

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

Anonymous ACL submission

## Abstract

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

## 1 Introduction

Many recent works propose to use large language models (LLMs) for tabular prediction (Slack and Singh, 2023; Hagselmann et al., 2023), where the tabular data is serialized as natural language and provided to LLMs with a short description of the task to solicit predictions. Despite the comprehensive examination of fairness considerations within conventional machine learning approaches applied

to tabular tasks (Bellamy et al., 2018), the exploration of fairness-related issues in the context of employing LLMs for tabular predictions remains a relatively underexplored domain.

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

Considering that tabular data finds extensive use in high-stakes domains (Grinsztajn et al., 2022) where information is typically structured in tabular formats as a natural byproduct of relational databases (Borisov et al., 2022), it is of paramount importance to thoroughly examine the fairness implications of utilizing LLMs for predictions on tabular data. In this paper, we conduct a series of investigations centered around this critical aspect, with the goal of discerning the underlying information sources upon which LLMs rely when making tabular predictions. Through this exploration, our investigation aims to ascertain whether, and to what degree, LLMs are susceptible to being influenced by social biases and stereotypes in the context of tabular data predictions.

Through experiments using GPT-3.5 to make predictions for tabular data in a zero-shot setting, we demonstrate that LLMs exhibit significant social biases (Section 4). This evidence confirms that LLMs inherit social biases from their training corpus and tend to rely on these biases when making predictions for tabular data.

084 Furthermore, we demonstrate that providing  
085 LLMs with few-shot examples (in-context learning)  
086 or fine-tuning them on the entire training dataset  
087 both exhibit moderate effects on bias mitigation  
088 (Sections 5.1 and 6.1). Nevertheless, the achieved  
089 fairness levels remain below what is typically at-  
090 tained with traditional machine learning methods,  
091 including Random Forests and shallow Neural Net-  
092 works, once again underscoring the presence of  
093 inherent bias in LLMs. Additionally, our investi-  
094 gation further reveals that flipping the labels of the  
095 in-context examples significantly narrows the gap  
096 in fairness metrics across different subgroups, but  
097 comes at the expected cost of a reduction in pre-  
098 dictive performance. This finding, in turn, further  
099 emphasizes and reaffirms the indication of inherent  
100 bias present in LLMs (Section 5.2). Additionally,  
101 we further show that while resampling the training  
102 set is a known and effective method for reducing  
103 biases in traditional machine learning methods like  
104 Random Forests and shallow Neural Networks, it  
105 proves to be less effective when applied to LLMs  
106 (Section 6.2).

107 These collective findings underscore the signif-  
108 icant influence of social biases on LLMs’ perfor-  
109 mance in tabular predictions. These biases sig-  
110 nificantly undermine fairness and pose substantial  
111 potential risks for using LLMs on tabular data, es-  
112 pecially considering that tabular data is extensively  
113 used in high-stakes domains, highlighting the need  
114 for more advanced and tailored strategies to address  
115 these biases effectively. Straightforward methods  
116 like in-context learning and data resampling may  
117 not be sufficient in this context.

## 118 2 Related work

119 **Fairness and Social Biases in LLMs** Fairness  
120 is highly desirable for ensuring the credibility and  
121 trustworthiness of algorithms. It has been demon-  
122 strated that unfair algorithms can reflect societal  
123 biases in their decision-making processes (Ben-  
124 der et al., 2021; Bommasani, 2021), primarily  
125 stemming from the biases present in their train-  
126 ing data (Caliskan et al., 2017; Zhao et al.,  
127 2017). LLMs, pre-trained on vast natural language  
128 datasets, are particularly susceptible to inheriting  
129 these social biases and have been shown to exhibit  
130 biases related to gender (Lucy and Bamman, 2021),  
131 religion (Abid et al., 2021b) and language vari-  
132 ants (Ziems et al., 2023; Liu et al., 2023a). These  
133 social biases can lead to the perpetuation of discrim-

134 ination and stereotype (Abid et al., 2021a; Bender  
135 et al., 2021; Weidinger et al., 2021). While re-  
136 cent literature has made strides in addressing these  
137 issues, there still exists a significant gap in com-  
138 prehensively assessing fairness in LLMs and its  
139 mitigation strategies for tabular data.

### 140 **Tabular Tasks and LLM for Tabular Data**

141 Tabular data extensively exist in many domains  
142 (Shwartz-Ziv and Armon, 2021). Previous works  
143 propose to utilize self-supervised deep techniques  
144 for tabular tasks (Yin et al., 2020; Arik and Pfister,  
145 2021), which, however, still underperform ensem-  
146 bles of gradient-boosted trees in the fully super-  
147 vised setting (Grinsztajn et al., 2022). This dispar-  
148 ity in performance can be attributed to the locality,  
149 sparsity, and mixed data types of tabular data. In  
150 recent times, LLMs have undergone intensive train-  
151 ing using vast amounts of natural language data,  
152 which has enabled them to exhibit impressive per-  
153 formance across various downstream tasks (Brown  
154 et al., 2020; OpenAI, 2023), even with little or no  
155 labeled task data. Therefore, recent approaches  
156 by (Hegselmann et al., 2023; Slack and Singh,  
157 2023) suggest serializing the tabular data as natural  
158 language, which is provided to LLM along with a  
159 short task description to generate predictions for  
160 tabular tasks. However, tabular data plays a crucial  
161 role in numerous safety-critical and high-stakes do-  
162 mains (Borisov et al., 2022; Grinsztajn et al., 2022),  
163 which makes fairness particularly crucial when em-  
164 ploying LLMs for making predictions on tabular  
165 data, especially considering the inherent social bi-  
166 ases present in LLMs. Despite the importance, this  
167 still remains largely unexplored. To the best of our  
168 knowledge, we regard our work as one of the most  
169 comprehensive investigations into the fairness is-  
170 sues arising when using LLMs for predictions on  
171 tabular data.

172 **In-Context Learning** Significant improvements  
173 for various tasks have been achieved by providing  
174 in-context examples to LLMs (Brown et al., 2020;  
175 Liu et al., 2022, 2023b). However, previous re-  
176 search by (Min et al., 2022; Wei et al., 2023b; Lyu  
177 et al., 2023) illustrate that the effective performance  
178 of in-context learning largely hinges on semantic  
179 priors rather than learning the input-label mapping  
180 (Akyürek et al., 2022; Xie et al., 2022; Von Os-  
181 wald et al., 2023) and the labels of the in-context  
182 examples might not play a crucial role in in-context  
183 learning, with flipped or random labels sometimes

having minimal impact on performance. Despite these findings, the predominant focus of existing investigation of in-context learning remains on conventional natural language processing tasks (Zhao et al., 2021; Min et al., 2022; Wei et al., 2023a,b), largely overlooking the domain of tabular data. Furthermore, the fairness of in-context learning and the impact of flipped labels on this fairness is yet to be thoroughly investigated.

### 3 Experimental Setup

In this section, we outline the general setup of the experiments conducted in our work.

#### 3.1 Models

In our work, we focus our experiments on GPT-3.5 (engine GPT-3.5-turbo) - an LLM released by OpenAI, trained with instruction tuning (Sanh et al., 2022; Wei et al., 2022) and reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), aligning LLMs with human preferences. Furthermore, we also compare its performance with conventional machine learning models in order to gain insight into the propagation of biases found within LLMs, which are likely mirrored in traditional models as well, consequently, offering valuable additional perspectives on the biases inherent in the training of LLMs. For this, we employ two widely used models for tabular data i.e., Random Forests (RF) and a shallow Neural Network (NN) of 3 layers. We provide additional implementation details in Appendix B.

#### 3.2 Datasets and Protected Attributes

To explore the fairness of LLMs in making predictions for tabular data, we utilize the following three widely used tabular datasets for assessing the fairness of traditional ML models: *Adult Income* (**Adult**) Dataset (Becker and Kohavi, 1996), **German Credit** Dataset (Dua and Graff, 2019), and *Correctional Offender Management Profiling for Alternative Sanctions* (**COMPAS**) Dataset (Larson et al., 2016). In this section, we introduce each dataset and discuss its associated protected attributes.

**Adult** The *Adult Income* dataset (Adult) is extracted from the 1994 U.S. Census Bureau database. The task is to predict whether a person earns more than \$50,000 per year based on their profile data (*greater than 50K* or *less than or equal to 50K*).

The original Adult Income Dataset contains 14 features. Following previous work (Slack and Singh, 2023), we retain only 10 features: “*workclass*”, “*hours per week*”, “*sex*”, “*age*”, “*occupation*”, “*capital loss*”, “*education*”, “*capital gain*”, “*marital status*”, and “*relationship*”. Our analysis on Adult primarily focuses on *sex* as the protected attribute, and *female* is acknowledged as a disadvantaged group.

**German Credit** The German Credit dataset is used to classify individuals based on their profile attributes as good or bad credit risks (*good* or *bad*). The raw dataset comprises 20 attributes. Consistent with previous work, we only retain the following features: “*age*”, “*sex*”, “*job*”, “*housing*”, “*saving accounts*”, “*checking account*”, “*credit amount*”, “*duration*”, and “*purpose*”. Same with Adult, *sex* is considered as a protected attribute in the German Credit dataset and *female* as the marginalized group.

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

A detailed description for each feature of the considered datasets is provided in Appendix A.

#### 3.3 Serialization and Prompt Templates

To employ the LLM for making predictions on these tabular datasets, each data point is first serialized as text. Following previous works on LLM for tabular predictions (Hegselmann et al., 2023; Slack and Singh, 2023), we format the feature names and values into strings as “ $f_1 : x_1, \dots, f_d : x_d$ ”, and

prompt to LLM along with a task description, as illustrated following:

```

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.

The example above is from the Adult dataset, where text in blue denotes the task description, text in green indicates optional few-shot examples (only used in in-context learning), and text in red is the test example. We provide the prompt templates for the other two datasets in Appendix C.

### 3.4 Evaluation Metrics

To assess fairness in the aforementioned datasets, we examine the disparity between different subgroups of protected attributes using the following common fairness metrics: accuracy, F1 score, statistical parity, and equality of opportunity. Here, we briefly explain each evaluation metric.

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

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

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

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

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

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

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

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

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

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

## 4 Zero-Shot Prompting for Tabular Data

To explore the fairness of LLMs when making predictions on tabular data, we first conduct experiments in a zero-shot setting. We assess the fairness metrics of the outcomes and examine whether LLMs without any finetuning or few-shot examples

would be influenced by social biases and stereotypes for tabular predictions. We run all the experiments 5 times and compute the mean and standard deviation.

In Tables 1-3, we present the evaluation of four fairness metrics, namely accuracy (ACC), F1 score (F1), statistical parity (SP), and equality of opportunity (EoO), for GPT-3.5 (engine GPT-3.5-turbo), RF and NN models on the **Adult**, **German Credit** and **COMPAS** datasets, respectively. For the Adult and German Credit datasets, the subgroups *female* and *male* are assessed regarding the protected attribute *sex*, identifying *female* as a disadvantaged group. In the COMPAS dataset, we evaluate *race* as protected attributes, recognizing African American (AA) as the disadvantaged group.

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

Exceptionally, GPT-3.5 is extremely biased for German Credit dataset where it classifies almost everything into ‘*good credit*’ class in the zero-shot setting, thus rendering the difference in SP and EoO for both subgroups to be near 0. The accuracy for each subgroup is near 50%, performing similar to random guessing. The possible reason might be that the German Credit dataset is too challenging for making tabular predictions with LLMs (especially, since the features of German Credit are ambiguous and vague). This also suggests that, when using LLM to make predictions on tabular data, a potential description of table feature names is favorable.

These findings demonstrate the tendency of LLMs to rely on social biases and stereotypes inherited from their training corpus when applied to tabular data. This implies that using LLMs for predictions on tabular data may incur significant

fairness risks, including the potential to disproportionately disadvantage marginalized communities as well as exacerbate social biases and stereotypes present in society. This is particularly concerning given the widespread application of tabular data in high-stake contexts, further magnifying the potential for harm.

## 5 Few-Shot Prompting for Tabular Data

As demonstrated in Section 4, employing LLMs for predictions on tabular data reveals significant social biases in a zero-shot setting. Instead of directly utilizing LLMs for zero-shot tabular predictions, this section explores whether including few-shot examples during prompting will reduce or amplify these biases. To delve deeper into the influence of few-shot examples during in-context learning (ICL), we not only consider the regular ICL approach as detailed in Section 5.1, but we also experiment by flipping the labels of the few-shot examples to further examine their effect on the biases, as discussed in Section 5.2. Again, for robustness, each experiment is conducted 5 times, with the mean and standard deviation reported.

### 5.1 Regular In-Context Learning

Previous works have demonstrated that LLMs can learn the input-label mappings in context (Akyürek et al., 2022; Xie et al., 2022; Von Oswald et al., 2023). However, the influence of in-context learning on fairness has not been thoroughly examined. For in-context learning, the test example and task description, along with a few-shot examples, are provided to the LLMs for generating the final predictions. The few-shot examples are inserted before the test example in the prompt, as outlined in Section 3.3. We set the number of in-context examples as 50. For each dataset, we randomly select the in-context examples from the training set for each test example.

In Tables 1-3, we demonstrate that for two of the evaluated datasets (except for COMPAS), the incorporation of few-shot examples brings about performance improvements. Additionally, we observe that incorporating few-shot examples into prompting reduces the fairness metric gap between different subgroups. However, a significant fairness issue still persists. Moreover, for the Adult and COMPAS datasets, the disparity in fairness metrics of in-context learning is more notable when compared to traditional models, such as RF and

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

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 ✓).

446 NN. This highlights the inherent biases embedded  
447 within LLMs, which are not solely derived from  
448 the task datasets.

## 449 5.2 Label Flipping

450 To delve deeper into the sources of biases within  
451 LLMs, we further examine the impact of the labels  
452 of in-context examples on fairness. As depicted  
453 in Tables 1-3, label flipping significantly reduces  
454 biases across all evaluated datasets. For all datasets,  
455 the difference in statistical parity (SP) and equal-  
456 ity of opportunity (EoO) is minimized with label-  
457 flipped ICL. For example, the absolute gap of EoO  
458 on the Adult dataset decreases from 0.483 in zero-  
459 shot prompting to 0.037, almost completely elimi-  
460 nating the bias. These findings further corroborate  
461 the existence of inherent biases in LLMs.

462 However, flipped labels lead to a significant drop

463 in predictive performance. Though previous re-  
464 search suggests that the effectiveness of ICL pre-  
465 dominantly stems from semantic priors, rather than  
466 learning the input-label mappings (Min et al., 2022;  
467 Wei et al., 2023b) and demonstrates that the perfor-  
468 mance of ICL is barely affected even with flipped  
469 or random labels for in-context examples, the focus  
470 of these works lies mainly on traditional natural  
471 language processing tasks.

472 In contrast, we observe that the labels of in-  
473 context examples hold substantial influence over  
474 predictive performance in our unique setup, where  
475 LLMs are deployed for predictions on tabular data.  
476 This could be attributed to the limited exposure of  
477 these models to tabular data during pre-training,  
478 thereby amplifying the role of input-label mapping  
479 of in-context examples.

|               |               |                | ACC                     | F1                      | SP                      | EoO                     |                        |
|---------------|---------------|----------------|-------------------------|-------------------------|-------------------------|-------------------------|------------------------|
| GPT-3.5-turbo | Zero-Shot     | <i>f</i>       | 0.471 <sub>0.011</sub>  | 0.359 <sub>0.021</sub>  | 0.980 <sub>0.011</sub>  | 1.000 <sub>0.000</sub>  |                        |
|               |               | <i>m</i>       | 0.556 <sub>0.000</sub>  | 0.357 <sub>0.000</sub>  | 0.984 <sub>0.000</sub>  | 0.972 <sub>0.000</sub>  |                        |
|               |               | <i>d</i>       | -0.084 <sub>0.011</sub> | 0.002 <sub>0.021</sub>  | -0.004 <sub>0.011</sub> | 0.028 <sub>0.000</sub>  |                        |
|               | Few-shot      | Regular        | <i>f</i>                | 0.610 <sub>0.013</sub>  | 0.593 <sub>0.013</sub>  | 0.348 <sub>0.027</sub>  | 0.453 <sub>0.029</sub> |
|               |               |                | <i>m</i>                | 0.606 <sub>0.007</sub>  | 0.603 <sub>0.008</sub>  | 0.337 <sub>0.007</sub>  | 0.450 <sub>0.012</sub> |
|               |               |                | <i>d</i>                | 0.003 <sub>0.012</sub>  | -0.010 <sub>0.011</sub> | 0.011 <sub>0.027</sub>  | 0.003 <sub>0.026</sub> |
|               |               | Label-flipping | <i>f</i>                | 0.614 <sub>0.011</sub>  | 0.606 <sub>0.012</sub>  | 0.695 <sub>0.011</sub>  | 0.842 <sub>0.000</sub> |
|               |               |                | <i>m</i>                | 0.559 <sub>0.013</sub>  | 0.538 <sub>0.011</sub>  | 0.638 <sub>0.013</sub>  | 0.672 <sub>0.023</sub> |
|               |               |                | <i>d</i>                | 0.056 <sub>0.021</sub>  | 0.067 <sub>0.021</sub>  | 0.057 <sub>0.012</sub>  | 0.170 <sub>0.023</sub> |
|               | Finetuning    | Regular        | <i>f</i>                | 0.571 <sub>0.067</sub>  | 0.567 <sub>0.062</sub>  | 0.619 <sub>0.101</sub>  | 0.711 <sub>0.186</sub> |
|               |               |                | <i>m</i>                | 0.548 <sub>0.011</sub>  | 0.539 <sub>0.023</sub>  | 0.532 <sub>0.123</sub>  | 0.569 <sub>0.098</sub> |
|               |               |                | <i>d</i>                | 0.024 <sub>0.079</sub>  | 0.029 <sub>0.085</sub>  | 0.087 <sub>0.022</sub>  | 0.141 <sub>0.088</sub> |
|               |               | Oversampling   | <i>f</i>                | 0.536 <sub>0.017</sub>  | 0.532 <sub>0.012</sub>  | 0.607 <sub>0.084</sub>  | 0.658 <sub>0.112</sub> |
|               |               |                | <i>m</i>                | 0.532 <sub>0.011</sub>  | 0.523 <sub>0.020</sub>  | 0.548 <sub>0.079</sub>  | 0.569 <sub>0.059</sub> |
|               |               |                | <i>d</i>                | 0.004 <sub>0.028</sub>  | 0.009 <sub>0.033</sub>  | 0.060 <sub>0.006</sub>  | 0.088 <sub>0.053</sub> |
|               |               | Undersampling  | <i>f</i>                | 0.548 <sub>0.034</sub>  | 0.547 <sub>0.033</sub>  | 0.571 <sub>0.034</sub>  | 0.632 <sub>0.074</sub> |
|               |               |                | <i>m</i>                | 0.556 <sub>0.000</sub>  | 0.555 <sub>0.000</sub>  | 0.444 <sub>0.000</sub>  | 0.500 <sub>0.000</sub> |
|               |               |                | <i>d</i>                | -0.008 <sub>0.034</sub> | -0.008 <sub>0.033</sub> | 0.127 <sub>0.034</sub>  | 0.132 <sub>0.074</sub> |
| RF            | Regular       | <i>f</i>       | 0.581 <sub>0.024</sub>  | 0.580 <sub>0.025</sub>  | 0.519 <sub>0.028</sub>  | 0.611 <sub>0.054</sub>  |                        |
|               |               | <i>m</i>       | 0.600 <sub>0.019</sub>  | 0.588 <sub>0.020</sub>  | 0.597 <sub>0.022</sub>  | 0.672 <sub>0.021</sub>  |                        |
|               |               | <i>d</i>       | -0.019 <sub>0.016</sub> | -0.008 <sub>0.016</sub> | -0.078 <sub>0.044</sub> | -0.062 <sub>0.061</sub> |                        |
|               | Oversampling  | <i>f</i>       | 0.576 <sub>0.018</sub>  | 0.575 <sub>0.018</sub>  | 0.505 <sub>0.018</sub>  | 0.589 <sub>0.021</sub>  |                        |
|               |               | <i>m</i>       | 0.568 <sub>0.032</sub>  | 0.552 <sub>0.034</sub>  | 0.616 <sub>0.025</sub>  | 0.661 <sub>0.037</sub>  |                        |
|               |               | <i>d</i>       | 0.008 <sub>0.034</sub>  | 0.023 <sub>0.035</sub>  | -0.111 <sub>0.013</sub> | -0.072 <sub>0.041</sub> |                        |
|               | Undersampling | <i>f</i>       | 0.586 <sub>0.024</sub>  | 0.585 <sub>0.024</sub>  | 0.533 <sub>0.024</sub>  | 0.632 <sub>0.047</sub>  |                        |
|               |               | <i>m</i>       | 0.575 <sub>0.031</sub>  | 0.555 <sub>0.037</sub>  | 0.635 <sub>0.033</sub>  | 0.683 <sub>0.022</sub>  |                        |
|               |               | <i>d</i>       | 0.011 <sub>0.024</sub>  | 0.031 <sub>0.031</sub>  | -0.102 <sub>0.041</sub> | -0.052 <sub>0.039</sub> |                        |
| NN            | Regular       | <i>f</i>       | 0.533 <sub>0.024</sub>  | 0.533 <sub>0.024</sub>  | 0.519 <sub>0.028</sub>  | 0.558 <sub>0.026</sub>  |                        |
|               |               | <i>m</i>       | 0.556 <sub>0.017</sub>  | 0.544 <sub>0.017</sub>  | 0.584 <sub>0.012</sub>  | 0.622 <sub>0.022</sub>  |                        |
|               |               | <i>d</i>       | -0.022 <sub>0.037</sub> | -0.012 <sub>0.036</sub> | -0.065 <sub>0.031</sub> | -0.064 <sub>0.026</sub> |                        |
|               | Oversampling  | <i>f</i>       | 0.548 <sub>0.040</sub>  | 0.547 <sub>0.040</sub>  | 0.552 <sub>0.028</sub>  | 0.611 <sub>0.026</sub>  |                        |
|               |               | <i>m</i>       | 0.562 <sub>0.026</sub>  | 0.547 <sub>0.024</sub>  | 0.603 <sub>0.048</sub>  | 0.644 <sub>0.057</sub>  |                        |
|               |               | <i>d</i>       | -0.014 <sub>0.037</sub> | 0.000 <sub>0.035</sub>  | -0.051 <sub>0.061</sub> | -0.034 <sub>0.065</sub> |                        |
|               | Undersampling | <i>f</i>       | 0.529 <sub>0.049</sub>  | 0.524 <sub>0.047</sub>  | 0.467 <sub>0.051</sub>  | 0.495 <sub>0.042</sub>  |                        |
|               |               | <i>m</i>       | 0.495 <sub>0.025</sub>  | 0.490 <sub>0.023</sub>  | 0.524 <sub>0.047</sub>  | 0.517 <sub>0.054</sub>  |                        |
|               |               | <i>d</i>       | 0.033 <sub>0.063</sub>  | 0.035 <sub>0.059</sub>  | -0.057 <sub>0.033</sub> | -0.022 <sub>0.061</sub> |                        |

Table 2: **Fairness evaluation for German Credit 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.

## 6 Finetuning for Tabular Data

### 6.1 Regular Finetuning

Finally, we extend our investigation to assess if finetuning the models on the entire training set could aid in diminishing the social biases in LLMs. For GPT-3.5, fine-tuning is executed using the publicly released API from OpenAI. For RF and NN, we provide the training details in Appendix B. We still run all the experiments 5 times and compute the mean and standard deviation. In Tables 1-3, we show that finetuning effectively reduces unfairness in all datasets, making them comparable and sometimes significantly better in terms of SP and EoO when compared to RF and NN. For example, the absolute difference in EoO after finetuning on the Adult dataset is 0.0714, which is lower than the 0.123 difference of an NN.

### 6.2 Resampling

We further explore the potential of resampling, a method frequently employed to enhance fairness in machine learning model training, particularly in scenarios where there is a significant class imbalance or bias in the data. To this end, we evaluate two approaches: oversampling the minority group and undersampling the majority group. As depicted in Tables 1-3, resampling fails to mitigate the social biases in LLMs when making tabular predictions, even though we demonstrate that oversampling generally reduces social biases for both RF and NN, except for a few instances such as oversampling in NN for adult dataset worsens the fairness.

Our finetuning experiments show that the social biases inherited from LLM’s pre-training data which are evident when making predictions on tabular data, can sometimes be mitigated through finetuning. Nevertheless, unlike the consistent out-

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

Table 3: **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 ✓).

comes typically seen in traditional machine learning models, data resampling does not consistently produce similar results for finetuning LLMs.

## 7 Conclusion

In this work, we thoroughly investigate the under-explored problem of fairness of large language models (LLMs) for tabular tasks. Our study unfolds in several phases. Initially, we assess the inherent fairness displayed by LLMs, comparing their performance in zero-shot learning scenarios against traditional machine learning models like random forests (RF) and shallow neural networks (NN). Furthermore, we investigate how LLMs learn and propagate social biases when subjected to few-shot in-context learning, label-flipped in-context learning, fine-tuning, and data resampling techniques. Our discoveries shed light on several key insights. We find that LLMs tend to heavily rely on the so-

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



## 553 Limitations

554 It is important to note that our study exclusively  
555 focuses on the GPT-3.5 model. Consequently, our  
556 conclusions are representative of GPT-3.5 alone  
557 and cannot be extrapolated to other LLMs, which  
558 might exhibit different behaviors or biases. This  
559 focus on a single model thus restricts the broader  
560 applicability of our findings.

561 Furthermore, for each experiment, we employed  
562 only one type of prompt. This approach limits  
563 the generalizability of our conclusions, as different  
564 prompts might yield varying results. The use of  
565 a singular prompt type does not capture the full  
566 spectrum of possible interactions and outcomes  
567 that might be observed with a diverse range of  
568 prompting strategies.

569 Looking ahead, we plan to broaden the scope of  
570 our research. This expansion will include experi-  
571 menting with additional models beyond GPT-3.5,  
572 thus offering a more comprehensive understand-  
573 ing of fairness for different LLMs. We also intend  
574 to explore a variety of prompting strategies, such  
575 as Chain of Thought (CoT) prompting, to assess  
576 how different methods may impact model bias and  
577 fairness. These future endeavors aim to provide  
578 a more nuanced and thorough exploration of the  
579 capabilities and limitations of LLMs in the context  
580 of fairness.

## 581 References

- 582 Abubakar Abid, Maheen Farooqi, and James Zou.  
583 2021a. [Large language models associate muslims with violence](#). *Nature Machine Intelligence*,  
584 3(6):461–463.  
585
- 586 Abubakar Abid, Maheen Farooqi, and James Zou.  
587 2021b. [Persistent anti-muslim bias in large language models](#). In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society, AIES '21*, page  
588 298–306, New York, NY, USA. Association for Com-  
589 puting Machinery.  
590
- 592 Ekin Akyürek, Dale Schuurmans, Jacob Andreas,  
593 Tengyu Ma, and Denny Zhou. 2022. What learning  
594 algorithm is in-context learning? investigations with  
595 linear models. *arXiv preprint arXiv:2211.15661*.
- 596 Sercan Ö Arik and Tomas Pfister. 2021. Tabnet: Atten-  
597 tive interpretable tabular learning. In *Proceedings of  
598 the AAAI conference on artificial intelligence*, vol-  
599 ume 35, pages 6679–6687.
- 600 Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain,  
601 Deep Ganguli, Tom Henighan, Andy Jones, Nicholas

- Joseph, Ben Mann, Nova DasSarma, Nelson El-  
hage, Zac Hatfield-Dodds, Danny Hernandez, Jack-  
son Kernion, Kamal Ndousse, Catherine Olsson,  
Dario Amodei, Tom Brown, Jack Clark, Sam Mc-  
Candlish, Chris Olah, and Jared Kaplan. 2021. [A  
general language assistant as a laboratory for align-  
ment](#). 602  
603  
604  
605  
606  
607  
608
- Christine Basta, Marta R. Costa-jussà, and Noe Casas.  
2019. [Evaluating the underlying gender bias in con-  
textualized word embeddings](#). In *Proceedings of the  
First Workshop on Gender Bias in Natural Language  
Processing*. 609  
610  
611  
612  
613
- Barry Becker and Ronny Kohavi. 1996. Adult.  
UCI Machine Learning Repository. DOI:  
<https://doi.org/10.24432/C5XW20>. 614  
615  
616
- Rachel K. E. Bellamy, Kuntal Dey, Michael Hind,  
Samuel C. Hoffman, Stephanie Houde, Kalapriya  
Kannan, Pranay Lohia, Jacquelyn Martino, Sameep  
Mehta, Aleksandra Mojsilovic, Seema Nagar,  
Karthikeyan Natesan Ramamurthy, John Richards,  
Diptikalyan Saha, Prasanna Sattigeri, Moninder  
Singh, Kush R. Varshney, and Yunfeng Zhang. 2018.  
[AI Fairness 360: An extensible toolkit for detecting,  
understanding, and mitigating unwanted algorithmic  
bias](#). 617  
618  
619  
620  
621  
622  
623  
624  
625  
626
- Emily M. Bender, Timnit Gebru, Angelina McMillan-  
Major, and Shmargaret Shmitchell. 2021. [On the  
dangers of stochastic parrots: Can language models  
be too big?](#) In *Proceedings of the 2021 ACM Confer-  
ence on Fairness, Accountability, and Transparency*. 627  
628  
629  
630  
631
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou,  
Venkatesh Saligrama, and Adam T Kalai. 2016. [Man  
is to computer programmer as woman is to home-  
maker? debiasing word embeddings](#). In *Advances in  
Neural Information Processing Systems*, volume 29.  
Curran Associates, Inc. 632  
633  
634  
635  
636  
637
- R Bommasani. 2021. Opportunities and risks of foun-  
dation models. [https://openai.com/reports/  
foundation-models/](https://openai.com/reports/foundation-models/). 638  
639  
640
- Vadim Borisov, Tobias Leemann, Kathrin SeBler, Jo-  
hannes Haug, Martin Pawelczyk, and Gjergji Kas-  
neci. 2022. [Deep neural networks and tabular data:  
A survey](#). *IEEE Transactions on Neural Networks  
and Learning Systems*, pages 1–21. 641  
642  
643  
644  
645
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020.  
[Language models are few-shot learners](#). In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. Curran Associates,  
Inc. 646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659

|     |  |     |
|-----|--|-----|
| 660 | Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. <a href="#">Semantics derived automatically from language corpora contain human-like biases</a> . <i>Science</i> , 356(6334):183–186.  |     |
| 661 |  |     |
| 662 |  |     |
| 663 |  |     |
| 664 | D. Dua and C. Graff. 2019. <a href="#">UCI machine learning repository</a> .   |     |
| 665 |  |     |
| 666 | Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Scott Johnston, Andy Jones, Nicholas Joseph, Jackson Kernian, Shauna Kravec, Ben Mann, Neel Nanda, Kamal Ndousse, Catherine Olsson, Daniela Amodei, Tom Brown, Jared Kaplan, Sam McCandlish, Christopher Olah, Dario Amodei, and Jack Clark. 2022. <a href="#">Predictability and surprise in large generative models</a> . In <i>2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22</i> , page 1747–1764, New York, NY, USA. Association for Computing Machinery.   |     |
| 667 |  |     |
| 668 |  |     |
| 669 |  |     |
| 670 |  |     |
| 671 |  |     |
| 672 |  |     |
| 673 |  |     |
| 674 |  |     |
| 675 |  |     |
| 676 |  |     |
| 677 |  |     |
| 678 |  |     |
| 679 |  |     |
| 680 | Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. 2022. <a href="#">Why do tree-based models still outperform deep learning on typical tabular data?</a> In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .   |     |
| 681 |  |     |
| 682 |  |     |
| 683 |  |     |
| 684 |  |     |
| 685 | Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David Sontag. 2023. <a href="#">TablIm: Few-shot classification of tabular data with large language models</a> . In <i>International Conference on Artificial Intelligence and Statistics</i> , pages 5549–5581. PMLR.  |     |
| 686 |  |     |
| 687 |  |     |
| 688 |  |     |
| 689 |  |     |
| 690 |  |     |
| 691 | Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. 2016. <a href="#">How we analyzed the compas recidivism algorithm</a> .   |     |
| 692 |  |     |
| 693 |  |     |
| 694 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What makes good in-context examples for GPT-3?</a> In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.   |     |
| 695 |  |     |
| 696 |  |     |
| 697 |  |     |
| 698 |  |     |
| 699 |  |     |
| 700 |  |     |
| 701 |  |     |
| 702 | Yanchen Liu, William Held, and Diyi Yang. 2023a. <a href="#">Dada: Dialect adaptation via dynamic aggregation of linguistic rules</a> .  |     |
| 703 |  |     |
| 704 |  |     |
| 705 | Yanchen Liu, Timo Schick, and Hinrich Schtze. 2023b. <a href="#">Semantic-oriented unlabeled priming for large-scale language models</a> . In <i>Proceedings of The Fourth Workshop on Simple and Efficient Natural Language Processing (SustaiNLP)</i> , pages 32–38, Toronto, Canada (Hybrid). Association for Computational Linguistics.  |     |
| 706 |  |     |
| 707 |  |     |
| 708 |  |     |
| 709 |  |     |
| 710 |  |     |
| 711 |  |     |
| 712 | Li Lucy and David Bamman. 2021. <a href="#">Gender and representation bias in GPT-3 generated stories</a> . In <i>Proceedings of the Third Workshop on Narrative Understanding</i> , pages 48–55, Virtual. Association for Computational Linguistics.  |     |
| 713 |  |     |
| 714 |  |     |
| 715 |  |     |
| 716 |  |     |
|     | Xinxi Lyu, Sewon Min, Iz Beltagy, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. <a href="#">Z-ICL: Zero-shot in-context learning with pseudo-demonstrations</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2304–2317, Toronto, Canada. Association for Computational Linguistics.  | 717 |
|     |  | 718 |
|     |  | 719 |
|     |  | 720 |
|     |  | 721 |
|     |  | 722 |
|     |  | 723 |
|     | Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Rethinking the role of demonstrations: What makes in-context learning work?</a> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.  | 724 |
|     |  | 725 |
|     |  | 726 |
|     |  | 727 |
|     |  | 728 |
|     |  | 729 |
|     |  | 730 |
|     |  | 731 |
|     | OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .   | 732 |
|     | Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . In <i>Advances in Neural Information Processing Systems</i> .  | 733 |
|     |  | 734 |
|     |  | 735 |
|     |  | 736 |
|     |  | 737 |
|     |  | 738 |
|     |  | 739 |
|     |  | 740 |
|     |  | 741 |
|     | Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. <a href="#">Multi-task prompted training enables zero-shot task generalization</a> . In <i>International Conference on Learning Representations</i> . | 742 |
|     |  | 743 |
|     |  | 744 |
|     |  | 745 |
|     |  | 746 |
|     |  | 747 |
|     |  | 748 |
|     |  | 749 |
|     |  | 750 |
|     |  | 751 |
|     |  | 752 |
|     |  | 753 |
|     |  | 754 |
|     |  | 755 |
|     |  | 756 |
|     |  | 757 |
|     | Ravid Shwartz-Ziv and Amitai Armon. 2021. <a href="#">Tabular data: Deep learning is not all you need</a> . In <i>8th ICML Workshop on Automated Machine Learning (AutoML)</i> .   | 758 |
|     |  | 759 |
|     |  | 760 |
|     |  | 761 |
|     | Dylan Slack and Sameer Singh. 2023. <a href="#">Tablet: Learning from instructions for tabular data</a> . <i>arXiv</i> .   | 762 |
|     |  | 763 |
|     | Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. <a href="#">Transformers learn in-context by gradient descent</a> . In <i>International Conference on Machine Learning</i> , pages 35151–35174. PMLR.  | 764 |
|     |  | 765 |
|     |  | 766 |
|     |  | 767 |
|     |  | 768 |
|     |  | 769 |
|     | Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. <a href="#">Finetuned language models are zero-shot learners</a> . In <i>International Conference on Learning Representations</i> .  | 770 |
|     |  | 771 |
|     |  | 772 |
|     |  | 773 |
|     |  | 774 |

|     |   |  |     |
|-----|---|--|-----|
| 775 | Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen,                               | <b>A Dataset Description</b>   | 824 |
| 776 | Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu, Denny                                   | We provide a detailed description of each feature  | 825 |
| 777 | Zhou, Tengyu Ma, et al. 2023a. Symbol tuning im-                                  | from the datasets evaluated in our paper.  | 826 |
| 778 | proves in-context learning in language models. <i>arXiv</i>                       |  |     |
| 779 | <i>preprint arXiv:2305.08298</i> .  |  |     |
| 780 | Jerry W. Wei, Jason Wei, Yi Tay, Dustin Tran, Albert                              | <b>A.1 Adult</b>   | 827 |
| 781 | Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,                                      | The original Adult Income Dataset contains 14 fea-   | 828 |
| 782 | Da Huang, Denny Zhou, and Tengyu Ma. 2023b.                                       | tures and the target <i>Income</i> , as described in Table 4.                                      | 829 |
| 783 | <a href="#">Larger language models do in-context learning dif-</a>                | Following prior work (Slack and Singh, 2023), we   | 830 |
| 784 | <a href="#">ferently</a> . <i>ArXiv</i> , abs/2303.03846.                         | omit <i>Education-Num</i> and <i>Fnlwgt</i> as they are not  | 831 |
| 785 | Laura Weidinger, John Mellor, Maribeth Rauh, Conor                                | crucial for income prediction, along with <i>Race</i> and  | 832 |
| 786 | Griffin, Jonathan Uesato, Po-Sen Huang, Myra                                      | <i>Native-Country</i> , to center our attention on <i>Sex</i> as                                   | 833 |
| 787 | Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,                                | the protected attribute.   | 834 |
| 788 | Zac Kenton, Sasha Brown, Will Hawkins, Tom  |  |     |
| 789 | Stepleton, Courtney Biles, Abeba Birhane, Julia                                   | <b>A.2 German Credit</b>   | 835 |
| 790 | Haas, Laura Rimell, Lisa Anne Hendricks, William                                  | The original German Credit Dataset contains 20   | 836 |
| 791 | Isaac, Sean Legassick, Geoffrey Irving, and Iason                                 | features, as detailed in Table 5. For simplicity and   | 837 |
| 792 | Gabriel. 2021. <a href="#">Ethical and social risks of harm from</a>              | consistency with prior work, only the features not   | 838 |
| 793 | <a href="#">language models</a> .   | shown in <i>italics</i> are retained in our work. Further-   | 839 |
| 794 | Sang Michael Xie, Aditi Raghunathan, Percy Liang,                                 | more, we extract <i>Sex</i> as an additional protected   | 840 |
| 795 | and Tengyu Ma. 2022. <a href="#">An explanation of in-context</a>                 | attribute from the <i>Personal Status and Sex</i> feature.   | 841 |
| 796 | <a href="#">learning as implicit bayesian inference</a> . In <i>International</i> |  |     |
| 797 | <i>Conference on Learning Representations</i> .                                   |  |     |
| 798 | Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Se-                                | <b>A.3 COMPAS</b>  | 842 |
| 799 | bastian Riedel. 2020. <a href="#">TaBERT: Pretraining for joint</a>               | The raw COMPAS Recidivism dataset contains   | 843 |
| 800 | <a href="#">understanding of textual and tabular data</a> . In <i>Proceed-</i>    | more than 50 attributes. Following the approach  | 844 |
| 801 | <i>ings of the 58th Annual Meeting of the Association</i>                         | of (Larson et al., 2016), we carry out the neces-  | 845 |
| 802 | <i>for Computational Linguistics</i> , pages 8413–8426, On-                       | sary preprocessing. More specifically, we group  | 846 |
| 803 | line. Association for Computational Linguistics.                                  | the <i>race</i> attribute into <i>African-American</i> and <i>Not</i>                              | 847 |
| 804 | Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-                                | <i>African-American</i> , and consider only the features   | 848 |
| 805 | donez, and Kai-Wei Chang. 2017. <a href="#">Men also like</a>                     | <i>sex</i> , <i>race</i> , <i>age</i> , <i>charge degree</i> , <i>priors count</i> , <i>risk</i> , | 849 |
| 806 | <a href="#">shopping: Reducing gender bias amplification using</a>                | and <i>two-year recid</i> (target). We frame the task  | 850 |
| 807 | <a href="#">corpus-level constraints</a> . In <i>Proceedings of the 2017</i>      | as predicting whether an individual will recidivate  | 851 |
| 808 | <i>Conference on Empirical Methods in Natural Lan-</i>                            | within two years ( <i>Did Not Reoffend</i> or <i>Reoffended</i> ),                                 | 852 |
| 809 | <i>guage Processing</i> , pages 2979–2989, Copenhagen,                            | based on their demographic and criminal history.   | 853 |
| 810 | Denmark. Association for Computational Linguis-                                   | Due to page limitations, we provide descriptions   | 854 |
| 811 | tics.   | for only the features used in our work in Table 6.   | 855 |
| 812 | Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and                                | <b>B RF and NN Hyperparameters</b>   | 856 |
| 813 | Sameer Singh. 2021. Calibrate before use: Improv-                                 | For RF, we fix the number of trees to 100 for all  | 857 |
| 814 | ing few-shot performance of language models. In <i>Inter-</i>                     | datasets as well as models. For NN, we use a 3   | 858 |
| 815 | <i>national Conference on Machine Learning</i> , pages                            | hidden-layered network with hyperparameters de-  | 859 |
| 816 | 12697–12706. PMLR.  | scribed in Table 7.  | 860 |
| 817 | Caleb Ziems, William Held, Jingfeng Yang, Jwala                                   | <b>C Prompt Templates</b>  | 861 |
| 818 | Dhamala, Rahul Gupta, and Diyi Yang. 2023. <a href="#">Multi-</a>                 | Beyond the Adult dataset, we provide the serializa-  | 862 |
| 819 | <a href="#">VALUE: A framework for cross-dialectal English</a>                    | tion and prompt templates utilized in our work for   | 863 |
| 820 | <a href="#">NLP</a> . In <i>Proceedings of the 61st Annual Meeting</i>            | the German Credit and COMPAS datasets here.  | 864 |
| 821 | <i>of the Association for Computational Linguistics</i>                           |  |     |
| 822 | <i>(Volume 1: Long Papers)</i> , pages 744–768, Toronto,                          | <b>C.1 German Credit</b>   | 865 |
| 823 | Canada. Association for Computational Linguistics.                                | <b>C.2 COMPAS</b>  | 866 |

| 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 4: Features in the original **Adult** dataset. Those not used in our work are shown in *italics*.

Predict the credit risk of a German bank customer based on their profile data. Answer with one of the following: bad | good.

Example 1 -  
 Age: 25 - 45  
 Sex: female  
 Job: highly skilled  
 Housing: rent  
 Saving accounts: little  
 Checking account: little  
 Credit amount: 2606  
 Duration: 21  
 Purpose: radio/TV  
 Answer: good

Age: 25 - 45  
 Sex: male  
 Job: skilled  
 Housing: own  
 Saving accounts: little  
 Checking account: little  
 Credit amount: 1345  
 Duration: 18  
 Purpose: radio/TV  
 Answer:

Figure 2: Prompt Template for **German Credit** Dataset.

Predict whether an individual will recidivate with in 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 3: Prompt Template for **COMPAS** Dataset.

| Feature                           | Type        | Description  |
|-----------------------------------|-------------|--|
| Credit Amount                     | Continuous  | The amount of credit requested by the applicant.   |
| Duration                          | Continuous  | The duration of the credit in months.  |
| <i>Installment Rate</i>           | Ordinal     | The installment rate in percentage of disposable income.   |
| <i>Residence Since</i>            | Ordinal     | The number of years the applicant has lived at their current residence.  |
| Age                               | Continuous  | The age of the applicant.  |
| <i>Number of Existing Credits</i> | Ordinal     | The number of existing credits at this bank.   |
| <i>Number of Dependents</i>       | Ordinal     | The number of dependents of the applicant.   |
| Checking Account Status           | Categorical | The status of the applicant’s checking account, such as “no checking,” “<0 DM,” “0-200 DM,” or “no known checking.”                          |
| <i>Credit History</i>             | Categorical | The credit history of the applicant, including categories like “critical/other existing credit,” “existing paid,” “delayed previously,” etc. |
| Purpose                           | Categorical | The purpose of the credit, such as “radio/tv,” “education,” “new car,” etc.  |
| Savings Account                   | Categorical | The status of the applicant’s savings account/bonds, including categories like “unknown/none,” “<100 DM,” “500-1000 DM,” etc.                |
| <i>Employment Since</i>           | Categorical | The duration of the applicant’s current employment, such as “unemployed,” “<1 year,” “4-7 years,” etc.                                       |
| <i>Personal Status and Sex</i>    | Categorical | The personal status and sex of the applicant, including categories like “male single,” “female div/dep/mar,” etc.                            |
| <i>Other Debtors/Guarantors</i>   | Categorical | Indicates the presence of other debtors/guarantors, such as “none,” “guarantor,” “co applicant.”   |
| <i>Property</i>                   | Categorical | Describes the type of property owned by the applicant, such as “real estate,” “life insurance,” “car or other,” etc.                         |
| <i>Other Installment Plans</i>    | Categorical | The presence of other installment plans.   |
| Housing                           | Categorical | The housing situation of the applicant, such as “own,” “for free,” and “rent.”   |
| Job                               | Categorical | The type of job held by the applicant, including categories like “skilled,” “unskilled resident,” “high qualif/self emp/mgmt,” etc.          |
| <i>Telephone</i>                  | Binary      | Indicates whether the applicant has a telephone (yes/no).  |
| <i>Foreign Worker</i>             | Binary      | Indicates whether the applicant is a foreign worker (yes/no).  |
| Risk (target)                     | Binary      | The target variable indicating credit risk (good/bad).   |

Table 5: Features in the original **German Credit** dataset. Those not used in our work are shown in *italics*. Additionally, from the original feature *Personal Status and Sex*, we extract *Sex* as a protected attribute.

| 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 6: Features in the **COMPAS** Recidivism Dataset (Preprocessed).

|                      | <b>h1</b> | <b>h2</b> | <b>h3</b> | <b>lr</b> | <b>batch size</b> | <b>epochs</b> |
|----------------------|-----------|-----------|-----------|-----------|-------------------|---------------|
| <b>Adult</b>         | 16        | 64        | 16        | 0.07      | 128               | 300           |
| <b>German Credit</b> | 64        | 64        | 32        | 0.07      | 128               | 300           |
| <b>COMPAS</b>        | 64        | 128       | 64        | 0.09      | 128               | 300           |

Table 7: 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 the number of epochs the models are trained for.