advocate for the creation and public release of such datasets, which are essential for advancing research on epistemic justice in AI and developing more inclusive generative models. By addressing these challenges, we aim to contribute to the development of generative AI systems that not only perform well on standard metrics but also uphold principles of epistemic justice, ensuring that the diversity of human knowledge and experience is respected and accurately represented in AI-generated content.

2 Related Work

The intersection of fairness, epistemic justice, and AI has garnered significant attention in recent years, particularly in the context of classification and regression tasks. However, less attention has been paid to fairness and epistemic considerations in more complex AI systems such as generative models and preference learning algorithms.

Fairness in Classification and Regression In traditional machine learning tasks, fairness metrics often focus on equalizing outcomes across different groups. [16] introduced the concept of Equal Opportunity, which aims to equalize true positive rates between protected and unprotected groups. Other measures include equalized odds, ensuring equal probability of positive outcomes across classes, and equal accuracy, which balances performance across groups [9]. These metrics, while valuable, primarily address fairness in standard classification tasks, where the goal is to ensure that the algorithm's output does not depend on sensitive attributes. However, in the context of Reinforcement Learning from Human Feedback (RLHF) and preference learning, the concept of fairness requires a different approach. In these scenarios, we acknowledge that different groups may have varying preferences, and thus, the output of the reward model should rightfully depend on the user. Our notion of fairness in this context focuses on ensuring that the model's accuracy in capturing these diverse preferences does not vary significantly across different groups and individual users.

Fairness in Preference Learning and Ranking While works like [26] and [7] have proposed fairness metrics for ranking tasks, their focus primarily remains on ensuring fair treatment of the items being ranked. [31] take a step further by discussing the equity of subgroup populations and exploring the trade-off between this equity and overall accuracy. However, there remains a crucial distinction between these approaches and ours. Standard fairness approaches in ranking typically aim to ensure that items from different groups (e.g., protected vs. unprotected) have equal opportunities

to be ranked highly. This focus on the fairness of outcomes for ranked items is important but does not address the full spectrum of fairness concerns in preference learning scenarios.

Our approach, in contrast, shifts the focus to the equity of the participants providing the rankings or expressing preferences. We argue that in preference learning and RLHF contexts, it's crucial to ensure that the model's ability to capture and represent preferences is equitable across all participants, regardless of their group membership. This means that while the content of preferences may vary across groups (which is expected and acceptable), the accuracy with which these preferences are captured and represented by the model should be consistent across all participants.

This distinction is particularly important in scenarios where diverse viewpoints and experiences are critical, such as in collaborative decision-making or in developing AI systems that need to be responsive to a wide range of user preferences. By focusing on the equity of participants rather than just the ranked items, we aim to develop models that are truly representative of diverse human preferences and experiences.

This perspective reveals a critical gap in the literature: most current approaches treat participants as interchangeable, assuming a universal preference or aggregating scores without considering individual differences [32, 6]. This "One-Truth" fallacy [2] fails to account for the diversity of human preferences, which can stem from personal, environmental, and socio-demographic factors [30].

Towards Diverse Preference Representation Recent work has begun to acknowledge the importance of diverse human preferences in AI alignment and preference learning. Approaches include developing multiple reward models [8, 6], multi-policy strategies [28], and consensus-based ranking [20]. However, there remains a lack of consistency in how diversity and model performance are measured and evaluated. Most studies rely on basic classification metrics like F1 score and accuracy across users [30, 37, 38] or average reward values [28, 6]. While these metrics provide some insight, they fail to capture the nuanced ways in which AI systems might perpetuate epistemic injustice by misrepresenting or undermining the subjective experiences of marginalized groups.

3 Background

This section introduces a mathematical framework for preference learning that allows us to quantify both the overall performance of a model and its fairness across diverse users.

Problem Definition Consider a set of k users $\mathcal{U} = \{1, \dots, k\}$ and a dataset $\mathcal{D} = \{(x_i, x'_i, s_i, u_i)\}_{i=1}^n$ composed of n pairwise comparisons. Each entry in the dataset represents a comparison where user u_i provided a score $s_i \in \mathcal{S}$. The score can be either: a binary variable (i.e. $\mathcal{S} = \{0, 1\}$ indicating which option was preferred, or a real value (i.e. $\mathcal{S} = \mathbb{R}$), where negative scores indicate preference for x_i , and positive scores preference for x'_i , and the magnitude of s_i reflects the strength of the preference. Our goal is to learn a model $f : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{Y}$ that can score any pair $(x, x') \in \mathcal{X} \times \mathcal{X}$. We define the error of the model as:

$$\mathcal{E}(f) = \mathbb{E}[\ell(f(x_i, x'_i), s_i)]$$

where ℓ is a loss function that computes the discrepancy between the model outputs and the ground truth score. The choice of loss function depends on the nature of the scores:

1. For real-valued scores, we can use the squared error:

$$\ell(f(x_i, x'_i), s_i) = (f(x_i, x'_i) - s_i)^2$$

2. For binary scores, where f predicts the probability that x_i is preferred over x'_i , we can use the binary cross-entropy (BCE) loss:

$$\ell(f(x_i, x'_i), s_i) = s_i \log(f(x_i, x'_i)) + (1 - s_i) \log(1 - f(x_i, x'_i))$$

3. Alternatively, we can use the 0-1 loss:

$$\ell(f(x_i, x'_i), s_i) = \begin{cases} 0 & \text{if } s_i = f(x_i, x'_i) \\ 1 & \text{else} \end{cases}$$

To evaluate the model’s performance for individual users, we define the user-specific error for user u :

$$\mathcal{E}_u(f) = \mathbb{E}[\ell(f(x_i, x'_i), s_i)|u]$$

Drawing inspiration from game theory literature on fair allocation, we introduce two key concepts that will help us evaluate both the overall performance and the fairness of our preference learning models.

Definition 1 (Efficiency). This efficiency measures the overall performance of the model across all users. For a model f , it is the mean of the errors across all users:

$$\bar{\mathcal{E}}(f) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \mathcal{E}_u(f)$$

Definition 2 (Equality). The equality measure helps us quantify the fairness of the model by looking at the worst-case performance. The equality of a model f is the maximum error among all users:

$$\mathcal{E}_{\max}(f) = \max_u \mathcal{E}_u(f)$$

These definitions formalize the trade-off between overall performance and fairness in preference learning models. While traditional machine learning approaches focus on minimizing average error, this can lead to performance disparities across users, potentially perpetuating epistemic injustices. By introducing equality alongside efficiency, we propose a framework for developing models that maintain consistent accuracy across all users.

Our focus on equality, particularly through the $\mathcal{E}_{\max}(f)$ metric, resonates with John Rawls’ maximin principle of justice [29, 19]. This approach prioritizes the welfare of the worst-off members, which in our context translates to minimizing the maximum error across all users. This Rawlsian perspective provides a philosophical justification for prioritizing outcomes for disadvantaged users or groups, ensuring that AI systems accurately represent all preferences, including those from marginalized or underrepresented populations.

4 Equality Metrics

We now introduce a comprehensive set of metrics to quantify the extent to which a model’s performance varies across users. These metrics, adapted from the economics literature on income inequality, allow us to measure different aspects of fairness and equality in AI systems.

The importance of these metrics lies in their ability to capture various manifestations of inequality in model performance. For instance, errors might be highly dissimilar across specific groups of users or may vary more gradually across the user population. By employing a range of metrics, each with distinct characteristics, we can gain a nuanced understanding of how fair our preference learning models is. All of the following metrics are non-negative and equal to zero only when the model’s performance is identical for all users, representing perfect equality.

Maximal Error Gap The Maximal Error Gap measures the largest discrepancy in model performance between any two users. This metric is particularly useful for identifying extreme cases of inequality and aligns with Rawlsian principles of justice by highlighting the worst-case scenario.

$$G_{\max}(f) = \max_{u, u' \in \mathcal{U}} (\mathcal{E}_u(f) - \mathcal{E}_{u'}(f)) \quad (1)$$

Standard Deviation of the Error This metric provides a measure of the overall spread of errors across users. A large standard deviation indicates significant variability in model performance, suggesting unequal representation of user preferences.

$$\sigma^2(f) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} (\mathcal{E}_u(f) - \bar{\mathcal{E}}(f))^2 \quad (2)$$

Gini Coefficient [14] The Gini Coefficient, widely used in economics to measure income inequality, it provides a holistic view of error distribution across users. It can be visualized using the Lorenz

curve, where the coefficient represents the area between the line of perfect equality and the actual error distribution curve. The Gini Coefficient is bounded between 0 (perfect equality) and 1 (extreme inequality).

$$G(f) = \frac{\sum_{u,u'} (|\mathcal{E}_u(f) - \mathcal{E}_{u'}(f)|)}{2|\mathcal{U}|^2 \bar{\mathcal{E}}(f)} \quad (3)$$

Generalised Entropy Index [33] This index offers flexibility through its α parameter, allowing us to focus on different parts of the accuracy distribution across users. Lower α values are sensitive to the existence of users with low accuracy. While higher α values are more sensitive to the existence of users with high accuracy.

$$G_\alpha(f) = \begin{cases} \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{\mathcal{E}_u(f)}{\bar{\mathcal{E}}(f)} \ln \frac{\mathcal{E}_u(f)}{\bar{\mathcal{E}}(f)} & \text{if } \alpha = 1 \\ -\frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \ln \frac{\mathcal{E}_u(f)}{\bar{\mathcal{E}}(f)} & \text{if } \alpha = 0 \\ \frac{1}{|\mathcal{U}|^{\alpha(\alpha-1)}} \sum_{u \in \mathcal{U}} \left[\left(\frac{\mathcal{E}_u(f)}{\bar{\mathcal{E}}(f)} \right)^\alpha - 1 \right] & \text{else} \end{cases} \quad (4)$$

Atkinson Index [4] Similar to the Generalized Entropy Index, the Atkinson Index uses an ϵ parameter to focus on inequalities at different ends of the accuracy distribution across users. As ϵ increases, the index becomes more sensitive to errors at the lower end of the distribution.

$$A_\epsilon(f) = \begin{cases} 1 - \frac{1}{\bar{\mathcal{E}}(f)} (\prod_{u \in \mathcal{U}} \mathcal{E}_u(f))^{1/N} & \text{if } \epsilon = 1 \\ 1 - \frac{1}{\bar{\mathcal{E}}(f)} \min_{u \in \mathcal{U}} \mathcal{E}_u(f) & \text{if } \epsilon = +\infty \\ 1 - \frac{1}{\bar{\mathcal{E}}(f)} \left(\frac{\sum_{u \in \mathcal{U}} \mathcal{E}_u(f)^{1-\epsilon}}{|\mathcal{U}|} \right)^{\frac{1}{1-\epsilon}} & \text{else} \end{cases} \quad (5)$$

Kutznets Ratio [21] Unlike the previous metrics that consider the entire error distribution, the Kuznets Ratio focuses on the extremes, comparing the errors of the top $\alpha\%$ of users to the bottom $\alpha\%$. This metric is particularly useful for identifying disparities between the best and worst-served users:

$$K_\alpha(f) = \frac{\sum_{\text{top } \alpha\%} \mathcal{E}_u(f)}{\sum_{\text{bottom } \alpha\%} \mathcal{E}_u(f)} \quad (6)$$

By employing this diverse set of metrics, we can comprehensively evaluate the fairness and equality of preference learning models, allowing us to focus on different ends of the distribution.

5 Experimental Setup

To rigorously evaluate our proposed equality metrics, we carefully selected two datasets that provide crucial individual-level annotation data. This granular information—detailing which user provided each annotation—enables us to precisely analyze variations in model performance across diverse users. Such user-specific data is instrumental in uncovering potential epistemic injustices in AI systems, yet it is often absent from many widely-used datasets in the field of Reinforcement Learning from Human Feedback (RLHF).

Notably, prominent datasets such as Safe RLHF [10], Helpful and Harmless [5], WebGPT [25], ImageReward [36], and AVA [24] lack user identification data. This omission precludes the differentiation of users with diverse preferences, rendering the computation of our proposed equality metrics impossible on these datasets. The absence of such critical information in these widely-used resources highlights a significant gap in the field’s ability to assess and address epistemic injustice in generative models.

We posit that the inclusion and release of user-specific annotation data is not merely beneficial but essential for the comprehensive evaluation of fairness in AI systems. By enabling the application of metrics like those proposed in this study, such data would significantly enhance our capacity to identify, quantify, and ultimately mitigate epistemic injustices in generative models. This underscores the urgent need for more nuanced and comprehensive datasets in the pursuit of truly fair and equitable AI systems.

To ensure robustness of the results, all experiments employ the following methodology: 1) Model selection was based on the best performance on a randomized validation set, using comparisons

not present in the training data. 2) Each model was independently trained five times to account for stochasticity and we report the average results across the five runs and the corresponding confidence intervals.

5.1 AI-EDI-Space Dataset

Dataset The first dataset consists of 7,833 street-view images representing a diverse set of public spaces from the Montreal Metropolitan Area. The dataset includes 19,990 pairwise comparisons, evaluated by 22 individuals who were carefully selected to maximize diversity and include underrepresented groups based on ethnicity, gender, sexuality and age. Each participant evaluated a minimum of 500 comparisons based on 35 different criteria designed to capture various qualities of public spaces. Participants provided a real value between -1 and 1 to avoid Arrow’s impossibility theorem [3], as explained by [1]. However, this approach introduces complexity into the voting patterns, as illustrated in Figure 1.

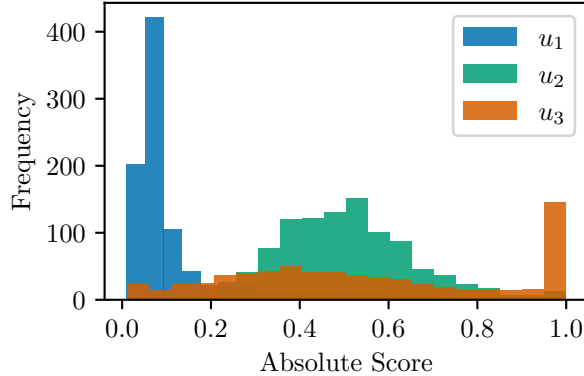


Figure 1: Voting patterns observed in the AI-EDI-Space Dataset. We plot the histogram of the absolute scores given by 3 participants with clearly distinct voting patterns. This figure highlights that users might have very different "taste profiles".

This dataset is ideal for testing algorithmic equity, as it includes a diverse set of participants with potentially divergent views on what makes a public space valuable, making the task highly subjective.

Model The model and training procedure used are similar to the one proposed in [11]. The model takes a single image as input and predicts the scores. The model consists of a feature extractor and a classifier head. The feature extractor is a pre-trained model. We experimented with several models, including VGG11 [34], EfficientNet [35], SqueezeNet [17], and DinoV2 [27]. The extracted features were then passed through a classifier head to predict a score for each of the 35 criteria. We observed that the model with EfficientNet, along with a double-layered classification head with 256 as the hidden dimension, gave the best results throughout all the experiments. Hence, all values for the AI-EDI-Space dataset are using an EfficientNet feature extractor. To train the model, we computed the scores for both images in each comparison and calculated the difference in scores between the two images. We then used the Mean Squared Error between the difference in scores and the ground truth score as the error to train the model.

Loss We use the 0-1 loss as the baseline metric to compute the various equality metrics, which corresponds to measuring the difference in accuracy across users.

5.2 Jester Jokes Dataset

Dataset We also test the proposed metrics on the Jester Jokes Dataset, originally developed for recommender systems research [15]. This dataset contains 100 jokes, each rated on a scale from -10 to +10 by 73,421 participants. The inherently subjective nature of humor appreciation makes this dataset particularly suitable for examining model performance across diverse user preferences. To ensure data quality and diversity, we only selected a subset of the annotations. We first filtered the

dataset to include only participants who rated all 100 jokes, ensuring comprehensive engagement from each user. From this filtered set, we then selected 1,000 users exhibiting the most diverse voting patterns. This selection was achieved through Principal Component Analysis (PCA) of users’ score vectors, followed by a uniform sampling of users maximally distant in the PCA space. This careful curation process resulted in a dataset that not only maintains balance during model training but also captures a wide spectrum of user preferences.

Model Since users directly provided scores for each joke, the problem can be formulated as a regression task. We trained a model that takes a single joke as input and predicts its score. The model consists of a feature extractor and a classifier head. The feature extractor is a pre-trained BART model [22]. The extracted features are then passed through a classifier head, which is either a single-layer or a double-layer perceptron, to predict the score. The Mean Squared Error (MSE) between the predicted score and the ground truth score was used as the loss function to train the model.

Loss We used the MSE loss as the baseline metric to compute the various equality metrics.

6 Methods

To address the challenge of inequality in model performance across users, we propose and evaluate two categories of techniques: pre-processing and in-processing. These approaches aim to enhance the fairness of preference learning models by ensuring more equitable representation of diverse user preferences.

6.1 Pre-Processing Techniques

Pre-processing techniques are applied to the data prior to model training. We investigate three scaling methods designed to normalize the distribution of scores across users:

1. **Min-Max Scaling:** This technique scales each participant’s scores to a range of $[-1, 1]$. It is applied individually to each user’s scores, preserving relative preferences within a user’s data while enabling comparability across users.
2. **Normalization Scaling:** This two-step process first applies standard normalization to each participant’s scores, adjusting the mean to 0 and standard deviation to 1. Subsequently, the scores are scaled to ensure they remain within the $[-1, 1]$ range. While this method, like Min-Max Scaling, is user-specific, it does not guarantee sparse unanimity.
3. **Mehestan Scaling [1]:** This more sophisticated approach considers the voting patterns of all participants when scaling an individual’s scores. The process involves: a) Converting raw comparison scores to individual scores using a Generalized Bradley-Terry Model [12]. b) Scaling and translating these scores using the BrMean primitive, which is designed to be resilient to potential manipulation by malicious voters. c) Preserving individual score distributions without final aggregation, maintaining the uniqueness of each participant’s preference pattern.

Mehestan Scaling is particularly effective in achieving sparse unanimity, a property that ensures the preservation of unanimous preferences even when user voting patterns differ significantly.

6.2 In-Processing Techniques

In-processing techniques are integrated into the model training process. We explore two primary approaches:

1. **User Embeddings:** By incorporating user-specific embeddings as additional input to the model, we aim to capture and adapt to individual voting patterns. This technique allows the model to learn user-specific features that may influence preference judgments.
2. **Contrastive Loss:** We employ contrastive loss in conjunction with least squares error (LSE). The contrastive loss works by increasing the distance between dissimilar scores, ensuring that the model’s outputs are not clustered too closely together. This helps prevent comparisons—calculated as the difference between the scores of two alternatives—from

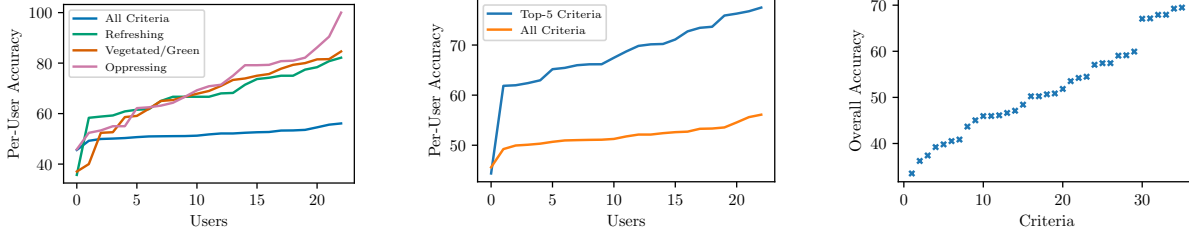
Experiment	Accuracy	G_{max}	σ^2	$G \times 10^{-2}$	$G_{\alpha=0} \times 10^{-4}$	$A_{\epsilon=\infty} \times 10^{-1}$	$K_{\alpha=20}$
Normalisation	51.1 ± 1.9	6.6 ± 1.1	1.7 ± 0.2	0.9 ± 0.1	5.5 ± 1.5	0.7 ± 0.1	1.37 ± 0.02
MinMax	50.2 ± 2.6	7.6 ± 4.0	1.8 ± 0.8	1.0 ± 0.4	7.1 ± 5.8	0.7 ± 0.6	1.38 ± 0.03
Mehestan	51.3 ± 4.4	7.1 ± 1.2	1.8 ± 0.3	1.1 ± 0.2	7.4 ± 3.4	0.7 ± 0.2	1.39 ± 0.04
Contrastive Loss	52.0 ± 0.6	7.7 ± 1.8	2.1 ± 0.6	1.1 ± 0.3	8.0 ± 4.7	0.8 ± 0.1	1.40 ± 0.04
User Emb.	50.2 ± 0.7	7.8 ± 1.0	1.9 ± 0.4	1.0 ± 0.2	7.0 ± 2.6	0.8 ± 0.4	1.39 ± 0.03

Table 1: Results for the AI-EDI Image Dataset over different experiments. Individual Users evaluated using Per-User Accuracy

being too close to zero, thereby improving the model’s ability to distinguish between different preferences.

7 Results

AI-EDI-Space Dataset After training a model on the AI-EDI-Space dataset, we computed the various equality metrics proposed earlier. Figure 2 shows the distribution of accuracy for all users across different criteria. Our analysis yielded several noteworthy observations: inequality appears to be higher for criteria where the model performs best overall, as shown in Table 1. This illustrates the potential trade-off between efficiency and equality. The observed inequality seems to be primarily driven by voting patterns, as the mean squared error (MSE) used to train the model penalizes discrepancies with the user’s comparisons. We also observe that users whose voting patterns cluster around 0 (a conservative voting approach) tend to achieve higher accuracy. Additionally, the number of comparisons annotated by each user may differ, potentially contributing to the observed inequalities. However, disentangling these effects requires further analysis to understand the influence of sample size on our observations.



(a) Sorted Per User Accuracy over i) all Criteria, ii) Top 3 Criteria - Oppressing, Vegetated/Green, Refreshing (in the descending order of accuracy).

(b) Sorted Per User Accuracy over i) All Criteria, ii) only 5 Criteria with highest overall accuracies (viz. Intimate, Regenerative, Refreshing, Vegetated/Green, Oppressing).

(c) Scatter Plot of accuracies for every criterion. The highest accuracy is for the criterion 'Oppressing', and the lowest accuracy is for the criterion 'Inviting/Welcoming'.

Figure 2: Experiment results on the AI-EDI Space Dataset

Jester Jokes Dataset Our analysis of the Jester Jokes dataset provided additional insights into the effectiveness of various scaling and translation techniques in addressing inequality. Table 2 presents the results of our equality metrics for different approaches. We made the following observations: 1) Normalization and MinMax Scaling achieved low Mean Squared Error (MSE) but failed to significantly improve equality. Notably, the Kuznets ratio remained close to 4, indicating substantial inequality between the top 20% and bottom 20% of participants in terms of model performance. 2) Mehestan Scaling, designed for sparse unanimity [1], yielded higher equality values despite a worse MSE. This finding suggests that techniques prioritizing unanimous preference recovery may contribute to greater epistemic fairness. 3) Given our careful selection of users who scored all 100 jokes, we can primarily attribute the observed inequality to differences in voting patterns rather than data imbalance. This reinforces the importance of considering diverse preference expressions in

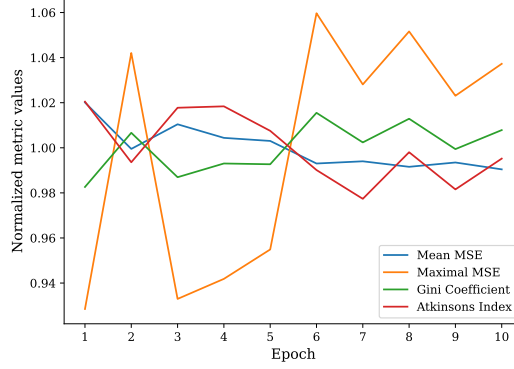


Figure 3: Performance and Metrics for Jester Jokes Dataset over epochs. A clear tradeoff in performance and equality can be seen from the figure where equality seems to increase over training epochs and performance decreases.

Experiment	MSE $\times 10^{-1}$	G_{max} $\times 10^{-1}$	σ^2 $\times 10^{-1}$	G $\times 10^{-1}$	$G_{\alpha=0}$ $\times 10^{-1}$	$A_{\epsilon=\infty}$ $\times 10^{-1}$	$K_{\alpha=20}$
Normalisation	2.66 ± 0.03	7.6 ± 0.3	1.38 ± 0.03	1.43 ± 0.02	1.55 ± 0.04	9.6 ± 0.1	4.81 ± 0.08
MinMax	2.61 ± 0.01	7.8 ± 0.5	1.35 ± 0.01	1.43 ± 0.01	1.55 ± 0.02	9.6 ± 0.3	4.84 ± 0.04
Mehestan	8.86 ± 2.12	21.9 ± 5.0	3.37 ± 0.98	1.02 ± 0.01	0.67 ± 0.15	6.7 ± 0.4	2.81 ± 0.37

Table 2: Results for the Jester Jokes Dataset over different experiments. Individual Users evaluated using Per-User Mean Squared Error

model development. Additionally, we can clearly see the tradeoff between performance and equality from Figure 3 where the average MSE falls over training epochs, indicating performance gains, whereas the Gini coefficient and Maximal MSE increase, indicating increased inequality.

8 Conclusion

In the era of generative AI, where algorithms increasingly shape decision-making processes, ensuring that these systems do not generate or amplify epistemic injustice is paramount. This study introduces a novel perspective on fairness in preference learning, focusing on the equitable representation of diverse human preferences and views. Our work provides a framework for quantifying and addressing epistemic fairness in AI models, contributing to the development of more just and inclusive AI technologies. Our findings underscore the potential for significant disparities in how well AI models capture preferences across different users. This raises critical questions about epistemic justice in AI systems and highlights the need for further research in several key areas: 1) Further investigation is needed to understand the sources of inequality in model performance and develop effective mitigation strategies. This includes examining how data characteristics, model architectures, and diverse human preferences interact to produce or exacerbate inequalities. 2) As Reinforcement Learning from Human Feedback (RLHF) becomes more prevalent, ensuring that the alignment process itself is equitable across diverse participant groups is crucial. This involves developing methods to capture a wide range of opinions and preferences, particularly from marginalized or underrepresented groups. 3) We advocate for the public release of more datasets that include annotation-level user information, specifically detailing which annotator or user provided each individual annotation. This granular data is crucial for conducting comprehensive evaluations of epistemic justice in AI models. Such datasets would enable researchers to track how different users' preferences and judgments are represented in model outputs, providing a more nuanced understanding of potential biases or inequalities in preference learning and generative AI systems. 4) Finally, the tension between individual fairness and overall system efficiency raises important ethical questions. Future work should explore how to balance these concerns in line with principles of distributive justice and epistemic fairness, particularly in the context of generative AI systems.

This study represents a step towards more equitable AI systems that respect and accurately represent diverse human preferences. As AI continues to play an increasingly significant role in society, addressing epistemic fairness will be crucial in ensuring that these systems serve all members of society equitably. Our work provides a foundation for future research in this critical area, aiming to develop AI technologies that are not only efficient but also just and inclusive.

References

- [1] Youssef Allouah, Rachid Guerraoui, Lê-Nguyen Hoang, and Oscar Villeda. Robust sparse voting. In *International Conference on Artificial Intelligence and Statistics*, pages 991–999. PMLR, 2024.
- [2] Lora Aroyo and Chris Welty. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine*, 36(1):15–24, 2015.
- [3] Kenneth J. Arrow. A Difficulty in the Concept of Social Welfare. *Journal of Political Economy*, 58(4):328–346, August 1950.
- [4] Anthony B Atkinson et al. On the measurement of inequality. *Journal of economic theory*, 2(3):244–263, 1970.
- [5] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.
- [6] Michiel A. Bakker, Martin J Chadwick, Hannah Sheahan, Michael Henry Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matthew Botvinick, and Christopher Summerfield. Fine-tuning language models to find agreement among humans with diverse preferences. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [7] Alex Beutel, Jilin Chen, Tulsee Doshi, Hai Qian, Li Wei, Yi Wu, Lukasz Heldt, Zhe Zhao, Lichan Hong, Ed H. Chi, and Cristos Goodrow. Fairness in Recommendation Ranking through Pairwise Comparisons. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2212–2220, Anchorage AK USA, July 2019. ACM.
- [8] Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Amrit Bedi, and Mengdi Wang. Maxmin-RLHF: Towards equitable alignment of large language models with diverse human preferences. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, 2024.
- [9] Andrew Cotter, Heinrich Jiang, Maya Gupta, Serena Wang, Taman Narayan, Seungil You, and Karthik Sridharan. Optimization with Non-Differentiable Constraints with Applications to Fairness, Recall, Churn, and Other Goals. *Journal of Machine Learning Research*, 20:1–59, 2019.
- [10] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe RLHF: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*, 2024.
- [11] Abhimanyu Dubey, Nikhil Naik, Devi Parikh, Ramesh Raskar, and César A. Hidalgo. Deep Learning the City: Quantifying Urban Perception at a Global Scale. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science, pages 196–212, Cham, 2016. Springer International Publishing.
- [12] Julien Fageot, Sadegh Farhadkhani, Lê-Nguyen Hoang, and Oscar Villeda. Generalized bradley-terry models for score estimation from paired comparisons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 20379–20386, 2024.

- [13] Miranda Fricker. *Epistemic injustice: Power and the ethics of knowing*. Oxford University Press, 2007.
- [14] Corrado Gini. On the measure of concentration with special reference to income and statistics, colorado college publication. *General series*, 208(1), 1936.
- [15] Ken Goldberg, Theresa Roeder, Dhruv Gupta, and Chris Perkins. Eigentaste: A Constant Time Collaborative Filtering Algorithm. *Information Retrieval*, 4(2):133–151, 2001.
- [16] Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. Equality of Opportunity in Supervised Learning. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [17] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [18] Jackie Kay, Atoosa Kasirzadeh, and Shakir Mohamed. Epistemic injustice in generative ai. *arXiv preprint arXiv:2408.11441*, 2024.
- [19] Patrik Joslin Kenfack, Samira Ebrahimi Kahou, and Ulrich Aïvodji. A survey on fairness without demographics. *Transactions on Machine Learning Research*, 2024.
- [20] Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cédric Colas, Peter Ford Dominey, and Pierre-Yves Oudeyer. Large language models as superpositions of cultural perspectives, 2023.
- [21] Simon Kuznets. Economic growth and income inequality. In *The gap between rich and poor*, pages 25–37. Routledge, 2019.
- [22] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online, July 2020. Association for Computational Linguistics.
- [23] John Stuart Mill. *On liberty*. 1859.
- [24] Naila Murray, Luca Marchesotti, and Florent Perronnin. Ava: A large-scale database for aesthetic visual analysis. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2408–2415, 2012.
- [25] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. Webgpt: Browser-assisted question-answering with human feedback, 2022.
- [26] Harikrishna Narasimhan, Andrew Cotter, Maya Gupta, and Serena Wang. Pairwise fairness for ranking and regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5248–5255, 2020.
- [27] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [28] Alexandre Ramé, Guillaume Couairon, Mustafa Shukor, Corentin Dancette, Jean-Baptiste Gaya, Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards, 2023.
- [29] John Rawls. *A theory of justice*. 1971.
- [30] Marta Sandri, Elisa Leonardelli, Sara Tonelli, and Elisabetta Jezek. Why don’t you do it right? analysing annotators’ disagreement in subjective tasks. In Andreas Vlachos and Isabelle Augenstein, editors, *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2428–2441, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.

- [31] Amartya Sanyal, Yaxi Hu, and Fanny Yang. How unfair is private learning ? In *The 38th Conference on Uncertainty in Artificial Intelligence*, May 2022.
- [32] Richard L. Sawyer, Nancy S. Cole, and James W. L. Cole. Utilities and the Issue of Fairness in a Decision Theoretic Model for Selection. *Journal of Educational Measurement*, 13(1):59–76, 1976.
- [33] Anthony F Shorrocks. The class of additively decomposable inequality measures. *Econometrica: Journal of the Econometric Society*, pages 613–625, 1980.
- [34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [35] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [36] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023.
- [37] Tianshu Yu, Ting-En Lin, Yuchuan Wu, Min Yang, Fei Huang, and Yongbin Li. Constructive large language models alignment with diverse feedback, 2023.
- [38] Dun Zeng, Yong Dai, Pengyu Cheng, Longyue Wang, Tianhao Hu, Wanshun Chen, Nan Du, and Zenglin Xu. On diversified preferences of large language model alignment, 2024.
- [39] 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.