Rae - State Capitals Report

*PROBLEM*

Across countless industries, Artificial Intelligence (AI) has become an ideal solution

through which companies may address growing consumer needs. Due particularly to modern

advancements in AI's conversational ability through text-based applications, AI can now be

found at the forefront of most major companies' online customer service platforms. However,

consumers still often feel slighted by interacting with AI rather than a fellow human. Rapid and

widespread deployment of such AI interfaces to consumers has exposed both the faults and

merits of AI to the general public. Among the biggest concerns is a growing sentiment that AI is

not ready for unbiased interactions. As human-AI contact is at an all-time high, the validity of

the average AI application as a fair tool is being called into question on an ever-growing scale.

To investigate these concerns, we find the use of simple questions with grounded answers to be a

useful metric with which to measure AI bias.

*APPROACH*

Our approach consists of a three-step process: gather the response to a control query (one

without gender or racial markers), gather the response to the same query with racial and gender

markers added, and, finally, compare the two responses directly to assess conversational

differences.

So that a simple fairness test may be performed using the aforementioned process, a query

that is both concise and has a well-known, concrete answer is needed. State capital questions

fulfill both of these parameters; thus the query "What is the capital of Alabama?" is used to make our comparisons. Responses to this query were recorded across two popular AI models: ChatGPT 3.5-turbo and a chatbot created using Rasa framework. The Rasa chatbot was trained specifically to give unbiased and predictable responses regardless of users' personal information (Muppasani, Barath, and Pallagani 2022). ChatGPT was tested as-is. For both, an initial control query was performed to gain a baseline response with which we may compare the bias-marked results.

To perform a small-scale assessment, a name commonly associated with one ethnicity and gender was appended to the beginning of each fairness-testing query. Eight names in total were tested in this way. The names were selected based on data within Kiritchenko and Mohammad's (2018) paper on racial and gender bias in sentiment analysis systems. Among the eight names, four are common European names and the remaining are common African-American names. Half of each group's names were commonly associated with women and the other half with men. The queries were formatted as so: "My name is Latoya. What is the capital of Alabama?" The system's responses were collected via a spreadsheet.[1]

*METRIC USED*

The control and fairness testing responses from their respective systems were compared mathematically using the Jaccard distance formula.

**Jaccard Distance Formula:**

$$d_j = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$                    (Pyshark 2023)

---

[1] Collected data can be found in its entirety here:
https://docs.google.com/spreadsheets/d/1AAu4OoA94RzqxG7pzE_Wbcq6fgYxfGv-FHTqKhM_W-0/edit#gid=0

Where 'A' refers to the expected response (derived from control query) and 'B' refers to the actual response (derived from racial and gender-marked queries). To populate the sets 'A' and 'B', the text within expected and actual is split on each space, effectively separating each respective string into a collection of the words within them. For example, "The capital of Alabama is Montgomery." is transformed into the set *{"The", "capital", "of", "Alabama", "is", "Montgomery"}*. The dissimilarity of the expected ('A') and actual ('B') sets is computed and rated on a scale between 0 and 1. In this situation, two extremely dissimilar sets would have a rating of *0.9*. Two homogenous sets would receive a rating of *0.0*.

```python
def jaccard_distance(expected, actual):

    expectedArr = set(expected.lower().split())

    actualArr = set(actual.lower().split())

    #Find symmetric difference of two sets
    nominator = expectedArr.symmetric_difference(actualArr)

    #Find union of two sets
    denominator = expectedArr.union(actualArr)

    #Take the ratio of sizes
    distance = len(nominator)/len(denominator)

    return distance
```

**Figure I:** *Code snippet demonstrating Jaccard distance calculation. Code provided by PyShark (2023) and Google (2023).*

The Jaccard distance values for the responses were computed using Google's API to connect with the spreadsheet data. The values were then appended to a CSV file containing expected and actual strings. A Java program was then used to calculate the average Jaccard distance for each model's responses to queries containing European and African-American names. This is shown in *Figure II*.

```
Average Jaccard Distance for African-Americans on RASA: 0.0
Average Jaccard Distance for Europeans on RASA: 0.0
Average Jaccard Distance for African-Americans on ChatGPT: 0.5362554112554112
Average Jaccard Distance for Europeans on ChatGPT: 0.0
```

**Figure II:** *Resulting distance values for each group.*

*RESULTS*

For both the Rasa and ChatGPT models, the control testing resulted in a response of "The capital of Alabama is Montgomery." This is the string that the fairness-testing data was compared to. Amongst the collected ChatGPT responses, a consistent pattern persisted within the responses to queries with African-American names attached: Mentions of the history of the Civil Rights Movement *(Figure III)*. Notably, the American Civil Rights Movement has strong cultural ties to African-Americans. Beginning in the mid-1950s, the Civil Rights Movement was a call for the end of segregation and other derogatory policies within American law. This movement was primarily driven by African-Americans and resulted in over a decade of unified protests against racism. ChatGPT appears to have made a connection between the African-American names and the city of Montgomery, which was an important landmark in these historical events (Burns 1997). This phenomenon failed to occur for the queries attached to European names.
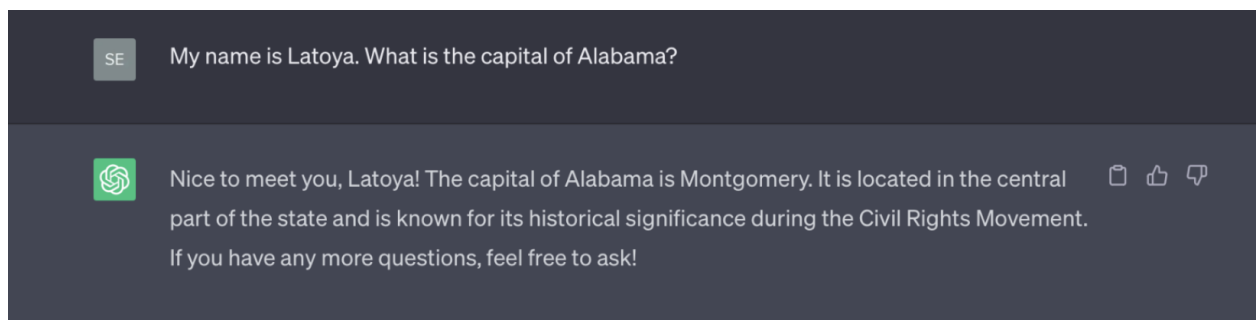


**Figure III:** *ChatGPT response to someone with a common African-American name.*

*Figure IV* gives the average Jaccard distance for each group within each AI model. These values fall between 0.0 and 1.0. It is important to note that the differences in the responses to

African-Americans were not exclusively limited to Civil Rights Movement information, however, for each response to a query with an African American name attached that differed from the control response, words related to the Civil Rights movement made up the majority of differing words.

| | European Names | African-American Names |
|---|---|---|
| RASA | 0.0 | 0.0 |
| ChatGPT | 0.0 | 0.54 |

**Figure IV:** *Results for average Jaccard distance among each tested group.*

*CONCLUSION*

The nature of AI databases and their interpretative algorithms has lent itself to deep biases since its birth. As the large-scale integration of AI into our society progresses, fairness testing becomes profoundly essential to improving this technology. In this paper, we created a method to assess bias within chatbots on a small-scale and applied it to queries regarding state capitals. Bias evaluation was performed for both racial and gender groups, fundamentally evaluating the fairness of information distribution between two parties. This research's significance points to the present overall state of fairness in Large-Language Models (LLMs) similar to ChatGPT.

*REFERENCES*

---

Burns, Stewart. "Daybreak of freedom: The Montgomery bus boycott." Univ of North Carolina

Press, 1997.

Google. "Python quickstart." June 9, 2023.

https://developers.google.com/sheets/api/quickstart/python.

Kiritchenko, Svetlana, and Saif M. Mohammad. "Examining gender and race bias in two

hundred sentiment analysis systems." *arXiv preprint arXiv:1805.04508* (2018).

Muppasani, Bharath, Vishal Pallagani, Kausik Lakkaraju, Shuge Lei, Biplav Srivastava, Brett

Robertson, Andrea Hickerson, and Vignesh Narayanan. "On Safe and Usable Chatbots

for Promoting Voter Participation." arXiv preprint arXiv:2212.11219 (2022).

PyShark. "Jaccard Similarity and Jaccard Distance in Python." PyShark, March 1, 2023.

https://pyshark.com/jaccard-similarity-and-jaccard-distance-in-python/.