
Investigating the Fairness of Large Language Models for Predictions on Tabular Data

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent literature has suggested the potential of using large language models (LLMs)
2 to make predictions for tabular tasks. However, LLMs have been shown to exhibit
3 harmful social biases that reflect the stereotypes and inequalities present in the
4 society. To this end, as well as the widespread use of tabular data in many high-
5 stake applications, it is imperative to explore the following questions: what sources
6 of information do LLMs draw upon when making predictions for tabular tasks;
7 whether and to what extent are LLM predictions for tabular tasks influenced
8 by social biases and stereotypes; and what are the consequential implications
9 for fairness? Through a series of experiments, we delve into these questions
10 and show that LLMs tend to inherit social biases from their training data which
11 significantly impact their fairness in tabular prediction tasks. Furthermore, our
12 investigations show that in the context of bias mitigation, though in-context learning
13 and fine-tuning have a moderate effect, the fairness metric gap between different
14 subgroups is still larger than that in traditional machine learning models, such
15 as Random Forest and shallow Neural Networks. This observation emphasizes
16 that the social biases are inherent within the LLMs themselves and inherited from
17 their pre-training corpus, not only from the downstream task datasets. Besides,
18 we demonstrate that label-flipping of in-context examples can significantly reduce
19 biases, further highlighting the presence of inherent bias within LLMs.

20 1 Introduction

21 Many recent works propose to use large language models (LLMs) for tabular prediction (Slack &
22 Singh, 2023; Hegselmann et al., 2023), where the tabular data is serialized as natural language and
23 provided to LLMs with a short description of the task to solicit predictions. Despite the comprehensive
24 examination of fairness considerations within conventional machine learning approaches applied
25 to tabular tasks (Bellamy et al., 2018), the exploration of fairness-related issues in the context of
26 employing LLMs for tabular predictions remains a relatively underexplored domain.

27 Previous research has shown that LLMs, such as GPT-3 (Brown et al., 2020), GPT-3.5, GPT-
28 4 (OpenAI, 2023) can exhibit harmful social biases (Abid et al., 2021a; Basta et al., 2019), which
29 may even worsen as the models become larger in size (Askell et al., 2021; Ganguli et al., 2022).
30 These biases are a result of the models being trained on text generated by humans that presumably
31 includes many examples of humans exhibiting harmful stereotypes and discrimination and reflects
32 the biases and inequalities present in society (Bolukbasi et al., 2016; Zhao et al., 2017), which can
33 lead to perpetuation of discrimination and stereotype (Abid et al., 2021a; Bender et al., 2021).

34 Considering that tabular data finds extensive use in high-stakes domains (Grinsztajn et al., 2022)
35 where information is typically structured in tabular formats as a natural byproduct of relational
36 databases (Borisov et al., 2022), it is of paramount importance to thoroughly examine the fairness

37 implications of utilizing LLMs for predictions on tabular data. In this paper, we conduct a series of
38 investigation centered around this critical aspect, with the goal of discerning the underlying informa-
39 tion sources upon which LLMs rely when making tabular predictions. Through this exploration, our
40 investigation aims to ascertain whether, and to what degree, LLMs are susceptible to being influenced
41 by social biases and stereotypes in the context of tabular data predictions.

42 Through experiments using GPT-3.5 to make predictions for tabular data in a zero-shot setting,
43 we demonstrate that LLMs exhibit significant social biases (Section 4). This evidence confirms
44 that LLMs inherit social biases from their training corpus and tend to rely on these biases when
45 making predictions for tabular data. Furthermore, we demonstrate that providing LLMs with few-shot
46 examples (in-context learning) or fine-tuning them on the entire training dataset both exhibit moderate
47 effect on bias mitigation (Sections 5 and 6). Nevertheless, the achieved fairness levels remain below
48 what is typically attained with traditional machine learning methods, including Random Forests
49 and shallow Neural Networks, once again underscoring the presence of inherent bias in LLMs.
50 Additionally, our investigation further reveals that flipping the labels of the in-context examples
51 significantly narrows the gap in fairness metrics across different subgroups, but comes at the expected
52 cost of a reduction in predictive performance. This finding, in turn, further emphasizes and reaffirms
53 the indication of inherent bias present in LLMs (Section 5). Additionally, we further show that while
54 resampling the training set is a known and effective method for reducing biases in traditional machine
55 learning methods like Random Forests and shallow Neural Networks, it proves to be less effective
56 when applied to LLMs (Section 6).

57 These collective findings underscore the significant influence of social biases on LLMs’ performance
58 in tabular predictions. These biases significantly undermines the fairness and poses substantial
59 potential risks for using LLMs on tabular data, especially considering that tabular data is extensively
60 used in high-stakes domains, highlighting the need for more advanced and tailored strategies to
61 address these biases effectively. Straightforward methods like in-context learning and data resampling
62 may not be sufficient in this context.

63 2 Related work

64 2.1 Fairness and Social Biases in LLMs

65 Fairness is highly desirable for ensuring the credibility and trustworthiness of algorithms. It has
66 been demonstrated that unfair algorithms can reflect societal biases in their decision-making pro-
67 cesses (Bender et al., 2021; Bommasani, 2021), primarily stemming from the biases present in
68 their training data (Caliskan et al., 2017; Zhao et al., 2017). LLMs, pre-trained on vast natural
69 language datasets, are particularly susceptible to inheriting these social biases and have been shown
70 to exhibit biases related to gender (Lucy & Bamman, 2021), religion (Abid et al., 2021b) and lan-
71 guage variants (Ziems et al., 2023; Liu et al., 2023). These social biases can lead to perpetuation
72 of discrimination and stereotype (Abid et al., 2021a; Bender et al., 2021; Weidinger et al., 2021).
73 While recent literature has made strides in addressing these issues, there still exists a significant gap
74 in comprehensively assessing fairness in LLMs and its mitigation strategies for tabular data.

75 2.2 Tabular Tasks and LLM for Tabular Data

76 Tabular data extensively exist in many domains (Shwartz-Ziv & Armon, 2021). Previous works
77 propose to utilize self-supervised deep techniques for tabular tasks (Yin et al., 2020; Arik & Pfister,
78 2021), which, however, still underperform ensembles of gradient boosted trees in the fully supervised
79 setting (Grinsztajn et al., 2022). Recent approaches by Hagselmann et al. (2023); Slack & Singh
80 (2023) suggests serializing the tabular data as natural language, which is provided to LLM along with
81 a short task description to generate predictions for tabular tasks. However, tabular data plays a crucial
82 role in numerous safety-critical and high-stakes domains (Borisov et al., 2022; Grinsztajn et al.,
83 2022), which makes the fairness particularly crucial when employing LLMs for making predictions
84 on tabular data, especially considering the inherent social biases present in LLMs. Despite the
85 importance, this still remains largely unexplored. To the best of our knowledge, we regard our work
86 as one of the most comprehensive investigations into the fairness issues arising when using LLMs for
87 predictions on tabular data.

88 3 Experimental Setup

89 **Models** In our work, we focus our experiments on GPT-3.5 (engine GPT-3.5-turbo). Furthermore,
90 we also compare its performance with conventional machine learning models in order to gain insight
91 into the propagation of biases. For this, we employ two widely used models for tabular data i.e.,
92 Random Forests (RF) and a shallow Neural Network (NN) of 3 layers. We provide additional
93 implementation details for these two models in the Appendix C.

94 **Datasets and Protected Attributes** To explore the fairness of LLMs in making predictions for
95 tabular data, we utilize the following widely used tabular datasets for assessing the fairness of
96 traditional ML models: *Adult Income (Adult)* Dataset (Becker & Kohavi, 1996) and *Correctional*
97 *Offender Management Profiling for Alternative Sanctions (COMPAS)* Dataset (Larson et al., 2016).
98 A detailed description for each dataset and each feature of the considered datasets is provided in
99 Appendix A.

100 **Serialization and Prompt Templates** To employ the LLM for making predictions on these tabular
101 datasets, each data point is first serialized as text. Following previous works on LLM for tabular
102 predictions (Hegselmann et al., 2023; Slack & Singh, 2023), we format the feature names and values
103 into strings as " $f_1 : x_1, \dots, f_d : x_d$ ", and prompt to LLM along with a task description.

104 **Evaluation Metrics** To assess fairness in the aforementioned datasets, we examine the disparity
105 between different subgroups of protected attributes using the following common fairness metrics:
106 accuracy, F1 score, statistical parity and equality of opportunity. We provide the detail for each
107 fairness metric in Appendix B

108 We run all the experiments 5 times and compute the mean and standard deviation.

109 4 Zero-Shot Prompting for Tabular Data

110 To explore the fairness of LLMs when making predictions on tabular data, we first conduct experi-
111 ments in a zero-shot setting. We assess the fairness metrics of the outcomes and examine whether
112 LLMs without any finetuning or few-shot examples would be influenced by social biases and stereo-
113 types for tabular predictions. In Tables 1 and 5, we present the evaluation of four fairness metrics, for
114 GPT-3.5 (engine GPT-3.5-turbo), RF and NN models on the Adult and COMPAS datasets, respec-
115 tively. For the Adult dataset, the subgroups *female* and *male* are assessed regarding the protected
116 attribute *sex*, identifying *female* as a disadvantaged group. In the COMPAS dataset, we evaluate *race*
117 as protected attributes, recognizing African American (AA) as the disadvantaged group.

118 It is notable that when utilizing LLMs to make predictions for tabular data directly, without any
119 fine-tuning or in-context learning, a significant fairness metric gap between the protected and non-
120 protected groups is observed for GPT-3.5 (highlighted in red). For instance, the EoO difference
121 between *male* and *female* on the *Adult* dataset reaches 0.483, indicating a substantial disadvantage for
122 the *female* group. Additionally, when compared with traditional methods like RF and NN, the bias in
123 zero-shot predictions made by GPT-3.5 is significantly larger for the Adult dataset. This observation
124 suggests an inherent gender bias in GPT-3.5. For COMPAS dataset, the racial bias in zero-shot setting
125 is comparatively lower than RF and NN but is still effectively high.

126 These findings demonstrate the tendency of LLMs to rely on social biases and stereotypes inherited
127 from their training corpus when applied to tabular data. This implies that using LLMs for predictions
128 on tabular data may incur significant fairness risks, including the potential to disproportionately
129 disadvantage marginalized communities as well as exacerbate social biases and stereotypes present in
130 society. This is particularly concerning given the widespread application of tabular data in high-stake
131 contexts, further magnifying the potential for harm.

132 5 Few-Shot Prompting for Tabular Data

133 Instead of directly utilizing LLMs for zero-shot tabular predictions, this section explores whether
134 including few-shot examples during prompting will reduce or amplify these biases. To delve deeper

			ACC	F1	SP	EoO	
GPT-3.5-turbo	Zero-Shot	<i>f</i>	0.898 _{0.001}	0.711 _{0.002}	0.065 _{0.001}	0.357 _{0.000}	
		<i>m</i>	0.742 _{0.002}	0.727 _{0.002}	0.464 _{0.003}	0.840 _{0.004}	
		<i>d</i>	0.157 _{0.002}	-0.016 _{0.002}	-0.399 _{0.003}	-0.483 _{0.004}	
	Few-shot	Regular	<i>f</i>	0.899 _{0.002}	0.735 _{0.003}	0.082 _{0.002}	0.429 _{0.000}
			<i>m</i>	0.781 _{0.003}	0.749 _{0.002}	0.339 _{0.003}	0.700 _{0.003}
			<i>d</i>	0.118 _{0.004}	-0.014 _{0.004}	-0.257 _{0.005} ↓	-0.271 _{0.003} ↓
		Label-flipping	<i>f</i>	0.682 _{0.004}	0.590 _{0.003}	0.396 _{0.006}	0.800 _{0.013}
			<i>m</i>	0.614 _{0.002}	0.605 _{0.002}	0.545 _{0.001}	0.763 _{0.003}
			<i>d</i>	0.068 _{0.004}	-0.015 _{0.004}	-0.148 _{0.006} ✓	0.037 _{0.014} ✓
	Finetuning	Regular	<i>f</i>	0.915 _{0.014}	0.773 _{0.036}	0.079 _{0.002}	0.476 _{0.048}
			<i>m</i>	0.799 _{0.005}	0.754 _{0.005}	0.269 _{0.036}	0.613 _{0.053}
			<i>d</i>	0.116 _{0.009}	0.020 _{0.039}	-0.190 _{0.035} ↓	-0.137 _{0.098} ↓
		Oversampling	<i>f</i>	0.913 _{0.016}	0.770 _{0.042}	0.081 _{0.004}	0.476 _{0.067}
			<i>m</i>	0.813 _{0.007}	0.780 _{0.003}	0.310 _{0.038}	0.702 _{0.048}
			<i>d</i>	0.100 _{0.013}	-0.010 _{0.041}	-0.229 _{0.030}	-0.226 _{0.077}
		Undersampling	<i>f</i>	0.912 _{0.015}	0.770 _{0.046}	0.086 _{0.006}	0.488 _{0.084}
			<i>m</i>	0.794 _{0.006}	0.751 _{0.001}	0.285 _{0.031}	0.631 _{0.044}
			<i>d</i>	0.118 _{0.021}	0.018 _{0.046}	-0.200 _{0.025}	-0.143 _{0.040}
RF	Regular	<i>f</i>	0.914 _{0.002}	0.767 _{0.006}	0.075 _{0.003}	0.457 _{0.010}	
		<i>m</i>	0.822 _{0.005}	0.783 _{0.005}	0.269 _{0.004}	0.652 _{0.004}	
		<i>d</i>	0.092 _{0.004}	-0.015 _{0.005}	-0.195 _{0.003}	-0.195 _{0.012}	
	Oversampling	<i>f</i>	0.912 _{0.006}	0.770 _{0.011}	0.084 _{0.005}	0.486 _{0.012}	
		<i>m</i>	0.824 _{0.002}	0.785 _{0.002}	0.270 _{0.003}	0.656 _{0.006}	
		<i>d</i>	0.087 _{0.005}	-0.015 _{0.01}	-0.185 _{0.004}	-0.170 _{0.011}	
	Undersampling	<i>f</i>	0.917 _{0.004}	0.776 _{0.011}	0.075 _{0.001}	0.471 _{0.018}	
		<i>m</i>	0.814 _{0.003}	0.771 _{0.004}	0.263 _{0.002}	0.627 _{0.009}	
		<i>d</i>	0.103 _{0.005}	0.005 _{0.011}	-0.187 _{0.001}	-0.156 _{0.018}	
NN	Regular	<i>f</i>	0.917 _{0.003}	0.778 _{0.019}	0.081 _{0.016}	0.490 _{0.068}	
		<i>m</i>	0.819 _{0.006}	0.773 _{0.015}	0.250 _{0.045}	0.614 _{0.079}	
		<i>d</i>	0.098 _{0.005}	0.006 _{0.009}	-0.169 _{0.032}	-0.123 _{0.033}	
	Oversampling	<i>f</i>	0.916 _{0.004}	0.794 _{0.013}	0.100 _{0.016}	0.562 _{0.058}	
		<i>m</i>	0.813 _{0.012}	0.774 _{0.008}	0.286 _{0.044}	0.663 _{0.056}	
		<i>d</i>	0.103 _{0.011}	0.020 _{0.018}	-0.186 _{0.030}	-0.102 _{0.038}	
	Undersampling	<i>f</i>	0.904 _{0.005}	0.748 _{0.014}	0.084 _{0.007}	0.452 _{0.030}	
		<i>m</i>	0.813 _{0.006}	0.774 _{0.005}	0.283 _{0.023}	0.659 _{0.031}	
		<i>d</i>	0.090 _{0.006}	-0.026 _{0.014}	-0.199 _{0.018}	-0.206 _{0.031}	

Table 1: **Fairness evaluation for Adult dataset.** This table depicts the evaluation of accuracy (ACC), F1 score (F1), statistical parity (SP), and equality of opportunity (EoO) metrics for the subgroup - *female* (*f*) and *male* (*m*) as well as the difference (*d*) between them. We list the protected group first. The significant fairness disparities are highlighted in red. Both in-context learning and finetuning can lead to bias reduction (indicated by ↓), and label-flipped in-context learning can further minimize bias (indicated by ✓).

135 into the influence of few-shot examples, we not only consider the regular in-context learning approach
136 in Section 5, but we also experiment by flipping the labels of the few-shot examples in Section 5.

137 **Regular In-Context Learning** Previous works have demonstrated that LLMs can learn the input-
138 label mappings in context (Akyürek et al., 2022; Xie et al., 2022; Von Oswald et al., 2023). However,
139 the influence of in-context learning on the fairness has not been thoroughly examined. For in-context
140 learning, the test example and task description, along with a few-shot examples, are provided to
141 the LLMs for generating the final predictions. The few-shot examples are inserted before the test
142 example in the prompt, as outlined in Section 3. We set the number of in-context examples as 50. For
143 each dataset, we randomly select the in-context examples from the training set for each test example.

144 In Tables 1, we demonstrate that the incorporation of few-shot examples brings about performance
145 improvements. Additionally, we observe that incorporating few-shot examples into prompting reduces
146 the fairness metric gap between different subgroups. However, a significant fairness issue still persists.
147 Moreover, the disparity in fairness metrics of in-context learning is more notable when compared to
148 traditional models, such as RF and NN. This highlights the inherent biases embedded within LLMs,
149 which are not solely derived from the task datasets.

150 **Label Flipping** To delve deeper into the sources of biases within LLMs, we further examine the
151 impact of the labels of in-context examples on fairness. As depicted in Tables 1 and 5, label flipping
152 significantly reduces biases across all evaluated datasets. And for all evaluated datasets, the difference
153 in statistical parity (SP) and equality of opportunity (EoO) is minimized with label-flipped in-context
154 learning. For example, the absolute gap of EoO on the Adult dataset decreases from 0.483 in zero-shot
155 prompting to 0.037, almost completely eliminating the bias. These findings further corroborates the
156 existence of inherent biases in LLMs.

157 However, flipped labels lead to a significant drop in predictive performance. Though previous research
158 suggests that the effectiveness of in-context learning predominantly stems from semantic priors,
159 rather than learning the input-label mappings (Min et al., 2022; Wei et al., 2023) and demonstrate
160 that the performance of in-context learning is barely affected even with flipped or random labels for
161 in-context examples, the focus of these works lies mainly on traditional natural language processing
162 tasks. In contrast, we observe that the labels of in-context examples hold substantial influence over
163 predictive performance in our unique setup, where LLMs are deployed for predictions on tabular data.
164 This could be attributed to the limited exposure of these models to tabular data during pre-training,
165 thereby amplifying the role of input-label mapping of in-context examples.

166 6 Finetuning for Tabular Data

167 Finally, we extend our investigation to assess if finetuning the models on the entire training set
168 could aid in diminishing the social biases in LLMs. For GPT-3.5, fine-tuning is executed using the
169 publicly released API from OpenAI. For RF and NN, we provide the training details in Appendix C.
170 In Tables 1 and 5, we show that finetuning effectively reduces unfairness in all datasets, making them
171 comparable and sometimes significantly better in terms of SP and EoO when compared to RF and
172 NN. For example, the absolute difference in EoO after finetuning on Adult dataset is 0.0714, which is
173 lower than 0.123 difference of a NN.

174 We further explore the potential of resampling, a method frequently employed to enhance fairness
175 in machine learning model training, particularly in scenarios where there is a significant class
176 imbalance or bias in the data. To this end, we evaluate two approaches: oversampling the minority
177 group and undersampling the majority group. As depicted in Tables 1 and 5, resampling fails to
178 mitigate the social biases in LLMs when making tabular predictions, even though we demonstrate
179 that oversampling generally reduces social biases for both RF and NN, except for a few instances
180 such as, oversampling in NN for adult dataset worsens the fairness.

181 Our finetuning experiments show that the social biases inherited from LLM’s pre-training data which
182 are evident when making predictions on tabular data, can sometimes be mitigated through finetuning.
183 Nevertheless, unlike the consistent outcomes typically seen in traditional machine learning models,
184 like RF and NN, data resampling does not consistently produce similar results for finetuning LLMs.

185 7 Conclusion

186 In this work, we thoroughly investigate the under-explored problem of fairness of large language
187 models (LLMs) for tabular tasks. We assess the inherent fairness displayed by LLMs, comparing
188 their performance in zero-shot learning scenarios against traditional machine learning models like
189 random forests (RF) and shallow neural networks (NN). Furthermore, we investigate how LLMs learn
190 and propagate social biases when subjected to few-shot in-context learning, label-flipped in-context
191 learning, fine-tuning, and data resampling techniques.

192 We find that LLMs tend to heavily rely on the social biases inherited from their pre-training data
193 when making predictions, which is a concerning issue. Moreover, we observe that few-shot in-context
194 learning can partially mitigate the inherent biases in LLMs, yet it cannot entirely eliminate them.
195 A significant fairness metric gap between different subgroups persists, and exceeds that observed
196 in RF and NN. This observation underscores the existence of biases within the LLMs themselves,
197 beyond just the task datasets. Additionally, label-flipping applied to the few-shot examples effectively
198 reverses the effects of bias, again corroborating the existence of inherent biases in LLMs. However,
199 as expected, this leads to a loss in predictive performance. Besides, our work reveals that while
200 fine-tuning can sometimes improve the fairness of LLMs, data resampling does not consistently yield
201 the same results, unlike what is typically observed in traditional machine learning models.

202 References

- 203 Abubakar Abid, Maheen Farooqi, and James Zou. Large language models associate muslims with
204 violence. *Nature Machine Intelligence*, 3(6):461–463, 2021a. doi: 10.1038/s42256-021-00359-2.
205 URL <https://doi.org/10.1038/s42256-021-00359-2>.
- 206 Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language
207 models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES
208 '21, pp. 298–306, New York, NY, USA, 2021b. Association for Computing Machinery. ISBN
209 9781450384735. doi: 10.1145/3461702.3462624.
- 210 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
211 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
212 2022.
- 213 Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *Proceedings of*
214 *the AAAI conference on artificial intelligence*, volume 35, pp. 6679–6687, 2021.
- 215 Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones,
216 Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernan-
217 dez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark,
218 Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for
219 alignment, 2021.
- 220 Christine Basta, Marta R. Costa-jussà, and Noe Casas. Evaluating the underlying gender bias in
221 contextualized word embeddings. In *Proceedings of the First Workshop on Gender Bias in Natural*
222 *Language Processing*, Sep 2019. doi: 10.18653/v1/w19-3805. URL [http://dx.doi.org/10.](http://dx.doi.org/10.18653/v1/w19-3805)
223 [18653/v1/w19-3805](http://dx.doi.org/10.18653/v1/w19-3805).
- 224 Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI:
225 <https://doi.org/10.24432/C5XW20>.
- 226 Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Kalpriya
227 Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, Seema Nagar,
228 Karthikeyan Natesan Ramamurthy, John Richards, Diptikalyan Saha, Prasanna Sattigeri, Moninder
229 Singh, Kush R. Varshney, and Yunfeng Zhang. AI Fairness 360: An extensible toolkit for
230 detecting, understanding, and mitigating unwanted algorithmic bias, October 2018. URL [https:](https://arxiv.org/abs/1810.01943)
231 [//arxiv.org/abs/1810.01943](https://arxiv.org/abs/1810.01943).
- 232 Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the
233 dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM*
234 *Conference on Fairness, Accountability, and Transparency*, Mar 2021. doi: 10.1145/3442188.
235 3445922. URL <http://dx.doi.org/10.1145/3442188.3445922>.
- 236 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T
237 Kalai. Man is to computer programmer as woman is to homemaker? debiasing
238 word embeddings. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett
239 (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates,
240 Inc., 2016. URL [https://proceedings.neurips.cc/paper_files/paper/2016/file/](https://proceedings.neurips.cc/paper_files/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf)
241 [a486cd07e4ac3d270571622f4f316ec5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf).
- 242 R Bommasani. Opportunities and risks of foundation models, 2021. [https://openai.com/](https://openai.com/reports/foundation-models/)
243 [reports/foundation-models/](https://openai.com/reports/foundation-models/).
- 244 Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji
245 Kasneci. Deep neural networks and tabular data: A survey. *IEEE Transactions on Neural Networks*
246 *and Learning Systems*, pp. 1–21, 2022. doi: 10.1109/TNNLS.2022.3229161.
- 247 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
248 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
249 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
250 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Ben-
251 jamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and

- 252 Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell,
253 M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp.
254 1877–1901. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)
255 [files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 256 Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. Semantics derived automatically from
257 language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017. doi: 10.1126/
258 [science.aal4230](https://www.science.org/doi/abs/10.1126/science.aal4230). URL <https://www.science.org/doi/abs/10.1126/science.aal4230>.
- 259 Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom
260 Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac
261 Hatfield-Dodds, Tom Henighan, Scott Johnston, Andy Jones, Nicholas Joseph, Jackson Kernian,
262 Shauna Kravec, Ben Mann, Neel Nanda, Kamal Ndousse, Catherine Olsson, Daniela Amodei,
263 Tom Brown, Jared Kaplan, Sam McCandlish, Christopher Olah, Dario Amodei, and Jack Clark.
264 Predictability and surprise in large generative models. In *2022 ACM Conference on Fairness,*
265 *Accountability, and Transparency, FAccT ’22*, pp. 1747–1764, New York, NY, USA, 2022. Associ-
266 ation for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533229. URL
267 <https://doi.org/10.1145/3531146.3533229>.
- 268 Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. Why do tree-based models still outperform
269 deep learning on typical tabular data? In *Thirty-sixth Conference on Neural Information Processing*
270 *Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Fp7__phQszn)
271 [Fp7__phQszn](https://openreview.net/forum?id=Fp7__phQszn).
- 272 Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David
273 Sontag. Tablml: Few-shot classification of tabular data with large language models. In *International*
274 *Conference on Artificial Intelligence and Statistics*, pp. 5549–5581. PMLR, 2023.
- 275 Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. How we analyzed the
276 compas recidivism algorithm, 2016. URL [https://www.propublica.org/article/](https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm)
277 [how-we-analyzed-the-compas-recidivism-algorithm](https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm).
- 278 Yanchen Liu, William Held, and Diyi Yang. Dada: Dialect adaptation via dynamic aggregation of
279 linguistic rules, 2023.
- 280 Li Lucy and David Bamman. Gender and representation bias in GPT-3 generated stories. In
281 *Proceedings of the Third Workshop on Narrative Understanding*, pp. 48–55, Virtual, June 2021.
282 Association for Computational Linguistics. doi: 10.18653/v1/2021.nuse-1.5. URL [https://](https://aclanthology.org/2021.nuse-1.5)
283 aclanthology.org/2021.nuse-1.5.
- 284 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
285 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In
286 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.
287 11048–11064, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
288 Linguistics. doi: 10.18653/v1/2022.emnlp-main.759. URL [https://aclanthology.org/2022.](https://aclanthology.org/2022.emnlp-main.759)
289 [emnlp-main.759](https://aclanthology.org/2022.emnlp-main.759).
- 290 OpenAI. Gpt-4 technical report, 2023.
- 291 Ravid Shwartz-Ziv and Amitai Armon. Tabular data: Deep learning is not all you need. In *8th ICML*
292 *Workshop on Automated Machine Learning (AutoML)*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=vdgtepS1pV)
293 [forum?id=vdgtepS1pV](https://openreview.net/forum?id=vdgtepS1pV).
- 294 Dylan Slack and Sameer Singh. Tablet: Learning from instructions for tabular data. *arXiv*, 2023.
- 295 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
296 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
297 *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- 298 Jerry W. Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,
299 Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differ-
300 ently. *ArXiv*, abs/2303.03846, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:257378479)
301 [257378479](https://api.semanticscholar.org/CorpusID:257378479).

- 302 Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra
303 Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins,
304 Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks,
305 William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of
306 harm from language models, 2021.
- 307 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
308 learning as implicit bayesian inference. In *International Conference on Learning Representations*,
309 2022. URL <https://openreview.net/forum?id=RdJVFCHjUMI>.
- 310 Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Sebastian Riedel. TaBERT: Pretraining for
311 joint understanding of textual and tabular data. In *Proceedings of the 58th Annual Meeting of the*
312 *Association for Computational Linguistics*, pp. 8413–8426, Online, July 2020. Association for
313 Computational Linguistics. doi: 10.18653/v1/2020.acl-main.745. URL <https://aclanthology.org/2020.acl-main.745>.
- 315 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Men also like
316 shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of*
317 *the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2979–2989,
318 Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.
319 18653/v1/D17-1323. URL <https://aclanthology.org/D17-1323>.
- 320 Caleb Ziems, William Held, Jingfeng Yang, Jwala Dhamala, Rahul Gupta, and Diyi Yang. Multi-
321 VALUE: A framework for cross-dialectal English NLP. In *Proceedings of the 61st Annual Meeting*
322 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 744–768, Toronto,
323 Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.44.
324 URL <https://aclanthology.org/2023.acl-long.44>.

325 **A Description for each Feature in each Dataset**

326 We provide a detailed description of each dataset evaluated in our paper.

327 **A.1 Adult**

328 The *Adult Income* dataset (Adult) is extracted from the 1994 U.S. Census Bureau database. The task
 329 is to predict whether a person earns more than \$50,000 per year based on their profile data (*greater
 330 than 50K or less than or equal to 50K*). The original Adult Income Dataset contains 14 features, as
 331 described in Table 2. Following previous work (Slack & Singh, 2023), we retain only 10 features:
 332 “*workclass*”, “*hours per week*”, “*sex*”, “*age*”, “*occupation*”, “*capital loss*”, “*education*”, “*capital
 333 gain*”, “*marital status*”, and “*relationship*”. Our analysis on Adult primarily focuses on *sex* as the
 334 protected attribute, and *female* is acknowledged as a disadvantaged group.

Feature	Type	Description
Age	Continuous	Represents the age of an individual.
Workclass	Categorical	Indicates the type of employment, such as private, self-employed, or government.
<i>Fnlwgt</i>	Continuous	Stands for “final weight” and is a numerical value used in sampling for survey data.
Education	Categorical	Specifies the highest level of education attained by the individual, such as high school, bachelor’s degree, etc.
<i>Education-Num</i>	Continuous	Represents the numerical equivalent of the education level.
Marital-Status	Categorical	Describes the marital status of the individual, including categories like married, divorced, or single.
Occupation	Categorical	Indicates the occupation of the individual, such as managerial, technical, or clerical work.
Relationship	Categorical	Specifies the individual’s role in the family, such as husband, wife, or child.
Race	Categorical	Represents the individual’s race or ethnic background.
Sex	Categorical	Indicates the gender of the individual, either male or female.
Capital-Gain	Continuous	Refers to the capital gains, which are profits from the sale of assets, of the individual.
Capital-Loss	Continuous	Represents the capital losses, which are losses from the sale of assets, of the individual.
Hours-Per-Week	Continuous	Denotes the number of hours worked per week by the individual.
<i>Native-Country</i>	Categorical	Specifies the native country or place of origin of the individual.
Income (target)	Binary	The target variable indicating whether an individual’s income exceeds a certain threshold, typically \$50,000 per year.

Table 2: Features in the original **Adult** dataset. Those not used in our work are shown in *italics*.

335 **A.2 COMPAS**

336 The COMPAS dataset comprises the outcomes from the *Correctional Offender Management Profiling
 337 for Alternative Sanctions* commercial algorithm, utilized to evaluate a convicted criminal’s probability
 338 of reoffending. Known for its widespread use by judges and parole officers, COMPAS has gained
 339 notoriety for its bias against African-Americans. The raw COMPAS Recidivism dataset contains
 340 more than 50 attributes. Following the approach of Larson et al. (2016), we perform necessary
 341 preprocessing, group “*race*” into *African-American* and *Not African-American*, and only consider
 342 the features “*sex*”, “*race*”, “*age*”, “*charge degree*”, “*priors count*”, “*risk*” and “*two year recid*”
 343 (target). We frame the task as predicting whether an individual will recidivate in two years (*Did Not
 344 Reoffend* or *Reoffended*) based on their demographic and criminal history. For the COMPAS dataset,

345 we consider *race* as the protected attribute. Due to page limitations, we provide descriptions for only
 346 the features used in our work in Table 3.

Feature	Type	Description
Sex	Categorical	The gender of the individual.
Race	Categorical	The race of the individual, grouped into <i>African-American</i> and <i>Not African-American</i> .
Age	Continuous	The age of the individual.
Charge Degree	Categorical	The degree of the charge against the individual.
Priors Count	Continuous	The number of prior convictions or charges.
Risk	Categorical	The risk assessment for recidivism.
Two-Year Recid (target)	Binary	The target variable indicating whether an individual recidivated within two years.

Table 3: Features in the **COMPAS** Recidivism Dataset (Preprocessed).

347 B Evaluation Metrics

348 Here, we briefly explain each evaluation metric for the fairness we consider in our work.

349 **Accuracy and F1** As the most basic metric, assessing accuracy among different subgroups ensures
 350 that the model delivers consistent performance across all groups, without undue favor to any particular
 351 subgroups. Considering that the evaluated datasets may be imbalanced, especially among different
 352 subgroups, the F1 Score computes the harmonic mean of precision and recall, offering a balanced
 353 perspective between these two metrics.

354 **Statistical Parity** Statistical parity is attained when *positive* decision outcomes (e.g., being pre-
 355 dicted as a good credit risk) are independent of the protected attributes. This metric assesses whether
 356 different subgroups receive similar treatment from the model. For each subgroup z_i of each protected
 357 attribute Z , we calculate

$$P(\hat{Y} = 1 | Z = z_i).$$

358 Then we calculate the Statistical Parity Difference (SPD) of this protected attribute as

$$SPD = P(\hat{Y} = 1 | Z = z_1) - P(\hat{Y} = 1 | Z = z_2),$$

359 where z_1 is the minority group and z_2 is the majority.

360 **Equality of Opportunity** Equality of opportunity requires that qualified individuals have an equal
 361 chance of being correctly classified by the model, regardless of their membership in a protected
 362 group. This metric ensures equal *true positive* rates between different subgroups, providing equal
 363 opportunities for each subgroup. Similar as statistical parity, for equality of opportunity, we calculate
 364 the Equal Opportunity Difference (EOD) as

$$EOD = P(\hat{Y} = 1 | Y = 1, Z = z_1) - P(\hat{Y} = 1 | Y = 1, Z = z_2).$$

365 Each of these metrics offers a different perspective on fairness. For each subgroup from each protected
 366 attribute, we will compute every aforementioned metric. A model demonstrating good fairness should
 367 show minimal gaps in these fairness metrics between different subgroups. Considering them together
 368 can provide a more comprehensive evaluation of the model’s fairness across different subgroups,
 369 ensuring that individuals are not unfairly disadvantaged based on their membership in a protected
 370 group.

371 C RF and NN hyperparameters

372 For RF, we fix number of trees to 100 for all datasets as well as models. For NN, we use a 3
 373 hidden-layered network with hyperparameters described in Table 4.

	h1	h2	h3	lr	batch size	epochs
Adult	16	64	16	0.07	128	300
German Credit	64	64	32	0.07	128	300
COMPAS	64	128	64	0.09	128	300

Table 4: Hyperparameters for all datasets for a 3 layer neural network, where h1, h2 and h3 represent the number of neurons in first, second and third hidden layers respectively, lr represents the learning rate, and is followed by the batch size and number of epochs the models are trained for.

374 **D Prompt Templates for each Dataset**

375 In this section, we provide the prompt templates we used in our work. The example below is from
 376 Adult dataset, where text in **blue** represents the task description, text in **green** denotes optional
 377 few-shot examples (only used in in-context learning), and text in **red** indicates the test example.

```

You must predict if income exceeds $50K/yr. Answer with one of the following:
greater than 50K | less than or equal to 50K.
Example 1 -
workclass: Private
hours per week: 20
sex: Male
age: 17
occupation: Other-service
capital loss: 0
education: 10th
capital gain: 0
marital status: Never-married
relationship: Own-child
Answer: less than or equal to 50K
...

workclass: Private
hours per week: 40
sex: Female
age: 24
occupation: Sales
capital loss: 0
education: Some-college
capital gain: 0
marital status: Never-married
relationship: Own-child
Answer:

```

Figure 1: Prompt Template for **Adult** Dataset.

378 Beyond the Adult dataset, we provide the serialization and prompt template utilized in our work for
 379 the COMPAS dataset, as following.

380 **D.1 COMPAS**

381 **E COMPAS Results**

382 **References**

Predict whether an individual will recidivate within two years based on demographic and criminal history. Answer with one of the following: Did Not Reoffend | Reoffended.

Example 1 -

sex: Male

race: African-American

age cat: 25 - 45

c charge degree: F

priors count: 0

risk: Low

Answer: Did Not Reoffend

sex: Male

race: African-American

age cat: 25 - 45

c charge degree: M

priors count: 13

risk: High

Answer:

Figure 2: Prompt Template for COMPAS Dataset.

				ACC	F1	SP	EoO				
GPT-3.5-turbo	Zero-Shot	AA	0.657	0.005	0.656	0.004	0.395	0.001	0.560	0.002	
		nAA	0.663	0.002	0.588	0.003	0.817	0.002	0.893	0.001	
		d	-0.006	0.005	0.068	0.006	-0.423	0.003	-0.334	0.002	
	Few-shot	Regular	AA	0.633	0.002	0.626	0.002	0.362	0.003	0.495	0.004
			nAA	0.642	0.001	0.623	0.002	0.614	0.002	0.709	0.002
			d	-0.008	0.003	0.003	0.003	-0.252	0.003	-0.214	0.005
		Label-flipping	AA	0.482	0.004	0.482	0.004	0.499	0.003	0.481	0.004
			nAA	0.412	0.003	0.408	0.003	0.471	0.002	0.404	0.003
			d	0.070	0.005	0.074	0.005	0.028	0.005	0.077	0.007
	Finetuning	Regular	AA	0.611	0.016	0.610	0.016	0.464	0.031	0.576	0.034
			nAA	0.616	0.013	0.586	0.016	0.657	0.032	0.724	0.029
			d	-0.005	0.017	0.024	0.024	-0.193	0.030	-0.148	0.027
		Oversampling	AA	0.609	0.007	0.608	0.007	0.494	0.071	0.605	0.066
			nAA	0.625	0.020	0.583	0.024	0.706	0.037	0.771	0.036
			d	-0.016	0.016	0.025	0.018	-0.212	0.037	-0.166	0.046
		Undersampling	AA	0.591	0.010	0.591	0.012	0.513	0.053	0.605	0.047
			nAA	0.641	0.008	0.612	0.009	0.663	0.035	0.749	0.037
			d	-0.050	0.016	-0.021	0.022	-0.150	0.033	-0.144	0.039
RF	Regular	AA	0.662	0.004	0.662	0.004	0.496	0.006	0.660	0.007	
		nAA	0.671	0.004	0.617	0.002	0.767	0.008	0.859	0.009	
		d	-0.009	0.007	0.045	0.005	-0.271	0.011	-0.199	0.014	
	Oversampling	AA	0.660	0.005	0.660	0.005	0.493	0.010	0.655	0.013	
		nAA	0.671	0.002	0.624	0.002	0.743	0.003	0.839	0.004	
		d	-0.010	0.006	0.037	0.006	-0.250	0.012	-0.184	0.016	
	Undersampling	AA	0.648	0.002	0.647	0.002	0.491	0.004	0.639	0.004	
		nAA	0.667	0.005	0.614	0.007	0.761	0.006	0.851	0.006	
		d	-0.020	0.007	0.033	0.008	-0.270	0.009	-0.211	0.008	
NN	Regular	AA	0.666	0.003	0.665	0.002	0.462	0.034	0.630	0.034	
		nAA	0.662	0.003	0.613	0.006	0.742	0.019	0.831	0.017	
		d	0.005	0.006	0.052	0.007	-0.280	0.019	-0.201	0.018	
	Oversampling	AA	0.656	0.001	0.653	0.012	0.507	0.090	0.665	0.101	
		nAA	0.643	0.013	0.580	0.034	0.757	0.107	0.828	0.091	
		d	0.013	0.014	0.073	0.043	-0.249	0.049	-0.163	0.046	
	Undersampling	AA	0.660	0.019	0.657	0.023	0.477	0.078	0.638	0.097	
		nAA	0.657	0.013	0.602	0.026	0.757	0.051	0.839	0.040	
		d	0.003	0.024	0.055	0.043	-0.280	0.041	-0.202	0.064	

Table 5: **Fairness evaluation for COMPAS dataset** for the subgroup - *African American* (AA), and *Non African American* (nAA) as well as the difference (*d*). The significant fairness disparities are highlighted in red. Both in-context learning and finetuning can lead to bias reduction (indicated by ↓), and label-flipped in-context learning can further minimize bias (indicated by ✓).