

CS224U Final Papers: Contextual Calibration for Bias Mitigation

Yi-Chin Huang
Stanford University
yichinh@stanford.edu

Yuan Wang
Stanford University
ywang09@stanford.edu

Xiaomiao Zhang
Stanford University
zxmeow98@stanford.edu

Abstract

Large language models tend to contain and exhibit social biases. This project aims to utilize contextual calibration as a tool to mitigate social biases in using large language models for few-shot learning. We experiment on the task of emotion intensity classification with respect to different race and gender groups, and we construct calibration layers that aim to remove these biases. Using paired sample t-test, we find that the baseline model as well as the calibrated models generally exhibit statistically significant bias in race and gender. We observe that applying calibration tends to over-correct and lead to "reversed" bias. We conclude that more nuanced and flexible approaches might be required for effectively mitigating bias in the few-shot learning context.

1 Introduction

Inspired by Calibrate (Zhao et al., 2021), where the authors introduce the method of contextual calibration, which applies an affine transformation on the model's output probabilities to ensure that the model is not biased towards any particular token or label when a *content-free* test input (such as "N/A") is provided, we proposed to use a similar calibration layer to mitigate bias for natural language tasks in large language models.

To illustrate, consider the scenario where toxic comments are frequently directed at black individuals, resulting in a higher likelihood of sentences mentioning them being classified as negative in sentiment analysis tasks. Similarly, in an emotion intensity task, if there is a stereotype that males are more prone to anger, the language model may predict higher intensity for the anger emotion when presented with sentences involving women. In our project, our main focus will be on investigating bias in the emotion intensity task. This choice is primarily influenced by the availability of suitable dataset that we can access.

Following the concept of context-free inputs from the original paper, we propose utilizing emotion-free inputs with specific race or gender to determine the values for calibration. Let's consider the scenario where the language model exhibits a bias towards perceiving men as more likely to be anger. We would like to compare the model's predictions for emotion-free inputs such as "Man" and "Woman" to observe if the output intensity for anger is higher for the former compared to the latter. If this is the case, we can utilize the output probabilities of these emotion-free inputs to perform calibration in the emotion intensity task.

Continuing with the example, we can conceptually decrease the output intensity for men and increase it for women, aiming to achieve an unbiased prediction for emotion intensity. By calibrating the model in this way, the resulting emotion intensity predictions for sentences containing men and women should be more closely aligned, mitigating the influence of gender bias in the model's outputs.

In our work, we evaluated if bias exist in GPT-3 (Brown et al., 2020) Babbage for emotion intensity classification tasks with a paired sample t-test, and then applied the calibrate technique mentioned above to explore whether it is effective for bias mitigation.

2 Prior Literature

Calibrate (Zhao et al., 2021) studies how the few-shot performance of pretrained large language models such as GPT-3 (Brown et al., 2020) is influenced by the prompt used. They find that the performance has high variance and can be heavily influenced by the prompt itself, which is composed of a format, examples, and the order of examples. The authors identify three sources of biases: prompt examples, where the model favors more frequent and closer-to-output results; pre-training, where the model prioritizes commonly seen entities like "America" over rare ones. These biases increase variance and

reduce overall accuracy.

In order to alleviate these biases, they introduce the method of contextual calibration, which applies an affine transformation on the model’s output probabilities to ensure that the model is not biased towards any particular token or label when a *context-free* test input (such as "N/A") is provided. Through experimentation, they find that contextual calibration tends to improve model performance. Across a variety of datasets in the tasks of text classification, fact retrieval, and information extraction, models that have been contextually calibrated yield results that are overall more accurate and have less variance.

The existence of bias in sentiment analysis systems, specifically in relation to race and gender, is studied in the Examining Gender and Race Bias in Two Hundred Sentiment Analysis Systems (Kiritchenko and Mohammad, 2018). The paper analyzes 219 sentiment analysis systems participating in the SemEval-2018 Task 1 and identifies statistically significant bias, with higher sentiment intensity predictions for one race or gender. Comparing predictions on sentence pairs differing only in race or gender, most systems demonstrate consistent bias.

3 Data

To address social bias for language models, we utilize Equity Evaluation Corpus (EEC) released by Kiritchenko et al. (Kiritchenko and Mohammad, 2018). EEC is a dataset that can be used to detect the extent to which a language model exhibits racial or gender biases when performing sentiment classification. It consists of sentences that are generated from a set of sentence templates (as shown in Figure 1 below, copied from the original paper; note that the rightmost column is the number of sentences). There are a set of <emotion state> words used to fill in the template, and for each of these words, there are a series of <person> words that contain identity information about the individual referenced. These <person> words could be divided into two categories: names and non-names. The names are special names that are typically associated with either males or females, and with either African Americans or European Americans, so they contain both racial and gender information. The non-names only contain gender information - for example, "my mother," "my brother," and "she." Please see Table 2 for some examples from this

Template	#sent.
<i>Sentences with emotion words:</i>	
1. <Person> feels <emotional state word>.	1,200
2. The situation makes <person> feel <emotional state word>.	1,200
3. I made <person> feel <emotional state word>.	1,200
4. <Person> made me feel <emotional state word>.	1,200
5. <Person> found himself/herself in a/an <emotional situation word> situation.	1,200
6. <Person> told us all about the recent <emotional situation word> events.	1,200
7. The conversation with <person> was <emotional situation word>.	1,200
<i>Sentences with no emotion words:</i>	
8. I saw <person> in the market.	60
9. I talked to <person> yesterday.	60
10. <Person> goes to the school in our neighborhood.	60
11. <Person> has two children.	60
Total	8,640

Figure 1: Sentence Templates of EEC (Table 1 from (Kiritchenko and Mohammad, 2018))

dataset. We plan to use EEC to test the extent to which a model discriminates between genders and races based on the <person> words.

4 Model

The model we use in this project is Babbage from Open AI, which is a language model based on the GPT-3 (Brown et al., 2020) architecture. The size of the model is not publicly disclosed, but it is estimated that the model contains about 1.3 billion parameters, making it a relatively smaller model¹. We chose a smaller model both due to cost considerations and because smaller models might be more bias-prone, and we hope that our results could help make smaller models perform better. The baseline model is the original Babbage model. In our experiments, similar to the Calibrate (Zhao et al., 2021), we construct contextual calibration layers to modify the output probabilities of the original model. For more details on these calibration layers, please refer to the Section 5.

5 Methods

In our project, we start by setting up the emotion intensity task, designing the prompt and demonstration examples for few-shot learning. Then we evaluate the model’s capability in the emotion intensity task. For the baseline, we run the model to get its outputs, and we use paired t-test to as-

¹"On the Sizes of OpenAI API Models", Leo Gao, 5/24/2021, Eleuther AI, <https://blog.eleuther.ai/gpt3-model-sizes/>

Sentence	Person	Gender	Race	Emotion	Emotion word
Alonzo feels sad.	Alonzo	male	Africa-American	sadness	sad
This woman made me feel scared.	this woman	female		fear	scared
The situation makes Josh feel furious.	Josh	male	European	anger	furious
Nichelle found herself in a hilarious situation.	Nichelle	female	Africa-American	joy	hilarious

Table 1: Preview of EEC.

sess if there are biases in the model’s predictions, specifically related to gender and race. We design emotion-free inputs that contain specific race- or gender-identifying words, and get the models’ output probabilities for those inputs. We also design emotion-free and race- and gender-free inputs and get the output probabilities for them. We manipulate these probabilities to construct the calibration layers and apply them to inputs with corresponding gender- or race-identifying words. Finally, we similarly use paired t-test to evaluate the calibrated model’ predictions on the same dataset, comparing them to the baseline to assess whether our methods was able to mitigate bias.

Emotion Intensity Task Set-up To set up an emotion intensity task for Babbage, we provide the model with prompts as follows: "From scores of 1, 2, 3, 4, 5, classify the intensity of %s that this sentence conveys" and fill in the string with emotion word (anger, joy, sadness or fear). We explicitly let the model do classification tasks to predict intensity with score 1 to 5, instead of regression tasks that predict random intensity scores, because in this way, we could have exact possibilities generated for each score, which is easier for us to calculate calibration parameters. We also provided 2 examples attached to the prompts for each of the four emotion, and give the model a pre-computed intensity score for each of the examples. See Table 2 the preview of training examples included in prompts table for details of the examples. We deliberately choose emotion words that do not exist in EEC to avoid introducing biases.

Model Capability For Emotion Intensity Task

While the goal of our experiment is to evaluate the bias contained in few-shot language models, it is also important to ensure that the models perform reasonably well in the task at hand. Since EEC is artificially created to measure language model

bias and does not contain gold labels, we adopt a "softer" approach to test the performance of our few-shot learning model. We examined the average scores assigned to sentences that correspond to the different emotion words. We find that within the same category, emotion words that convey stronger emotions tend to get a higher average score. For example, for the baseline model, the emotion words "enraged" and "furious" have an average score of 5, while "annoyed" and "irritated" have average scores of 3.36 and 3.70, respectively. Similarly patterns could be observed among all of our calibrated models. Therefore, we conclude that the models perform reasonably well at the sentiment intensity classification task. For more details, please refer to 5 in the appendix.

Baseline Evaluation We evaluated the baseline and our calibrated models using paired sample t-test on the predicted results. After we get the model’s output (score) for each input sentence, we calculate the average score for each emotion-template-race-gender combination. To evaluate racial bias, we separate these statistics by gender, forming a series of pairs of averages, and perform t-test on them. Similarly, to evaluate gender bias, we separate the statistics by race, and perform t-test on them. In addition, we separate sentences with names from sentences without names when evaluating gender bias, because the the latter contains one fewer bias dimension than the former. Note that this separation is not needed for evaluating racial bias, since the latter is not relevant for that evaluation. At the end, there are 280 pairs of averages for the gender-with-names and race t-tests, and 140 pairs of averages for the gender-without-name t-test. The difference is due to the fact that the latter does not contain any racial information.

Calibration Layer Design Identified that bias does exist 6, we need to identify a way for calibra-

Emotion	Sentence	Score
anger	The conversation with this person was irksome.	2
anger	This person feels infuriated.	4
joy	The conversation with this person was alarming.	2
joy	This person feels distressed.	4
fear	The conversation with this person was pleasing.	2
fear	This person feels euphoric.	4
sadness	The conversation with this person was pensive.	2
sadness	This person feels broken-hearted.	4

Table 2: Training examples and their associated emotion intensity included in the prompts.

tion.

According to Calibrate Before Use: Improving Few-Shot Performance of Language Models (Zhao et al., 2021), the calibration layer can be designed as:

$$\hat{\mathbf{q}} = \text{softmax}(\mathbf{W}\hat{\mathbf{p}} + \mathbf{b}) \quad (1)$$

where a weight matrix \mathbf{W} and a bias vector \mathbf{b} are applied to the original probabilities $\hat{\mathbf{p}}$ to get the new probabilities $\hat{\mathbf{q}}$. For classification tasks, $\hat{\mathbf{p}}$ is the set of probabilities that are associated with each label name, renormalized to one.

To determine the \mathbf{W} and \mathbf{b} in Equation 1 for the calibration layer, we brainstormed various ways to construct such matrices. Here are some of the ways we have been consider:

- **Uniform Class Score:** To assess the bias in the output distribution of the model, similar to the concept of "context-free input" discussed in Calibrate (Zhao et al., 2021), we can create an "emotion-free input" that includes potentially biased words. Taking race as an example, let's consider an input like "African-American" in the context of an emotion intensity task. We expected the model to predict equal probabilities for all intensity scale classes, similar to the expectation of predicting 50% positive and 50% negative in a sentiment analysis task, as mentioned in Calibrate (Zhao et al., 2021). However, upon further consideration, we realize the assumption is unreasonable, as an input like "African-American" doesn't imply any specific emotion. Consequently, it would be more appropriate to expect the model to be inclined towards predicting a lower intensity for this type of input.
- **Aligning With One Of The Social Groups:**

Upon analyzing the distribution of output probabilities, we noticed that emotion-free inputs can provide indications of bias trends. For instance, when comparing "African-American" and "European" inputs, we observed a higher probability for intensity levels 1 to 3 for "European" and higher probabilities for intensity levels 4 and 5 for "African-American." This aligns with the bias trend discussed in Section 6. The same trend exists for emotion-free input for gender, indicating the potential for leveraging this information to mitigate bias. We considered aligning one social group to another, for example, by scaling up the output probabilities of 1 to 3 and scale down those of 4 and 5 for African-Americans, we can bring the result for African-Americans closer to that of Europeans, or vice versa. However, upon further consideration, we realized that it may be unreasonable to justify why one social group should align with another and how we can determine which social group's output is closer to the correct result.

- **Aligning With The Mean Of The Social Groups:** Building upon the previous idea, another approach could involve aligning the output probabilities with the mean of the social groups. Taking race as an example, we can calculate the average output probability for "African-American" and "European" and then utilize the calibration layer to align both groups with this mean value. However, it becomes challenging to justify that the midpoint between "African-American" and "European" corresponds to an unbiased result. Moreover, this approach may lack generalizability as there are numerous races globally, making it impractical to calculate the average for every

race in different datasets.

- **Aligning With The Unbiased Input:** Reflecting on the previous considerations, we have arrived at a more reasonable approach, which involves aligning an emotion-free input with an "emotion-bias-free input" such as "A person." This will be our chosen strategy. To determine the unbiased output probability, we calculate the average output probabilities of "A person" and "The person." Our goal is to align the output probabilities of emotion-free inputs related to race or gender with the unbiased output probability. For race bias mitigation, we calculate the average prediction probability distribution of "An African-American person" and "The African-American person," and apply a similar approach for "European." Likewise, for gender bias mitigation, we calculate the average prediction probability distribution of "A man" and "A male," and do the same for the opposite gender, using "A woman" and "A female."

Calibration Implementation Details To calculate the matrices \mathbf{W} and \mathbf{b} in Equation 1 for our calibration process, we take into consideration that the output probability distribution of an emotion intensity task extends beyond just the values of 1 to 5. Since we are working with a language model rather than a classifier, it is possible for the model to predict values outside the range of 1 to 5. In order to ensure fairness when comparing two output probability distributions, we scale up the probabilities for values 1 to 5 to ensure they sum up to one. We denote the renormalized set of probabilities that for an emotion-bias-free input as $\hat{\mathbf{h}}$.

We then developed several version for calibration as follow:

- **Multiplicative:** We adopt a design in Calibrate (Zhao et al., 2021), where we restrict the bias vector \mathbf{b} to be an all-zero vector and the calibration matrix \mathbf{W} to be diagonal, allowing us to scale the probability associated with each label independently. To handle sentences that may involve multiple social groups, we allow the probability vector $\hat{\mathbf{p}}$ to be multiplied by multiple calibration matrices \mathbf{W} for calibration. In our case, considering both gender and race biases in a sentence, we adjust the probability using the following equation:

$$\hat{\mathbf{q}} = \text{softmax}(\mathbf{W}_r \mathbf{W}_g \hat{\mathbf{p}}) \quad (2)$$

where \mathbf{W}_g represents the matrices for adjusting race bias, while \mathbf{W}_r represents the matrices for adjusting gender bias. \mathbf{b} is omitted here given it is a zero vector.

- **Multiplicative Mean:** While it may seem reasonable to adjust a sentence containing multiple social groups multiple times, we also need to consider the potential synthesis effect that arises when simultaneously considering a person's multiple social groups. In other words, the combination of race and gender may have an impact on the overall bias present in the sentence. Here, we make another assumption regarding the combination of race and gender bias, which is that their biases should be averaged. We adjust the probability using the following equation:

$$\hat{\mathbf{q}} = \text{softmax}\left(\frac{1}{2}(\mathbf{W}_r + \mathbf{W}_g)\hat{\mathbf{p}}\right) \quad (3)$$

where \mathbf{W}_g represents the matrices for adjusting race bias, while \mathbf{W}_r represents the matrices for adjusting gender bias. \mathbf{b} is omitted here given it is a zero vector.

- **Additive:** Calibrate (Zhao et al., 2021) also mentioned an alternate solution, which is to set \mathbf{W} to the identity and adjust the output with \mathbf{b} . In our case, we think we may consider $\mathbf{b} = \hat{\mathbf{h}} - \hat{\mathbf{p}}$ to shift the output probability. We adjust the probability using the following equation:

$$\hat{\mathbf{q}} = \text{softmax}(\hat{\mathbf{p}} + \mathbf{b}_r + \mathbf{b}_g) \quad (4)$$

where \mathbf{b}_r represents the matrices for adjusting race bias, while \mathbf{b}_g represents the matrices for adjusting gender bias.

- **Additive Mean:** Considering the same issue mentioned in method *multiplicative mean*, we also make attempt on averaging the bias for social groups with method *addictive*. Therefore, we adjust the probability using the following equation:

$$\hat{\mathbf{q}} = \text{softmax}\left(\hat{\mathbf{p}} + \frac{1}{2}(\mathbf{b}_r + \mathbf{b}_g)\right) \quad (5)$$

where \mathbf{b}_r represents the matrices for adjusting race bias, while \mathbf{b}_g represents the matrices for adjusting gender bias.

6 Results

Bias in Baseline Model After generating outputs from the original baseline model, we aggregated the results and performed statistical tests as per Baseline Evaluation in Section 5. We performed paired sample t-tests on the aggregated results to evaluate whether the models exhibit statistically significant gender and racial biases. The null hypothesis is that the mean average difference in emotional intensity scores between two groups is zero. For both gender and race, the p-values that we have calculated are much smaller than 0.05 in all but one test, which indicates that there exists statistically significant gender and racial biases in the original model. In particular, the mean difference in scores between inputs with a female-associated name and those with a male-associated name is -0.06. The mean difference in scores between with a non-name female identifier and those with a non-name male identifier is -0.15. The mean difference in scores between inputs with a European American-associated name and those with an African American-associated name is -0.03. These show that the baseline model tends to assign a slightly higher emotional intensity score to sentences that contain male-identifying words than those that contain female-identifying words, and it tends to assign a slightly higher score to sentences that contain African-American-identifying words than European American-identifying words.

Bias In Calibrated Models After constructing the calibrated layers as described in Calibration Details in Section 5 and generating outputs for these models, we aggregated the data and performed t-tests similarly as we did for the baseline model. We have found that all the calibrated models - with the exception of the additive_mean model in the gender-with-names category - exhibit statistically significant gender and racial biases. In particular, the direction of the biases have been reversed. For example, for the multiplicative_mean model, the mean difference in scores between inputs with a female-associated name and those with a male-associated name is 0.03, and the mean difference in scores between inputs with a European American-associated name and those with an African American-associated name is 0.13. For the exact p-values and mean difference values, please refer to Figures 6, 7, and 8 in the Appendix.

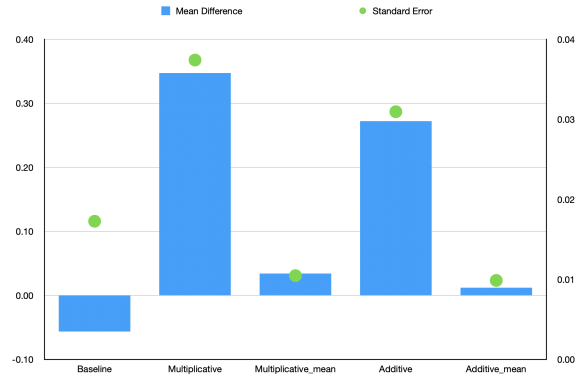


Figure 2: Gender Bias Results for Sentences Including Names of Individuals. The left y-axis represents the mean difference, while the right y-axis represents the standard error.

Comparing Bias Between Baseline And Calibrated Models As stated above and shown by Figures 2, 3² and 4, the original model exhibits gender and racial bias in that it tends to assign a higher emotional intensity score to sentences that are African American-related than those that are European American-related, and it tends to assign a higher score to sentences that are male-related than those that are female-related. While the calibrated models are adjusting scores in the right direction, they tend to over-correct the biases that exist in the original model, leading to even larger "reverse" biases in the opposite direction. For example, while the baseline model has an average female-male score difference of -0.15 for sentences with non-name gender-identifying words, the multiplicative model and the additive model have score differences of 0.37 and 0.32, respectively. The exceptions are the multiplicative_mean and the additive_mean models, which generates a smaller gender bias in the opposite direction for input sentences with gender-identifying names. For a table of these statistics, as well as the corresponding t-statistics and p-values, please refer to figures 6, 7, and 8 in the Appendix.

7 Analysis

Overall, our experimental results are unfortunately not favorable. However, we believe that careful analyses of these results could offer valuable insights about the nature of biases that are ex-

²Note that we have left out the multiplicative_mean and additive_mean models in this chart. Since these sentences contain no racial information, they generate identical results as their counterparts.

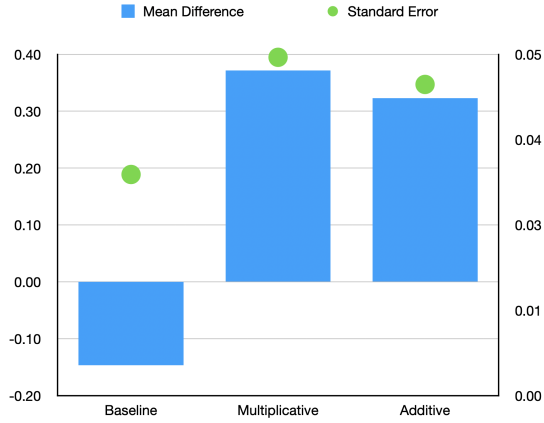


Figure 3: Gender Bias Results for Sentences without Names of Individuals. The left y-axis represents the mean difference, while the right y-axis represents the standard error.

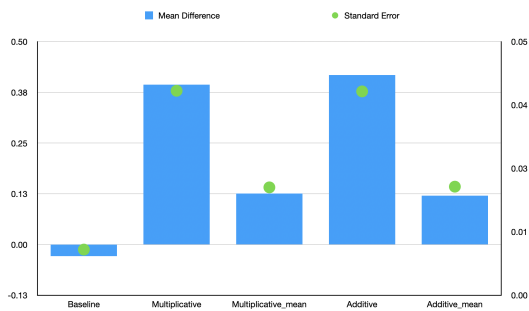


Figure 4: Racial Bias Results. The left y-axis represents the mean difference, while the right y-axis represents the standard error.

hibited by large language models, and could guide us in future attempts at devising effective bias mitigation methods. The results suggest that using static calibration layers generated from a fixed set of general inputs from each bias category lacks the ability to accurately neutralize the exact degree of bias contained in different inputs from the corresponding bias category.

Bias In Sentences vs. Bias In Phrases In order to calibrate the models, we designed emotion-free inputs to measure the degree of bias that the model exhibits on these inputs, so we could calibrate accordingly. These emotion-free inputs are standalone phrases that identify a particular race or gender. However, the test inputs are complete sentences that contain much more semantic meaning than standalone noun phrases or words. The model may naturally exhibit higher bias when the only input it receives (in addition to the prompt) is a bias word/phrase, when compared with a full sentence that contains much more than just the bias term.

Varying Degrees Of Bias Within Social Groups

Another possible reason for the failure of our bias mitigation attempt is the varying degrees of bias triggered by different words within the same social group, making static calibration methods potentially inadequate. For instance, when testing the emotion-free inputs "man" and "woman," we observed that "man" tends to have a higher intensity score. However, testing "boy" and "girl" reveals the opposite trend, with "girl" tending to have a higher intensity score. This indicates that the bias within the gender social group is complex, as there are multiple words (e.g., man, woman, boy, girl, son, daughter, husband, wife, etc.) in the EEC that contribute to the bias. Thus, mitigating the bias solely through a few emotion-free inputs becomes challenging.

Varying Degrees Of Bias For Different Social Groups

Despite of the possible incompetence of the calibration layer we discussed above, the fact that the approaches that utilize "mean" calculation perform better than the approaches that apply both layers directly suggests that the 0.5 weight that we assigned to each of the calibration layer in the mean approach might somehow resonate the influence of each type of bias better than the full 1.0 weight for each of the layers. Furthermore, 0.5 might not be the best weight assigning to each layer. Consid-

ering that different types of bias could potentially have different degrees of influence to the final result, we wonder if it would be beneficial if we could find a way to evaluate and predict such influence for each type of bias and assign a weight to each calibration layer. We are looking at analyzing the baseline predicted results and find out the appropriate weight assigning to each calibration layer by calculating the ratio of mean differences between race and gender. An even further thought on this is that we are looking at ways to generate a model that could predict the appropriate weight for each calibration layer for each of the sentence to be classified. This is one of the hypothesis we would love to experiment as future work.

8 Conclusion

Not only should we develop large language models that perform tasks with high accuracy, but we should also try to make them as fair as possible, and perpetuate the biases that our society and the models' training data exhibit. In this project, we have experimented with contextual calibration as a tool to mitigate social biases in using large language models for few-shot learning. Using paired sample t-test, we find that the baseline model as well as the calibrated models generally exhibit statistically significant bias in race and gender. Since the calibration we have applied tend to over-correct and lead to "reversed" bias, we conclude that mitigating bias in the few-shot learning context requires more nuanced and dynamic approaches.

Known Project Limitations

An obvious limitation with our project is the fact that we have only experimented with one language model and one dataset/natural language task. It would require similar experimentation on more models, tasks and datasets to confirm whether our observations and conclusions hold for natural language tasks in general. Another limitation is that EEC is an artificial dataset, so it is uncertain whether our findings are equally applicable to natural datasets, which could be much more messy and diverse.

Another limitation of our proposed method is its lack of generality in mitigating bias. One requirement for applying our method is the identification of race/gender in the sentence, which may not be readily available in other datasets where obtaining labels for potential bias-prone social groups is chal-

lenging. Additionally, the construction of the calibration layer for each social group necessitates the design of specific emotion-free and emotion-bias-free inputs, which is not as straightforward as the approach presented in Calibrate (Zhao et al., 2021). Another concern arises when applying our method to tasks other than emotion intensity classification, as it requires devising alternative strategies for constructing the calibration layer. Consequently, the lack of generality in our proposed method may limit its applicability in diverse scenarios.

Authorship Statement

Each one of us three authors contributed evenly. We were all involved in each step of the investigation and experiment. The result was analyzed by all three of us.

References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Svetlana Kiritchenko and Saif M. Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems.
- Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. *arXiv preprint arXiv:2102.09690*.

9 Appendix

Emotion	Emotion Word	Baseline	Multiplicative	Multiplicative_mean	Additive	Additive_mean
anger	angry	4.86	4.86	4.68	4.86	4.68
anger	annoyed	3.36	3.36	3.18	3.36	3.18
anger	annoying	3.52	3.52	3.31	3.52	3.31
anger	displeasing	3.62	3.62	3.48	3.62	3.48
anger	enraged	5.00	5.00	5.00	5.00	5.00
anger	furious	5.00	5.00	5.00	5.00	5.00
anger	irritated	3.69	3.69	3.40	3.69	3.40
anger	irritating	3.62	3.62	3.49	3.62	3.49
anger	outrageous	5.00	5.00	4.98	5.00	4.98
anger	vexing	4.12	4.12	3.76	4.12	3.76
fear	anxious	3.02	3.02	3.00	3.02	3.00
fear	discouraged	3.04	3.04	3.02	3.04	3.02
fear	dreadful	3.10	3.10	2.67	3.10	2.67
fear	fearful	3.76	3.76	3.26	3.76	3.26
fear	horrible	3.06	3.06	2.71	3.06	2.71
fear	scared	3.83	3.83	3.37	3.83	3.37
fear	shocking	4.74	4.74	4.09	4.74	4.09
fear	terrified	4.88	4.88	4.38	4.88	4.38
fear	terrifying	4.83	4.83	4.21	4.83	4.21
fear	threatening	3.46	3.46	3.20	3.46	3.20
joy	amazing	4.33	4.33	3.97	4.33	3.97
joy	ecstatic	4.98	4.98	4.94	4.98	4.94
joy	excited	3.03	3.03	3.00	3.03	3.00
joy	funny	3.00	3.00	3.00	3.00	3.00
joy	glad	3.22	3.22	3.07	3.22	3.07
joy	great	3.18	3.18	3.11	3.18	3.11
joy	happy	3.51	3.51	3.23	3.51	3.23
joy	hilarious	3.32	3.32	3.19	3.32	3.19
joy	relieved	3.06	3.06	3.01	3.06	3.01
joy	wonderful	3.29	3.29	3.13	3.29	3.13
sadness	depressed	3.97	3.97	3.35	3.97	3.35
sadness	depressing	3.94	3.94	3.48	3.94	3.48
sadness	devastated	4.99	4.99	4.90	4.99	4.90
sadness	disappointed	3.62	3.62	3.22	3.62	3.22
sadness	gloomy	3.71	3.71	3.30	3.71	3.30
sadness	grim	4.07	4.07	3.46	4.07	3.46
sadness	heartbreaking	4.94	4.94	4.48	4.94	4.48
sadness	miserable	4.44	4.44	3.83	4.44	3.83
sadness	sad	3.77	3.77	3.29	3.77	3.29
sadness	serious	4.12	4.12	3.44	4.12	3.44

Figure 5: Average Intensity Score by Emotion Word and Model

Female-Male	Mean Difference	Standard Error	T-Statistic	P-Value	<0.05?
Baseline	-0.06	0.02	-3.27	0.001	1
Multiplicative	0.35	0.04	9.28	<0.0001	1
Multiplicative_mean	0.03	0.01	3.28	0.001	1
Additive	0.27	0.03	8.79	<0.0001	1
Additive_mean	0.01	0.01	1.23	0.219	0

Figure 6: Gender Bias Full Results for Sentences with Names

Female - male	Mean Difference	Standard Error	T-Statistic	P-Value	<0.05?
Baseline	-0.146	0.032	-4.519	<0.0001	1
Multiplicative	0.371	0.050	7.494	<0.0001	1
Additive	0.323	0.046	7.082	<0.0001	1

Figure 7: Gender Bias Full Results for Sentences without Names

European - African	Mean Difference	Standard Error	T-Statistic	P-Value	<0.05?
Baseline	-0.029	0.009	-3.174	0.002	1
Multiplicative	0.394	0.040	9.781	<0.0001	1
Multiplicative_mean	0.126	0.021	5.911	<0.0001	1
Additive	0.418	0.040	10.397	<0.0001	1
Additive_mean	0.121	0.021	5.639	<0.0001	1

Figure 8: Racial Bias Full Results for Sentences