

CS 232 Final Project: Analyzing the Potential Bias in GPT-J Language Model for Media Across the World

I. Introduction

Many aspects of culture in different areas of the world have been influenced by western media and culture. According to Stephen Galloway, an executive editor that featured for The Hollywood Reporter, claimed ten years ago that “compared to Hollywood productions, foreign films don’t play that well in their home markets”. Although foreign films have gained much attention, especially as the popularity of Korean culture increases, Hollywood is still dominating the box offices globally as it did ten years ago. As artificial intelligence is becoming more accessible to the public, I wanted to dive deeper into the potential bias that the language models can have in terms of culture. More specifically, I looked into the potential bias that may be present in language models for television shows and movies in different countries.

Previous research has already shown the bias present in language models. In Sheng’s paper, *Societal Biases in Language Generation: Progress and Challenges*, her team addresses the societal biases in language generation. Her team discusses the issue of current language generation technology that can negatively affect marginalized groups. She gives examples of chatbots that may produce more negative responses for certain groups of people that may discourage them from using this new technology, putting them at a technological disadvantage compared to other groups of people. Thus, she highlights the importance of quantifying the bias and societal impact of that bias of language generation models. In order to do this, her team focused on the techniques used for natural language generation (NLG) tasks. One negative impact of this bias is vulnerability. She discusses how privacy related issues and misinformation in texts generated from NLG can put a certain group at a disadvantage. Sheng addresses different bias mitigation strategies such as data-based mitigation and training-based mitigation but also mentions the main challenges of using these mitigation methods. Sheng and her team provided a commentary on the current progress of the biases in language generation and the associated societal impacts of it.

Another research team looked into the pitfalls of benchmarks, originally created to address the harm induced by the reproduction of hate speech and stereotype of language models. In Blodgett’s paper, *Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets*, her research team uses measurement modeling to analyze benchmark datasets that benchmark sentence templates that may have hate speech or stereotypes to measure how much natural language processing systems reproduce those languages. Her team dives into the details of what and how each of these benchmark dataset measures the reproduction of hate speech and stereotypes. They analyzed four different benchmark datasets to assess the level of stereotyping in a language model: StereoSet which looks into intra- and inter- sentence prediction tests, CrowS-Pairs which looks only at intra-sentence prediction tests, WinoBias (WB) which analyzes the pronouns that the model fills into sentence, and Winogender that assigns a gender based on the context of the sentence. Her team

highlighted the potential implicit decisions on what the benchmark deems as stereotyping and suggests that the reasonings behind those decisions are made explicit. They address the pitfalls of benchmarks that have the purpose of reducing harm due to the implicit conceptualization of stereotyping.

In Ted Underwood's response to the "Stochastic Parrots" paper, he discusses how neutral language models can potentially highlight the shared aspects of culture and behavior. Although models will most definitely have inherent bias, we can take advantage of this behavior instead of refraining from modeling by using it to capture the different perspectives of the world. For example, he shares an example of how machine learning can help historians measure the difference in their perspectives of history by training models. He also mentions the possibility of building lighter-weight models that can be more accessible to institutions or people that don't have the resources to train numerous models. For example, models can be designed to capture the social structure in the training dataset. He argues that the authors of "Stochastic Parrots" are not envisioning the possible benefits of neutral language models by refraining from it because it's dangerous.

Similar to the research teams that analyzed societal biases in language models and stereotyping of these models that can target marginalized groups, I will be analyzing the potential bias present in the GPT-J language completion model on sentence frames that focus on television shows, genres, and movies in different parts of the world. I hypothesize that due to Hollywood's great presence in media around the world and the highly likely, biased dataset that the GPT-J language model trained on, the model will have inherent bias towards westernized culture. However, it's important to keep in mind that the results from this project can also capture the different perspectives of the world even with the inevitable bias as Underwood explained.

II. Probe Task

In order to test my hypothesis, I focused on the GPT-J model's sentence completion predictions of the most likely next words given the context of the sentence frame. I constructed a dataset of 32 frame sentences that I used to operationalize the construct that I have chosen: tv show and movies. I created variants of each of these 32 frame sentences in 5 different countries to analyze my interest in media across different parts of the world.

I have show an example of a frame sentence I constructed down below along with its variants:

1. Frame sentence - I'm a twenty year old woman living in BLANK. My favorite television show genre is
 - a. Variant 1 - I'm a twenty year old woman living in California. My favorite television show genre is
 - b. Variant 2 - I'm a five year old girl living in Busan. My favorite television show genre is

- c. Variant 3 - I'm a five year old girl living in London. My favorite television show genre is
- d. Variant 4 - I'm a five year old girl living in Tecate. My favorite television show genre is
- e. Variant 5 - I'm a five year old girl living in Beijing. My favorite television show genre is
- f. neutral Version - I'm a five year old girl living in a city. My favorite television show genre is

For the example above, I chose to look at the most likely words that the model completes the sentence with to see whether the neutral version predicts itself in a specific location by comparing the probability of its predicted words compared to the probabilities of the predicted, country-specific words. If the neutral version predicted the genre 'crime', for example, and the first variant that locates the model in the United States predicts 'crime' with the closest probability, this will be an indication that there is bias in the model towards American culture.

More specifically, I constructed this dataset in a TSV file that my code can read in and calculate the probabilities of the next 5 most likely words based on the GPT-J model. I formatted my dataset as a TSV file according to the following fields below:

NUMBER EXAMPLE SENTENCE COUNTRY

- NUMBER is an ID number that is unique to each of the frame sentences.
- EXAMPLE SENTENCE is the country-specific sentence or neutral sentence for that frame sentence.
- COUNTRY is the name of the country that is being targeted.

In this dataset, I constructed example sentences that looked into the genres that different age groups watch, the most popular television shows in different countries, television shows with the most awards, television shows that made people most happy or angry or sad, television shows that people watch regularly at different times of the day, the most watched and least watched shows, and so on. Then, I took these frame sentences and its country-specific variants to get the probability distributions of the words that the GPT-J model generates to complete the sentence based on the context of the sentence.

Below, I have shown a few examples of the frame sentences I constructed in my TSV data set file:

13 A	I'm a high school student that lives in California. I like to go to the movie theater with my friends to watch movies in this genre:	US
13 B	I'm a high school student that lives in Busan. I like to go to the movie theater with my friends to watch movies in this genre:	South Korea
13 C	I'm a high school student that lives in London. I like to go to the movie theater with my friends to watch movies in this genre:	UK
13 D	I'm a high school student that lives in Tecate. I like to go to the movie theater with my friends to watch movies in this genre:	Mexico
13 E	I'm a high school student that lives in Beijing. I like to go to the movie theater with my friends to watch movies in this genre:	China
13 F	I'm a high school student that lives in the city. I like to go to the movie theater with my friends to watch movies in this genre:	Neutral

For frame sentence #13, I looked at the probability distributions of the next 5 likely words generated by the GPT-J model in 5 different countries: U.S., South Korea, United Kingdom, Mexico, and China. Here, I focused on high school students in these different countries and prompted the model to show the genre of movies they will watch with their friends at the movie theater in order to see what this age group in different parts of the world like to watch at the theater with their friends. Depending on the results, we can see whether the neutral sentence has bias towards a specific culture based on the probability distribution difference between the neutral version compared to the other country-specific versions.

2	A	I'm a twenty year old woman living in California. My favorite television show genre is	US
2	B	I'm a twenty year old woman living in Busan. My favorite television show genre is	South Korea
2	C	I'm a twenty year old woman living in London. My favorite television show genre is	UK
2	D	I'm a twenty year old woman living in Tecate. My favorite television show genre is	Mexico
2	E	I'm a twenty year old woman living in Beijing. My favorite television show genre is	China
2	F	I'm a twenty year old woman living in a city. My favorite television show genre is	Neutral

For frame sentence #2, I focused on a more specific age group, twenty year olds, and their favorite television show genre to test my hypothesis on whether the model held bias towards American culture in a specific age group of people in these 5 different countries.

After running this dataset using the GPT-J token query model that returns the next 5 likely words and its associated probabilities, I was able to get the following results for frame sentence #13 and frame sentence #2:

Frame sentence #13:

13	A	I'm a high school student that lives in California. I like to go to the movie theater with my friends to watch movies in this genre:	US	horror 0.19415296614170074	action 0.142045333
98151398		Horror 0.1056458030078888	comedy 0.08053655177354813	OTHER 0.3719734475016594	
13	B	I'm a high school student that lives in Busan. I like to go to the movie theater with my friends to watch movies in this genre:	South Korea	action 0.1694764494895935	horror 0.157385811
20967865		comedy 0.0795300155878067	action 0.06422026455402374	OTHER 0.468259047716856	
13	C	I'm a high school student that lives in London. I like to go to the movie theater with my friends to watch movies in this genre:	UK	horror 0.18967776000499725	action 0.136507511
13891602		Horror 0.10321082174777985	comedy 0.07868026196956635	OTHER 0.39763951301574707	
13	D	I'm a high school student that lives in Tecate. I like to go to the movie theater with my friends to watch movies in this genre:	Mexico	action 0.19499222934246063	horror 0.164065003
39508057		Horror 0.10697876662015915	action 0.0802225648164749	OTHER 0.376751147210598	
13	E	I'm a high school student that lives in Beijing. I like to go to the movie theater with my friends to watch movies in this genre:	China	action 0.2120305299758911	horror 0.116326496
00505829		comedy 0.0939332898616791	action 0.08834542202949524	OTHER 0.4487123563885689	
13	F	I'm a high school student that lives in the city. I like to go to the movie theater with my friends to watch movies in this genre:	Neutral	horror 0.19845937192440033	action 0.161578848
95801544		Horror 0.1027902215719223	comedy 0.08648684620857239	OTHER 0.366305448114872	

Frame sentence #2:

2	A	I'm a twenty year old woman living in California. My favorite television show genre is	US	comedy 0.1571657657623291	drama 0.14007385075092316	horror 0.0781232938170433	sc
ience		0.06154699996113777	reality 0.0517851822078228	OTHER 0.5113049075007439			
2	B	I'm a twenty year old woman living in Busan. My favorite television show genre is	South Korea	drama 0.15778659284114838	comedy 0.12127791345119476	romance 0.077789478	
00397873		horror 0.06281478703022003	action 0.05833357200026512	OTHER 0.521997656673193			
2	C	I'm a twenty year old woman living in London. My favorite television show genre is	UK	comedy 0.14707258343696594	drama 0.1446731835603714	crime 0.0761747658252716	ho
error		0.0755090636014930	the 0.05175481364130974	OTHER 0.5047737471759319			
2	D	I'm a twenty year old woman living in Tecate. My favorite television show genre is	Mexico	comedy 0.16785606741905212	drama 0.16785606741905212	action 0.06308538466691971	ho
error		0.05530757084488869	reality 0.049699850380420685	OTHER 0.4961950592696667			
2	E	I'm a twenty year old woman living in Beijing. My favorite television show genre is	China	comedy 0.1239340752363205	drama 0.1068817749619484	reality 0.07755596935749054	
		the 0.06579384952783585	crime 0.06211326643824577	OTHER 0.563721064478159			
2	F	I'm a twenty year old woman living in a city. My favorite television show genre is	Neutral	drama 0.15962855517864227	comedy 0.1557384580373764	horror 0.07071831077337265	cr
ime		0.059993188828229904	sitcom 0.053910382091999054	OTHER 0.5000111050903797			

III. Evaluation Metric

In order to quantify and measure the cultural bias of the GPT-J language model based on my constructed frame sentences, I used an evaluation metric that compares the probability distributions of the top 5 sentence completions that the GPT-J model produces. I measured the cultural biases of the GPT-J model by comparing a country-specific sentence to the neutral frame sentence using a list of the most likely or top 5 predicted words according to the GPT-J model. By using a metric that will evaluate

the top 5 sentence completions that GPT-J produces, I was able to calculate the probability differences between the country-specific frame sentences and the neutral version sentence frames for each of those frame sentences. For each frame sentence, GPT-J assigned the probabilities to all of the possible words that it can generate based on the context of the sentence, and took the top 5 words with the highest probability. In my evaluation metric, I took the sum of the for each country-specific frame sentence and the sum of the probabilities of the neutral version then took the difference between them to identify which country had the closed probability distribution to the neutral version. The country with the lowest difference is the country with the culture that the model is biased towards.

I further analyzed the probability distribution differences of the 5 countries I studied by changing the randomness of the sentence completions of the GPT-J model. To do this, I changed the temperature of the token query of the GPT-J model to see if this affected the results. The closer the temperature is to 1, the higher the randomness and the opposite is true the closer the temperature is to 0.

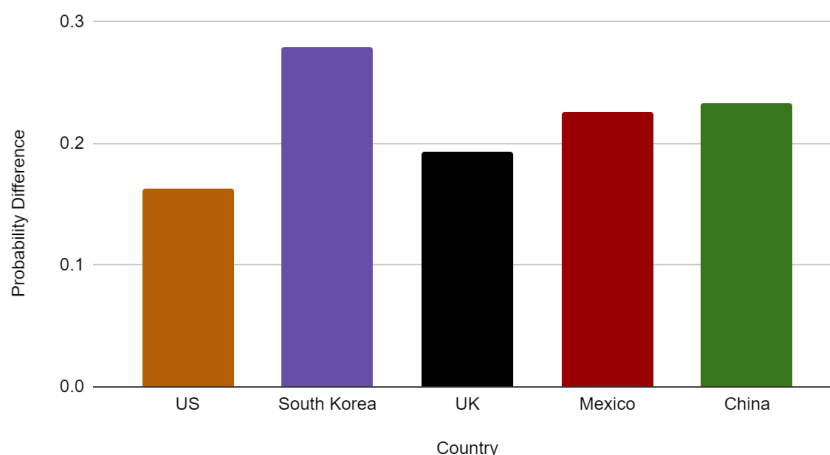
IV. Results

Below are my results with the different temperatures I looked at:

Temperature = 0.95

```
US: 0.16354458371642977
South Korea: 0.27938912849640474
UK: 0.1930003816378303
Mexico: 0.2264174473239109
China: 0.23272115021245554
```

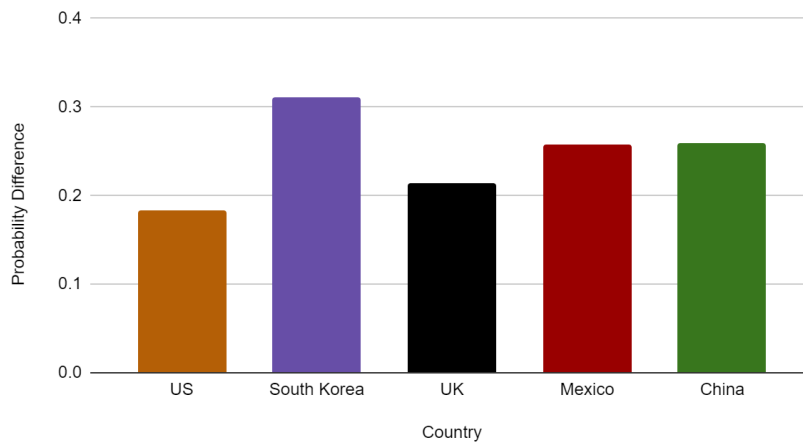
Probability Difference for Temperature = 0.95



Temperature = 0.85

```
US: 0.18271992803784087
South Korea: 0.3103788517182693
UK: 0.21466674684779719
Mexico: 0.257005006831605
China: 0.2588839901727624
```

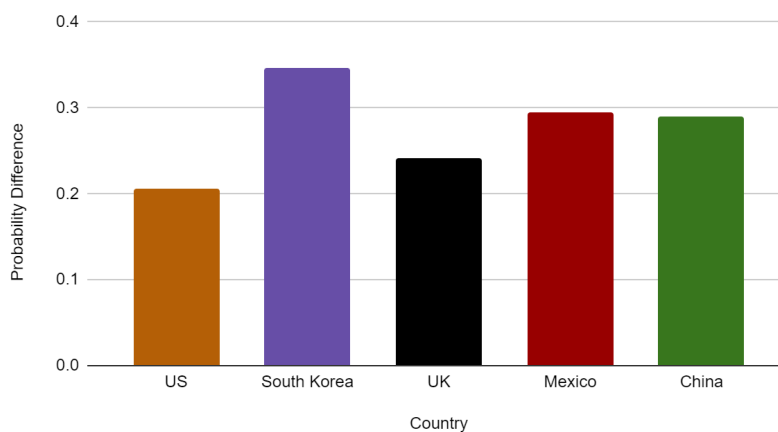
Probability Difference for Temperature = 0.85



Temperature = 0.75

```
US: 0.20627940533449873
South Korea: 0.3470157040574122
UK: 0.24084802376455627
Mexico: 0.29426130309002474
China: 0.29059078931459226
```

Probability Difference for Temperature = 0.75



Even across different temperatures, the GPT-J model showed consistent results. The model was biased towards American culture, validating my initial hypothesis. As we can see, the results clearly show that the U.S. had the least probability distribution different compared to the neutral frame sentences than the other 4 countries I analyzed, showing that the neutral frame sentences predicted words that were similar to the words predicted when the BLANK country was filled in with the U.S. to complete the frame sentences.

In the results above, we can see that the U.S. has a probability difference of 0.16 for a model temperature of 0.95 and when the temperature decreased (less randomness of generated words), that difference increased slightly, to about 0.20 with a temperature of 0.75, but was still the minimum difference compared to the difference of other countries.

It's important to note that not all of the top 5 next likely words that were generated for each variation of the frame sentence were not always words. Frequently, new line characters ('\n') or articles (the, and, etc) were generated, skewing the results' accuracy. However, since I looked at the probability distribution between these words and the neutral words, this wasn't as great of an issue.

V. Conclusions

In Sheng's paper, her research team discussed the potential disadvantage that language models can put marginalized people in due to its inherent bias. According to the results of my project, I was able to confirm that there was, indeed, bias in the GPT-J language model that leaned towards American culture. This was verified by running the model across different temperatures, which are basically numbers between 0-1 that determines the randomness of the sentence completion words that are generated by the language model. Although my results were able to show that there was bias in the GPT-J model in one area of study (television shows and movies), this may be an indicator, along with previous research from teams like Sheng's, that these models can induce potential harm on cultures that aren't part of the westernized, American culture.

VI. References

→ View Archive. "Mapping the Latent Spaces of Culture." The Stone and the Shell, 21 Oct. 2021, tedunderwood.com/2021/10/21/latent-spaces-of-culture.

Blodgett, Su Lin, et al. "Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets." Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, 2021, <https://doi.org/10.18653/v1/2021.acl-long.81>.

Galloway, Stephen. "How Hollywood Conquered the World (All Over Again)." Foreign Policy, 24 Feb. 2012, foreignpolicy.com/2012/02/24/how-hollywood-conquered-the-world-all-over-again

Sheng, Emily, et al. "Societal Biases in Language Generation: Progress and Challenges." Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, 2021, <https://doi.org/10.18653/v1/2021.acl-long.330>.