

A Study of Implicit Ranking Unfairness in Large Language Models

Anonymous ACL submission

Abstract

Content Warning: This paper contains examples of misgendering and erasure that could be offensive and potentially triggering.

Recently, Large Language Models (LLMs) have demonstrated a superior ability to serve as ranking models. However, concerns have arisen as LLMs will exhibit discriminatory ranking behaviors based on users' sensitive attributes (e.g., gender). Worse still, in this paper, we identify a subtler form of discrimination in LLMs, termed *implicit ranking unfairness*, where LLMs exhibit discriminatory ranking patterns based solely on non-sensitive user profiles, such as user names. Such implicit unfairness is more widespread but less noticeable, threatening the ethical foundation. To comprehensively explore such unfairness, our analysis will focus on three research aspects: (1) We propose an evaluation method to investigate the severity of implicit ranking unfairness. (2) We uncover the reasons for causing such unfairness. (3) To mitigate such unfairness effectively, we utilize a pair-wise regression method to conduct fair-aware data augmentation for LLM fine-tuning. The experiment demonstrates that our method outperforms the existing methods regarding ranking fairness. Lastly, we emphasize the need for the community to identify and mitigate the implicit unfairness, aiming to avert the potential deterioration in the reinforced human-LLMs ecosystem deterioration.

1 Introduction

Large language models (LLMs), represented by ChatGPT (Wu et al., 2023b) have empowered ranking tasks (Wu et al., 2023a), which plays an important role in filtering overload information to users (Liu et al., 2009). However, ensuring that LLMs do not pose ethical risks becomes crucial. Recently, various evaluation methods have been introduced to assess the degree of discrimination in LLMs (Zhang et al., 2023b; Kasneci et al., 2023;

Chang et al., 2023), showing that LLMs frequently exhibit pronounced ranking discriminatory behaviors against explicit sensitive attributes, such as gender (Zhang et al., 2023b; Tamkin et al., 2023).

Although a massive amount of work focuses on addressing unfairness when explicitly using sensitive attributes in ranking tasks (Dai et al., 2024), our investigation reveals the persistence of *implicit ranking unfairness*: LLMs even generate substantial discriminatory ranking behaviors when using non-sensitive yet personalized user profiles (e.g., user names and email addresses). *Implicit ranking unfairness* in LLMs highlights new and more urgent risks towards LLMs-based application (e.g., recommendation, search) because (1) such implicit unfairness is often inconspicuous because it only depends on non-sensitive user profiles; and (2) implicit ranking unfairness is more widespread since these non-sensitive user profiles can be easily acquired and used by existing platforms, such as user names or email addresses. To comprehensively analyze the problem, in this paper, we will focus on three research aspects regarding implicit ranking unfairness in LLMs.

Firstly, we propose an evaluation method to investigate how serious the implicit ranking unfairness is in existing LLMs. Specifically, following the practice in (Zhang et al., 2023b), we design a ranking task prompt template (Figure 1). Then we give substantial empirical evidence to confirm the existence of implicit ranking unfairness. Finally, we find that the degree of implicit ranking unfairness is nearly 2-4 times more serious than explicit unfairness, and the unfairness is caused by collaborative information. Empirical evidence is in Section 4).

Secondly, since this implicit unfairness is more severe and more hidden, we aim to investigate the reasons behind its occurrence. Specifically, we identify that the LLMs can probe sensitive attributes exclusively from these personalized and

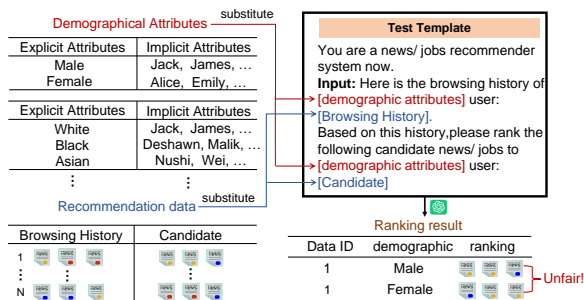


Figure 1: Overall workflow of our evaluation. The ranking list outputs by LLMs should be the same when replacing different sensitive attributes in prompts.

non-sensitive user profiles. Then we also show that the word embeddings of certain non-sensitive user profiles are more closely aligned with the sensitive attribute. Such phenomena contribute to the collection of unfair datasets during the pre-training phases (see evidence in Section 5).

Finally, we aim to propose a method to mitigate such implicit ranking unfairness. Previous research proposed to mitigate user unfairness either by employing privacy policies that hide sensitive attributes (Xiao et al., 2023; Brown et al., 2022; Kandpal et al., 2022), utilizing certain prompts to instruct LLMs to disregard sensitive attributes (Hua et al., 2023) or add counterfactual sample to enhance fairness (Ghanbarzadeh et al., 2023). However, they show limited effectiveness in mitigating implicit ranking unfairness (See Section 6).

In this paper, we propose a fair-aware data argumentation method to mitigate such unfairness. Specifically, we incorporate counterfactual samples that contain certain implicit attributes to help the model produce fair ranking results. Due to the massive and noisy characteristic of the non-sensitive features, we employ a pair-wise regression method to choose hard and informational non-sensitive features to conduct data argumentation. The experiments demonstrate that our method outperforms the existing methods on two ranking datasets.

Major Contributions: In summary, we have the following major findings: (1) We uncover that the LLMs-based ranking system demonstrates substantial implicit unfairness. (2) We analyze the reasons for causing such implicit unfairness. (3) We propose a new fair-aware data argumentation method to mitigate the implicit ranking unfairness effectively. Our code is available at https://anonymous.4open.science/r/Implicit_Rank_Unfairness-3C71.

2 Preliminary

Previous work shows LLMs’ powerful ability to serve as an information retriever (Dai et al., 2023; Bao et al., 2023b). Figure 1 shows the overall workflow of our settings.

In LLMs-based ranking applications, let \mathcal{U} be the user set. A user $u \in \mathcal{U}$ will have non-sensitive features v_u (e.g., user names) and sensitive features $s_u \in \mathcal{S}$ (e.g., user gender). In our work, we define the set \mathcal{S} to represent sensitive attribute types such as gender, race, or continent, and s_u is selected from options [Male, Female], [White, Black, Asian], or [Asian, Africa, Americas, Europe, Oceania]. When a user u engages ranking systems, a personalized prompt p_u will be used to instruct LLMs to make ranking results. Given the prompt p_u and optional user features v_u and s_u , the LLMs-based ranking model will output a ranking list $L_K(u) = \{i_1, i_2, \dots, i_K\}$, where K is the fixed ranking size and i_j is the j -th given item.

We consider the measurement as counterfactual fairness in individual-level (Wu et al., 2019; Li et al., 2023), i.e., the ranking list $L_K(u)$ outputs by LLMs should be the same in the counterfactual world as in the real world. For example, if we modify a user’s sensitive attribute from “male” (real world) to “female” (counterfactual world) while keeping all other characteristics constant (e.g., browsing histories), the ranking list should remain unchanged. Formally, given the same personalized prompt p_u and features v_u, s_u of the user, the general ranking model $f : L_K(u) = f(p_u, v_u, s_u)$ is counterfactually fair if for any $s', s \in \mathcal{S}$:

$$P(L_K(u)|s_u = s) = P(L_K(u)|s_u = s'), \quad (1)$$

where $P(L_K(u))$ is the distribution of $L_K(u)$.

Previous works (Zhang et al., 2023b) have found that when we explicitly take the sensitive feature s_u as input user features, recommender model f often does not meet the criteria outlined in the Equation (1). Formally, we can define:

Explicit ranking unfairness: $L_K(u) = f(p_u, v_u, s_u)$, which do not satisfy Equation (1).

However, we discover that even if we mask s_u as an input in the LLMs-based ranking model f , it still yields significantly discriminatory output distributions when categorized based on different sensitive attributes s_u . Formally, we can define:

Implicit ranking unfairness: $L_K(u) = f(p_u, v_u)$, and $L_K(u)$ do not satisfy the Equation (1). Because non-sensitive attribute v_u may

Table 1: Statics of different user names, where $|\mathcal{N}_s|$ denotes the number of user names belonging to the demographic group s .

s	Gender		Race		
	Male	Female	White	Black	Asian
$ \mathcal{N}_s $	1068	1040	1175	256	463
s	Continent				
	Asia	Americas	Africa	Europe	Oceania
$ \mathcal{N}_s $	463	374	136	1075	60

have a strong correlation with sensitive attribute s_u learned in the pre-training phase of LLMs.

3 Evaluation Settings

In this section, we will describe our evaluation settings including the datasets and some details.

3.1 Non-sensitive Attribute Selection

Specifically, we collect first names by choosing the most popular first names in 2014 from 229 countries (regions) across different genders, races, nationalities groups¹. The detailed statistic information is in Table 1.

3.2 Discrimination Measurement

Following (Gallegos et al., 2023), we utilize the metric U-Metric to measure the discrimination degree under the previous evaluation settings:

$$U(\mathcal{S}) = \sum_{s \in \mathcal{S}} |\text{Metric}(s) - \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \text{Metric}(s)| / |\mathcal{S}|,$$

where $\text{Metric}(s)$ is the evaluation metric under s group, which can be either $\text{NDCG}@K = \frac{1}{N} \sum_{j=1}^N \frac{\sum_{k=1}^K (2^{r_k} - 1) / (\log_2(j+1))}{(2^{\text{rank}_j} - 1) / (\log_2(\text{rank}_j + 1))}$, or other ranking metric such as MRR (Dai et al., 2023), where rank_j is the rank of the first correct answer in the ranking list $L_K(s, j)$ for user u within the top K recommendations, and r_k is a relevance score of the item with the k -th rank, which is 1 if it is a positive sample otherwise 0.

3.3 Other Settings

In this section, we will describe our evaluation settings including the datasets and some details.

Dataset. We utilize the two common-used ranking datasets: **MIND** (Wu et al., 2020) collected user news click behaviors on the Microsoft platform, which comprises 15,777,377 impression logs from a total of 1 million users; **CareerBuilder** is

collected based on their previous online job applications, and work history. The data covers the records of 321,235 users applying for 365,668 jobs from April 1 to June 26, 2012.

Following the practice in (Dai et al., 2023; Zhang et al., 2023c), we also apply the filter criteria where both the impression list and history list are required to have more than 5 items each and sample 300 data uniformly to evaluate the LLMs in every trial.

LLM Settings. In all the experiments, we utilize the ChatGPT series (gpt-3.5-turbo-xxx)² and Llama2 (Touvron et al., 2023). The numbers "xxx" refer to the release or revision dates. In all LLMs, we set the maximum generated token number to 2048, the nucleus sampling ratio is 1, the temperature is 0.2, the penalty for frequency is 0.0, and the penalty for presence is 0.0.

4 Implicit Unfairness of LLMs

In this section, we aim to evaluate the implicit unfairness. Note that we average the different ChatGPT versions and Llama 2 (Touvron et al., 2023) results to conduct the analysis.

4.1 Existence of Implicit Unfairness

Specifically, we design N topic sentences, where several keywords of certain topics are formed into a topic sentence. Suppose T_1, T_2, \dots, T_N denotes the constructed topic sentence, where N denotes the topic number. The topic distribution $P(L_K(s))$ of group s is defined as $[S_1, S_2, \dots, S_N] = \text{Softmax}([Z_1, Z_2, \dots, Z_N])$, where $Z_j = \sum_{n \in \mathcal{N}_s} \sum_{i \in L_K(n)} e^{(T_j)^T e(i)}$.

Gender Discrimination. From the sub-figures in Figures 2(a) and 2(c), we can observe that LLMs tend to provide noticeably different responses for different genders. For example, in news recommendations, ChatGPT will deliver more political news to male users while giving more life, health, art, and sports related news to female users. In the context of job recommendations, ChatGPT tends to suggest a higher number of service-related positions to male users and an increased number of medical-related jobs to female users.

Race Discrimination. From the sub-figures in Figure 2(b) and 2(d), we find that LLMs also give different category ratios for different races. For example, LLMs will deliver more political but less art news to black users. As for job recommendations, LLMs tend to recommend more service-related but

¹<https://forebears.io/forenames/most-popular>

²<https://platform.openai.com/>

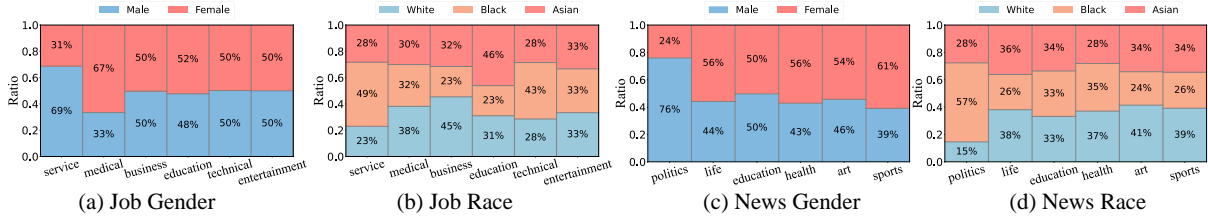


Figure 2: The discriminatory behaviors against certain topics of LLMs under job and news domain for user names belonging to different Gender and Race groups.

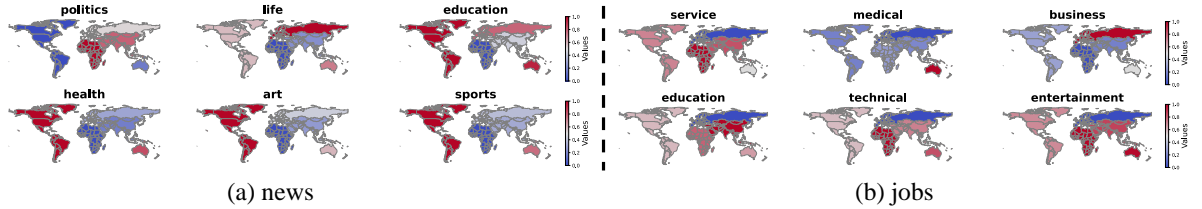


Figure 3: The discriminatory ranking behaviors against certain topics of LLMs under job and news domain for user names belonging to different Continent groups. A deeper red color indicates that LLMs are more likely to assign this type of news or jobs to users in the continent, while a deeper blue color suggests that LLMs are less likely to assign this type of news or jobs to users in the continent.

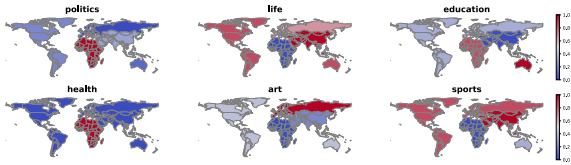


Figure 4: The discriminatory ranking behaviors against certain topics of LLMs under the news domain for user emails. A deeper red/blue color indicates that LLMs are more/less likely to assign this type of news.

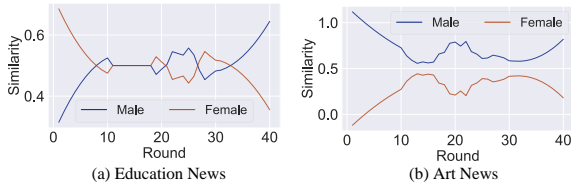


Figure 5: Similarity curves of different gender groups *w.r.t.* interaction rounds. Higher similarity denotes the LLMs will deliver more items related to topics to users.

less educational jobs to black users. Meanwhile, LLMs are likely to give more business and educational jobs to white and Asian users, respectively.

Continent Discrimination. From Figure 3 we can observe that LLMs reveal stereotype bias at the geographical level. Similarly, LLMs will deliver more political news to African users while more education, health, art, and sports-related news to users in America. In the realm of job recommendations, there is a tendency for LLMs to suggest a greater number of service-oriented positions to African users, whereas it leans toward proposing more educational jobs to Asian users.

Influences for Other Attributes. We also examine whether LLMs can exhibit implicit ranking unfairness when email addresses are used as non-sensitive features. Specifically, we choose the continental top 10 university email domain address ³.

From Figure 4, we can observe a similar discriminatory ranking pattern compared to the implicit ranking fairness when utilizing user names (see Figure 3). For example, LLMs will deliver more political and healthy news to users whose email domain addresses are African universities and more life and sports news to users whose email domain addresses are America’s universities. The experiments also verified different non-sensitive features can all cause serious implicit user unfairness.

Implicit Unfairness During Conversation. Next, to investigate the implicit unfairness degree during the conversation process, following the practice in (Zhang et al., 2023a), we will give a simulation interactive process between the user and ranking models every round. For each round, the LLMs will give a ranking list L_K with size K according to a user’s browsing history. Next, the user will select an item whose is in the first position of L_K , to serve as their browsing history for the next interaction round, since previous research has indicated that users tend to view items in higher positions (Craswell et al., 2008).

From Figure 5 (a) and (b), we can observe that in the long term, LLMs exhibit a higher tendency to

³<https://www.usnews.com/education/best-global-universities/>

recommend unipolar news. For example, it tends to recommend more art and education news to male users than female users gradually, causing information bubbles for male and female groups.

The experiment confirmed that implicit ranking unfairness in LLMs-based ranking models may lead to more reinforced unipolar ranking results, which pose a threat to diversity and potentially trap different user groups within information bubbles.

4.2 Implicit Ranking Unfairness Degree

In this section, our objective is to investigate how is implicit ranking unfairness compared with explicit unfairness and unfairness caused by the collaborative filtering information.

Comparison with Explicit Unfairness. In Figure 7, we compare the discrimination degrees (U-NDCG@3 and U-NDCG@5) under three demographic types with the explicit and implicit ranking unfairness utilizing different versions of ChatGPT and Llama2.

From Figure 7, we discern that in the evaluation at the Continent level, both the explicit and implicit ranking unfairness exhibit similar averaged discrimination measurements. However, when comparing the Gender and Race levels, we find that explicit unfairness is often lower than the implicit fairness degree by about 2-4 times. These experiments also confirm that when utilizing common demographic terms such as “Male” and “White”, LLMs are more likely to cause implicit fairness.

Influence of Collaborative Filtering. Previous research indicates that collaborative filtering information utilized in ranking during pre-training may also contribute to unfairness (Yao and Huang, 2017). Therefore, we aim to conduct a simulation to investigate the unfairness degree raised by collaborative filtering (CF) information. We choose DCN (Wang et al., 2017) and GRU4Rec (Tan et al., 2016) as two commonly used ranking models for learning CF information.

Specifically, owing to the privacy policy, the dataset does not include any sensitive attributes of users. Therefore, for every user, we utilized the point-wise probing described in Section 5 to predict the sensitive attributes of a user. Specifically, for at time t , we utilized the historical clicked item sequence $[i^{t-H}, i^{t-H+1}, \dots, i^{t-1}]$ to simulate, i.e. $\hat{s}_u = \arg \max_{s \in S} \sum_{h=1}^H \left(\hat{z}_s^{\text{point}}(i^{t-h}) / \tilde{z} \right)$, where H is the pre-defined maximum history length. Given the simulated sensitive attribute as the user

Table 2: Testing accuracy for probing using ChatGPT and Llama2 on news and job recommendation tasks.

	demographic	gender	race	continent
news	ChatGPT	0.667	0.659	0.510
	Llama2	0.833	0.777	0.466
jobs	ChatGPT	0.552	0.645	0.505
	Llama2	0.916	0.666	0.533
	random	0.500	0.333	0.200

context, trained a ranking model based on this context. In the inference phase, we mixed the data both in real-world and counterfactual world (Wu et al., 2019; Kusner et al., 2017), i.e. keeping other features constant, we replaced the user-sensitive attributes to assess the performance variation among different groups, considering this difference as a measure of unfairness degree.

From the reported results in Table 3, we can see the degree of implicit ranking unfairness in LLMs significantly outperforms all of the unfairness learned with CF information. The experiment verifies that implicit ranking unfairness does not rely on much on collaborative information but contributes to the correlation between non-sensitive attributes and sensitive attributes.

5 Implicit Ranking Unfairness Traceback

In this section, our objective is to investigate why the implicit ranking unfairness exist.

5.1 Inferring Sensitive Attribute Ability

Firstly, we utilize the probing technique (Vulić et al., 2020; Gurnee and Tegmark, 2023) under two most-performing LLMs ChatGPT (Roumeliotis and Tselikas, 2023) and Llama-2 7B (Touvron et al., 2023) to investigate whether LLMs can inference the sensitive attribute from the non-sensitive attribute in terms of their wide world knowledge.

To validate the effectiveness of pair-wise regression, we also compare the commonly used point-wise probing (Gurnee and Tegmark, 2023) to predict the appropriate demographic attribute utilizing non-sensitive attributes:

$$l^{\text{point}} = \mathbb{E}_j \left[\sum_{s \in S} \sum_{n \in \mathcal{N}_s} \sum_{i \in L_K(n, j)} \mathbf{CE}(z, \hat{z}^{\text{point}}(i)) \right], \quad (2)$$

where $\hat{z}^{\text{point}}(i) = \mathbf{MLP}(e(i); \theta^{\text{point}})$ and $\mathbf{CE}(\cdot)$ denotes the cross entropy loss.

From Table 2, we can observe that probing ability on ChatGPT and Llama2 are reliable, as they

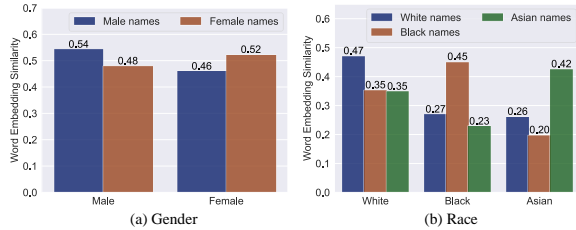


Figure 6: Word embeddings similarities between user names and sensitive attribute words.

consistently outperform random probing with a substantial margin. The experiment also verifies that different LLMs both have the ability to inference sensitive attributes from the non-sensitive attribute in terms of their wide world knowledge.

5.2 Word Embedding Similarities.

Secondly, we aim to investigate whether LLMs learn a close embedding between popular names and their sensitive attributes to determine if LLMs capture their relationships at a more fine-grained level. Since we cannot get embeddings from black-box LLMs ChatGPT, we only utilize the white-box LLM Llama 2 to conduct the experiments. We extract the word embeddings from the embedding table and average the sub-word embeddings.

We compute the distance of two embeddings based on cosine similarities $\cos(\cdot)$. Formally, the similarities between the sensitive attribute s and all non-sensitive attributes $[\mathcal{N}_s]_{s' \in \mathcal{S}}$ are: $\text{Softmax}([\cos(e_s, \sum_{n \in \mathcal{N}_{s'}} e_n / |\mathcal{N}_{s'}|)]_{s' \in \mathcal{S}})$, where e_s, e_n denote the word embeddings of sensitive attribute s and non-sensitive attribute n .

From Figure 6, it is evident that at the word level, non-sensitive attributes such as user names exhibit a significant correlation with sensitive attributes. This suggests that during the pre-training phase, LLMs can effectively learn and exploit these correlations, resulting in unfair ranking outcomes.

6 Implicit Ranking Unfairness Mitigation

In this section, we propose a fair data argumentation method to mitigate implicit ranking unfairness. We employ the 2SLS procedure (Kmenta, 2010) to remove the noise in non-sensitive attributes. After that, we can conduct data augmentation effectively by utilizing the top-N different feature sets that exhibit the most serious unfair behaviors in ranking.

Table 3: Unfairness degree compared ranking models learned collaborative information from and the implicit ranking unfairness of different versions of ChatGPT. The metric is U-NDCG@5. ‘‘Improv.’’ denotes the percentage of ChatGPT’s implicit user unfairness exceeding the highest degree of unfairness brought from collaborative information.

Models	DCN	GRU4Rec	ChatGPT	Improv.
News				
Gender	0.104	0.016	0.203	95.1%
Race	0.158	0.231	0.319	38.1%
Continent	0.324	0.158	0.711	119.4%
Jobs				
Gender	0.08	0.137	0.220	60.6%
Race	0.043	0.110	0.479	335%
Continent	0.139	0.115	0.798	474.1%

6.1 Stage-1.

In the first stage, we utilize pair-wise regression to train a RankNet (Borges et al., 2005), which aims to select user names that can be easily inferred from their demographic information.

In the ranking tasks, we take into account the order of the generated text within the ranking list. Ranking task implies a higher position in the ranking list L_K signifies greater importance for the associated item (Craswell et al., 2008). Therefore, we aim to investigate how LLM can infer demographic attributes through the patterns of ranking orders. Similarly, we also formulate this problem as a multi-classification task, where the class number corresponds to the demographic size $|\mathcal{S}|$.

Then, every item pair (i_j^n, i_m^n) is constructed from the ranking list $L_K(n, l)$, which takes n as a proxy for the demographic attribute in the prompts (Figure 1), where $i_j^n, i_m^n \in L_K(n, l)$ is the item in the j -th and m -th position of the ranking list, respectively with $m > j$. The pair reveals the ranking patterns in the ranking list.

Given the training data, we train the pair-wise regression network using the RankNet (Borges et al., 2005) with the loss function as

$$l^{\text{pair}} = \mathbb{E}_l \left[\sum_{s \in \mathcal{S}} \sum_{n \in \mathcal{N}_s} \sum_{j=1}^{K-1} \sum_{m=j+1}^K \text{CE} \left(z, \hat{z}^{\text{pair}}(i_j^n, i_m^n) \right) \right], \quad (3)$$

where the loss is a expectation among different sample i and $z \in \mathbb{R}^{|\mathcal{S}|}$ is the one-hot encoding representation of true demographic label s , and $\hat{z}^{\text{pair}}(i_j, i_m) \in \mathbb{R}^{|\mathcal{S}|}$ is computed through RankNet:

$$\hat{z}^{\text{pair}}(i_j, i_m) = \text{MLP} \left(e(i_j) \| e(i_m); \theta^{\text{pair}} \right),$$

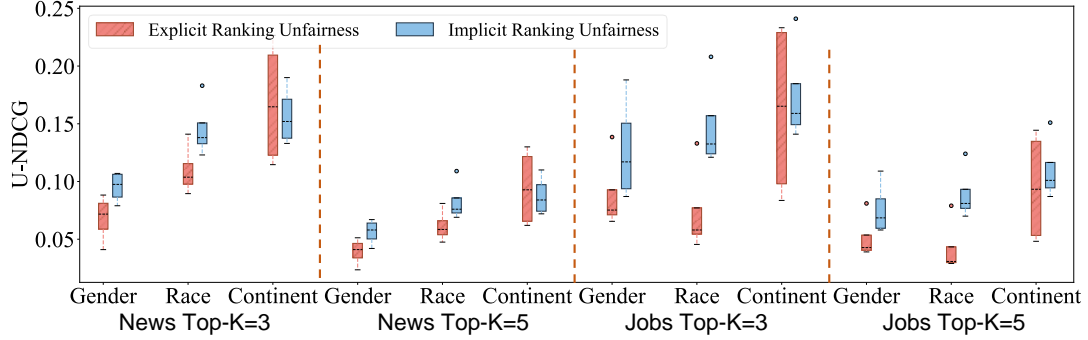


Figure 7: Comparing the averaged discrimination degrees (U-NDCG@3 and U-NDCG@5) of different versions of ChatGPT and Llama 2 under three demographic types (Gender, Race, and Continent) for news and job domain.

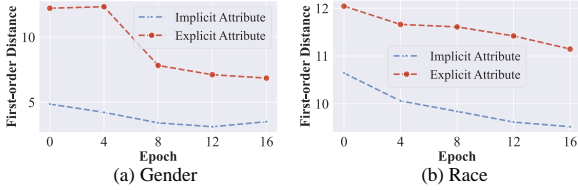


Figure 8: The first-order distance between embeddings of implicit attributes (such as user names) and embeddings of explicit attributes is measured during the tuning epochs of our method on News datasets.

where \parallel is the concat function for two vectors and θ^{pair} is the parameter of MLP network and $e(i)$ can be obtained by averaging the hidden embeddings of Llama2 to encode the textual item i as a vector.

6.2 Stage-2

In the second stage, after deciding the parameters of RankNet, we will decide the \mathcal{N}'_s for all sensitive group s to conduct data argumentation. Specifically, we will replace each non-sensitive attribute $n \in \mathcal{N}'_s$ to the “[demographic]” placeholder in Figure 1. In this way, one ranking sample can be augmented into $\sum_{s \in \mathcal{S}} |\mathcal{N}'_s|$ samples and feed these samples into instruction tuning phases of LLMs-based ranking tasks (Bao et al., 2023a). Specifically, we will choose the N non-sensitive attributes n of each sensitive group s . \mathcal{N}'_s is defined as:

$$\mathcal{N}'_s = \arg \max_{\mathcal{L} \in \mathcal{N}_s, |\mathcal{L}|=N} \sum_{n \in \mathcal{L}} \mathbb{E}_{j < m} [\text{CE}(z, \hat{z}^{\text{pair}}(i_j^n, i_m^n))]. \quad (4)$$

7 Experiments

In this section, we will conduct experiments to show the effectiveness of our methods.

7.1 Settings

The dataset and evaluation details are the same as Section 3.3. Due to the constraint of ChatGPT-

series API, we only utilize Lora (Hu et al., 2021) techniques to conduct instruction tuning for ranking tasks (Bao et al., 2023a) on Llama2 by employing different fairness strategies. The experiments were conducted under four NVIDIA A5000.

For the baseline, we compare four common-used types of methods to mitigate unfairness in LLMs: (1) **Self-Align**: following the practice in Sun et al. (2023), we utilize ChatGPT-3.5 (stronger LLM) to generate more reliable and fair responses to user’s queries and fine-tune the original Llama2 with the high-quality self-aligned responses. (2) **Re-Weight**: following (Jiang et al., 2024), during the tuning phase, we set the weight to be inversely proportional to the popularity of the item. (3) **Data-Argument** (Ghanbarzadeh et al., 2023): we replace the “[Demographic]” placeholder with explicit sensitive attribute as illustrated in Section 3.3. (4) **Prompt-Tuning** (Chisca et al., 2024): we utilize the prompt-tuning techniques to learn a fairness prompt to decrease the unfairness behaviors.

7.2 Experimental Results

In our experimental results, we mainly compare the unfairness of the most common sensitive attributes: gender and race. For the continent, we also observe a similar tendency.

In Table 4, it becomes evident that our method significantly outperforms the baselines across all datasets and sensitive attributes, encompassing different top-K ranking sizes. The experiments conclusively demonstrate that our method can mitigate the implicit ranking unfairness effectively.

7.3 Experimental Analysis

In this section, we will analyze why our method can mitigate implicit ranking unfairness. In Figure 8, we use TSNE (Van der Maaten and Hinton, 2008) to reduce the dimensionality of the vectors

Table 4: Unfairness degree (U-NDCG) compared between different models. “Improv.” denotes the percentage of implicit ranking unfairness exceeding the highest degree of implicit unfairness of baselines. Bold numbers mean the improvements over the best baseline are statistically significant (t-tests and p -value < 0.05).

model/domain	News				Jobs			
	gender		race		gender		race	
	top-3	top-5	top-3	top-5	top-3	top-5	top-3	top-5
Self-Align	0.0671	0.0379	0.0848	0.0471	0.0814	0.0464	0.1069	0.0627
Re-Weight	0.0751	0.0412	0.0807	0.0475	0.0536	0.0297	0.0501	0.0267
Data-Argument	0.0886	0.0498	0.0620	0.0363	0.0471	0.0264	0.0434	0.0235
Prompt-Tuning	0.0504	0.0276	0.0534	0.0297	0.0580	0.0344	0.0805	0.0459
Ours	0.0424*	0.0219*	0.0526*	0.0287*	0.0406*	0.0226*	0.0356*	0.0190*
Improv.	52.14%	56.02%	37.97%	39.57%	50.12%	51.29%	66.69%	69.69%

and calculate the distances between them to assess whether the large model reduces the distance between different groups of sensitive attributes in the ranking task.

From Figure 8, we can observe that using implicit attributes for data augmentation not only reduces the embedding distances between different implicit attributes but also brings embeddings of explicit attributes (such as “Male, Female”) closer together. In this way, the LLM-based ranking model will find it difficult to infer demographic attributes from user names, thereby effectively achieving ranking fairness.

8 Related Work

Recently, researchers have discovered that LLMs can exhibit discriminatory behaviors (Gallegos et al., 2023). In previous discrimination evaluation settings, researchers often measure stereotype sentence pairs that only differ in the sensitive attribute. For example, they often adapt terms “Male” and “Female” (Nangia et al., 2020; Delobelle et al., 2022; Gallegos et al., 2023) and for Race, they often substitute terms “Black”, “White” and “Asian” (Zhang et al., 2023b; Tamkin et al., 2023). Among the allocational harms, previous studies found that LLMs often exhibit discrimination against certain groups. For example, Salinas et al. (2023); de Vassimon Manela et al. (2021); McGee (2023); Thakur et al. (2023); Bolukbasi et al. (2016) discovered that LLMs will generate discriminatory content for disadvantaged gender. (Zhang et al., 2023b) show recommendation outcomes may discriminate against certain groups, see also (Rozado, 2023; Hutchinson et al., 2020). In our research, we mainly utilize the counterfactual fairness concept to measure the *implicit unfairness*

of LLMs-based recommendation.

There are some works that try to mitigate unfairness problems in LLMs. For example, RLHF (Ouyang et al., 2022) and RLAIIF (Bai et al., 2022) try to utilize reinforcement learning to align LLMs with human values. Generally, to address the imbalance in the original dataset against certain groups, some work (Ghanbarzadeh et al., 2023; Zhang et al., 2023b; Lu et al., 2020) create matched pairs (e.g., male or female) to ensure a more equitable dataset and other methods (Dixon et al., 2018; Sun et al., 2022) add non-toxic examples for groups. Other approaches (Orgad and Belinkov, 2022; Deldjoo and di Noia, 2024) suggest the use of down-weighting samples containing social group or discriminated information as a re-sampling strategy. While some method proposes to utilize the prompt-tuning method to learn a fair-aware prompt (Hua et al., 2023; Chisca et al., 2024). Moreover, other studies (Raffel et al., 2020; Ngo et al., 2021) propose to filter out and remove discriminated or taxonomic content from datasets.

9 Conclusion

In conclusion, our findings show that LLMs-based ranking models exhibit serious implicit unfairness. This implies that, even when sensitive attributes are not explicitly provided, LLMs can still exhibit discriminatory ranking behaviors. Regarding the root causes, we find that LLMs’ capability to deduce sensitive attributes from non-sensitive attributes contributes to the collection of unfair datasets during pre-training. Finally, we propose a new method to mitigate such unfairness effectively by utilizing fair-aware data augmentation. Our research emphasizes the necessity of identifying and moderating implicit ranking unfairness in existing LLMs.

581 Limitations

582 Finally, in our paper, we mainly utilize ChatGPT,
583 and Llama2 as our evaluation LLMs and only test
584 the discrimination behaviors against demographic
585 information in recommendation tasks. Meanwhile,
586 we currently only select user names and user emails
587 as the implicit attribute. However, different LLMs
588 and different discrimination behaviors may exhibit
589 different forms of implicit unfairness. This paper
590 serves as a valuable illustration to the community,
591 emphasizing the importance of careful considera-
592 tion when assessing the discrimination behaviors
593 in LLMs.

594 Ethics Statement

595 This study is a retrospective analysis conducted on
596 publicly available datasets with research-oriented
597 licenses, involving neither human participants nor
598 the requirement for informed consent. All results
599 generated by LLMs are utilized for offline analysis
600 by the authors and remain invisible to real-world
601 users, ensuring no actual social impact. User pro-
602 files used in the experiments, including names, gen-
603 ders, races, and nationalities, are simulated, and all
604 user identities have been completely anonymized.
605 The primary objective of this study is to enhance
606 the fairness of LLMs, aligning with the principles
607 of responsible and ethical usage.

608 References

609 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
610 Amanda Askell, Jackson Kernion, Andy Jones,
611 Anna Chen, Anna Goldie, Azalia Mirhoseini,
612 Cameron McKinnon, et al. 2022. Constitutional
613 ai: Harmlessness from ai feedback. *arXiv preprint*
614 *arXiv:2212.08073*.

615 Keqin Bao, Jizhi Zhang, Wenjie Wang, Yang Zhang,
616 Zhengyi Yang, Yancheng Luo, Fuli Feng, Xiang-
617 naan He, and Qi Tian. 2023a. A bi-step grounding
618 paradigm for large language models in recommenda-
619 tion systems. *arXiv preprint arXiv:2308.08434*.

620 Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang,
621 Fuli Feng, and Xiangnan He. 2023b. Tallrec: An
622 effective and efficient tuning framework to align
623 large language model with recommendation.
624 *arXiv preprint arXiv:2305.00447*.

625 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou,
626 Venkatesh Saligrama, and Adam T Kalai. 2016. Man
627 is to computer programmer as woman is to home-
628 maker? debiasing word embeddings. *Advances in*
629 *neural information processing systems*, 29.

Hannah Brown, Katherine Lee, Fatemehsadat
Mireshghallah, Reza Shokri, and Florian Tramèr.
2022. What does it mean for a language model
to preserve privacy? In *Proceedings of the 2022*
ACM Conference on Fairness, Accountability, and
Transparency, pages 2280–2292. 630
631
632
633
634
635

Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier,
Matt Deeds, Nicole Hamilton, and Greg Hullender.
2005. Learning to rank using gradient descent. In
Proceedings of the 22nd international conference on
Machine learning, pages 89–96. 636
637
638
639
640

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
Cunxiang Wang, Yidong Wang, et al. 2023. A sur-
vey on evaluation of large language models. *ACM*
Transactions on Intelligent Systems and Technology. 641
642
643
644
645

Andrei-Victor Chisca, Andrei-Cristian Rad, and
Camelia Lemnaru. 2024. [Prompting fairness: Learning prompts for debiasing large language models](#). In
Proceedings of the Fourth Workshop on Language
Technology for Equality, Diversity, Inclusion, pages
52–62, St. Julian’s, Malta. Association for Computa-
tional Linguistics. 646
647
648
649
650
651
652

Nick Craswell, Onno Zoeter, Michael Taylor, and Bill
Ramsey. 2008. An experimental comparison of click
position-bias models. In *Proceedings of the 2008*
international conference on web search and data
mining, pages 87–94. 653
654
655
656
657

Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu,
Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang,
and Jun Xu. 2023. Uncovering chatgpt’s capabilities
in recommender systems. *Recsys*. 658
659
660
661

Sunhao Dai, Chen Xu, Shicheng Xu, Liang Pang, Zhen-
hua Dong, and Jun Xu. 2024. Unifying bias and
unfairness in information retrieval: A survey of chal-
lenges and opportunities with large language models.
arXiv preprint arXiv:2404.11457. 662
663
664
665
666

Daniel de Vassimon Manela, David Errington, Thomas
Fisher, Boris van Breugel, and Pasquale Minervini.
2021. Stereotype and skew: Quantifying gender bias
in pre-trained and fine-tuned language models. In
Proceedings of the 16th Conference of the European
Chapter of the ACL: Main Volume, pages 2232–2242. 667
668
669
670
671
672

Yashar Deldjoo and Tommaso di Noia. 2024. [Cfair-llm: Consumer fairness evaluation in large-language model recommender system](#). 673
674
675

Pieter Delobelle, Ewoenam Kwaku Tokpo, Toon
Calders, and Bettina Berendt. 2022. Measuring fair-
ness with biased rulers: A comparative study on bias
metrics for pre-trained language models. In *NAACL*
2022: the 2022 Conference of the North American
chapter of the Association for Computational Lin-
guistics: human language technologies, pages 1693–
1706. 676
677
678
679
680
681
682
683

684	Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In <i>Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society</i> , pages 67–73.	738
685		739
686		
687		
688		
689	Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2023. Bias and fairness in large language models: A survey .	
690		
691		
692		
693		
694	Somayeh Ghanbarzadeh, Yan Huang, Hamid Palangi, Radames Cruz Moreno, and Hamed Khanpour. 2023. Gender-tuning: Empowering fine-tuning for debiasing pre-trained language models. <i>arXiv preprint arXiv:2307.10522</i> .	
695		
696		
697		
698		
699	Wes Gurnee and Max Tegmark. 2023. Language models represent space and time .	
700		
701	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models .	
702		
703		
704		
705	Wenyue Hua, Yingqiang Ge, Shuyuan Xu, Jianchao Ji, and Yongfeng Zhang. 2023. Up5: Unbiased foundation model for fairness-aware recommendation. <i>arXiv preprint arXiv:2305.12090</i> .	
706		
707		
708		
709	Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in nlp models as barriers for persons with disabilities. In <i>Proceedings of the 58th Annual Meeting of the ACL</i> , pages 5491–5501.	
710		
711		
712		
713		
714	Meng Jiang, Keqin Bao, Jizhi Zhang, Wenjie Wang, Zhengyi Yang, Fuli Feng, and Xiangnan He. 2024. Item-side fairness of large language model-based recommendation system. In <i>Proceedings of the ACM on Web Conference 2024</i> , pages 4717–4726.	
715		
716		
717		
718		
719	Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In <i>International Conference on Machine Learning</i> , pages 10697–10707. PMLR.	
720		
721		
722		
723	Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. <i>Learning and individual differences</i> , 103:102274.	
724		
725		
726		
727		
728		
729		
730	Jan Kmenta. 2010. Mostly harmless econometrics: An empiricist’s companion.	
731		
732	Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. <i>Advances in neural information processing systems</i> , 30.	
733		
734		
735	Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, Juntao Tan, Shuchang Liu, and Yongfeng Zhang. 2023. Fairness in recommendation: Foundations, methods, and applications. <i>ACM Transactions on Intelligent Systems and Technology</i> , 14(5):1–48.	738
736		739
737		
	Tie-Yan Liu et al. 2009. Learning to rank for information retrieval. <i>Foundations and Trends® in Information Retrieval</i> , 3(3):225–331.	740
		741
		742
	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2020. Gender bias in neural natural language processing. <i>Logic, language, and security: essays dedicated to Andre Scedrov on the occasion of his 65th birthday</i> , pages 189–202.	743
		744
		745
		746
		747
	Robert W McGee. 2023. Gender discrimination arguments and non sequiturs: A chatgpt essay. <i>Available at SSRN 4413432</i> .	748
		749
		750
	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. <i>arXiv preprint arXiv:2010.00133</i> .	751
		752
		753
		754
	Helen Ngo, Cooper Raterink, João GM Araújo, Ivan Zhang, Carol Chen, Adrien Morisot, and Nicholas Frosst. 2021. Mitigating harm in language models with conditional-likelihood filtration. <i>arXiv preprint arXiv:2108.07790</i> .	755
		756
		757
		758
		759
	Hadas Orgad and Yonatan Belinkov. 2022. Blind: Bias removal with no demographics. <i>arXiv preprint arXiv:2212.10563</i> .	760
		761
		762
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	763
		764
		765
		766
		767
		768
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	769
		770
		771
		772
		773
		774
	Konstantinos I Roumeliotis and Nikolaos D Tselikas. 2023. Chatgpt and open-ai models: A preliminary review. <i>Future Internet</i> , 15(6):192.	775
		776
		777
	David Rozado. 2023. The political biases of chatgpt. <i>Social Sciences</i> , 12(3):148.	778
		779
	Abel Salinas, Parth Vipul Shah, Yuzhong Huang, Robert McCormack, and Fred Morstatter. 2023. The unequal opportunities of large language models: Revealing demographic bias through job recommendations. <i>arXiv preprint arXiv:2308.02053</i> .	780
		781
		782
		783
		784
	Hao Sun, Zhixin Zhang, Fei Mi, Yasheng Wang, Wei Liu, Jianwei Cui, Bin Wang, Qun Liu, and Minlie Huang. 2022. Moraldial: A framework to train and evaluate moral dialogue systems via moral discussions. <i>arXiv preprint arXiv:2212.10720</i> .	785
		786
		787
		788
		789

790	Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin	Yijia Xiao, Yiqiao Jin, Yushi Bai, Yue Wu, Xianjun	845
791	Zhang, Zhenfang Chen, David Cox, Yiming Yang,	Yang, Xiao Luo, Wenchao Yu, Xujiang Zhao, Yanchi	846
792	and Chuang Gan. 2023. Principle-driven self-	Liu, Haifeng Chen, et al. 2023. Large language mod-	847
793	alignment of language models from scratch with min-	els can be good privacy protection learners. <i>arXiv</i>	848
794	imal human supervision.	<i>preprint arXiv:2310.02469.</i>	849
795	Alex Tamkin, Amanda Askill, Liane Lovitt, Esin	Sirui Yao and Bert Huang. 2017. Beyond parity: Fair-	850
796	Durmus, Nicholas Joseph, Shauna Kravec, Karina	ness objectives for collaborative filtering. <i>Advances</i>	851
797	Nguyen, Jared Kaplan, and Deep Ganguli. 2023.	<i>in neural information processing systems</i> , 30.	852
798	Evaluating and mitigating discrimination in language	An Zhang, Leheng Sheng, Yuxin Chen, Hao Li, Yang	853
799	model decisions. <i>arXiv preprint arXiv:2312.03689.</i>	Deng, Xiang Wang, and Tat-Seng Chua. 2023a. On	854
800	Yong Kiam Tan, Xinxing Xu, and Yong Liu. 2016. Im-	generative agents in recommendation. <i>arXiv preprint</i>	855
801	proved recurrent neural networks for session-based	<i>arXiv:2310.10108.</i>	856
802	recommendations. In <i>Proceedings of the 1st work-</i>	Jizhi Zhang, Keqin Bao, Yang Zhang, Wenjie Wang,	857
803	<i>shop on deep learning for recommender systems</i> ,	Fuli Feng, and Xiangnan He. 2023b. Is chatgpt fair	858
804	pages 17–22.	for recommendation? evaluating fairness in large	859
805	Himanshu Thakur, Atishay Jain, Praneetha Vaddamanu,	language model recommendation. <i>arXiv preprint</i>	860
806	Paul Pu Liang, and Louis-Philippe Morency. 2023.	<i>arXiv:2305.07609.</i>	861
807	Language models get a gender makeover: Mitigating	Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin	862
808	gender bias with few-shot data interventions. <i>arXiv</i>	Zhao, Leyu Lin, and Ji-Rong Wen. 2023c. Recom-	863
809	<i>preprint arXiv:2306.04597.</i>	mendation as instruction following: A large language	864
810	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	model empowered recommendation approach. <i>arXiv</i>	865
811	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>preprint arXiv:2305.07001.</i>	866
812	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti		
813	Bhosale, et al. 2023. Llama 2: Open founda-		
814	tion and fine-tuned chat models. <i>arXiv preprint</i>		
815	<i>arXiv:2307.09288.</i>		
816	Laurens Van der Maaten and Geoffrey Hinton. 2008.		
817	Visualizing data using t-sne. <i>Journal of machine</i>		
818	<i>learning research</i> , 9(11).		
819	Ivan Vulić, Edoardo Maria Ponti, Robert Litschko,		
820	Goran Glavaš, and Anna Korhonen. 2020. Prob-		
821	ing pretrained language models for lexical semantics.		
822	<i>arXiv preprint arXiv:2010.05731.</i>		
823	Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang.		
824	2017. Deep & cross network for ad click predictions.		
825	In <i>Proceedings of the ADKDD'17</i> , pages 1–7.		
826	Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan		
827	Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing Xie,		
828	Jianfeng Gao, Winnie Wu, et al. 2020. Mind: A		
829	large-scale dataset for news recommendation. In		
830	<i>ACL</i> , pages 3597–3606.		
831	Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang,		
832	Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu,		
833	Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen.		
834	2023a. A survey on large language models for rec-		
835	ommendation.		
836	Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang		
837	Liu, Qing-Long Han, and Yang Tang. 2023b. A brief		
838	overview of chatgpt: The history, status quo and		
839	potential future development. <i>IEEE/CAA Journal of</i>		
840	<i>Automatica Sinica</i> , 10(5):1122–1136.		
841	Yongkai Wu, Lu Zhang, and Xintao Wu. 2019. Coun-		
842	terfactual fairness: Unidentification, bound and algo-		
843	rithm. In <i>Proceedings of the twenty-eighth interna-</i>		
844	<i>tional joint conference on Artificial Intelligence.</i>		