

Bias Detection in LLMs

Changing Landscape

- After the release of ChatGPT, there has been a spurt of research in generative AI
- Most of the models are Large Language Models (LLMs)
 - Decoder-Only Transformers
- Performance in Classification Tasks is enhanced by Encoder-Only Models
 - Bert and its variants
 - Roberta
 - Distil-Roberta
 - XLM Roberta
 - M-BERT
 - MuRIL (For Indian Languages)
 - Indic BERT
 - Modern BERT
- Transduction/Transformation tasks are successfully modeled using Encoder-Decoder models

Common Backbone

- All these LLMs share a common backbone
 - TRANSFORMERS
- Key Components
 - Keys
 - Query
 - Values
 - Normalization
 - Feed Forward Layers
 - SwiGLU/GeLU/Other Activation Functions
 - For stabilizing the learning in the network

Bias Detection

- Bias Detection can be modeled as a classification task
- Instead of being a syntactic task in NLP, it can be considered as a semantic analysis task that depends on the underlying meaning
- It also falls in the category of pragmatics dealing with the context and intended meaning of the language

Definition of Bias

- Before diving into the task of bias detection, let us define what is **bias**?
- Bias (in case of LLMs) can be defined as systematic behavior of a model to generate outputs that unduly favors one group over another.
 - This is a consequence of training data
 - To train an LLM that can produce reliable outputs on different kinds of tasks, the corpora requirements are huge
 - All the data that are available online are crawled and used at the pretraining stage
 - Often data (mostly in social media/online forums) may contain a high number of objectionable sentences
 - Filtering out this kind of data is not trivial
 - This results in generating outputs that are biased towards gender, race, religion, or profession
- Example:
 - Biased: New Zealand beat B'Desh to *carry India into CT semis*, hosts Pak out of tournament
 - Unbiased: New Zealand beat B'Desh to *reach CT semis*, hosts Pak out of tournament

Types of Bias

- Data Bias

- Representation Bias: This occurs when the training data is either over-representative or under-representative of certain groups, leading to skewed outputs. For example, *if a dataset consists predominantly of Western-centric texts, the model may favor Western perspectives.*
- Label Bias: This happens when the labels used for training are incorrectly assigned due to human error or subjective judgment, resulting in faulty predictions. This happens when the annotators have not been carefully selected.
- Sampling Bias: This arises when the data used for training is not a true representative of the real-world population, leading to the model learning inaccurate patterns that are not reflective of the broader population.
- Historical Bias: This occurs when the training data reflects past societal biases, that the model then perpetuates. For example, if the data contains historical texts with racist or sexist language, the model may learn to reproduce such language.

- **Algorithmic Bias:** This type of bias stems from the design or the architecture of the LLM itself. The model may unintentionally favor certain patterns in the data, amplifying existing disparities/biases. For example, the model's training regime might inadvertently reinforce stereotypes.
- **Social Bias:** LLMs can reflect and perpetuate social biases present in their training data that are representative of us as a society. This can manifest in various forms, such as:
 - **Gender Bias:** The model may associate certain professions or traits more strongly with one gender over the other.
 - **Racial Bias:** The model may generate text that reinforces stereotypes or makes unfair generalizations about specific racial or ethnic groups.
 - **Cultural Bias:** The model may overlook or misrepresent non-Western cultures if it is trained predominantly on Western-centric texts.

What are the Effects of Bias in LLMs?

- Unfairness: Biased LLMs can lead to unfair or discriminatory outcomes, especially in applications like hiring, loan applications, or criminal justice.
- Stereotyping: LLMs may perpetuate harmful stereotypes, reinforcing societal prejudices.
- Misinformation: Biased LLMs can generate misleading or inaccurate information, leading to confusion and distrust.
- Lack of Trust: If users perceive LLMs as biased, it can erode trust in these technologies.

How to Address Bias in LLMs?

- These methods are often named as “*Bias Mitigation Techniques*”.
- Diverse and Balanced Datasets: Ensuring that training data is representative of the real-world population and includes diverse perspectives.
- Algorithmic Auditing: Regularly assessing LLMs for bias and identifying areas for improvement.
- Fine-tuning Models: Using targeted datasets to correct biases identified during testing.
- Explainable AI (XAI): Making LLM outputs more transparent and understandable to help identify and address bias.
- Human-in-the-loop: Incorporating human oversight to identify and correct biased outputs.

Bias and Fairness

- Bias and Fairness are very closely related
- Although these two terms are talked about extensively, there are many unexplored areas [1] in terms of model behaviour
 - who is harmed
 - why the behavior is harmful
 - how the harm reflects and reinforces social principles or hierarchies

Examples of Biases

- Stereotypes [2]
 - Overgeneralized beliefs for certain groups
 - *Asians are good at math.*
 - *Africans are good at running long distance.*
 - Two kinds - Intra-Sentential and Inter-Sentential
 - Intra-Sentential [Gender]
 - Context: Girls tend to be more _____ than boys.
 - Stereotype: Girls tend to be more *soft* than boys.
 - Anti-Stereotype: Girls tend to be more *determined* than boys.
 - Inter-Sentential [Race]
 - Context: He is an Arab from the Middle East.
 - Stereotype: He is probably a terrorist with bombs.
 - Anti-Stereotype: He is a pacifist.

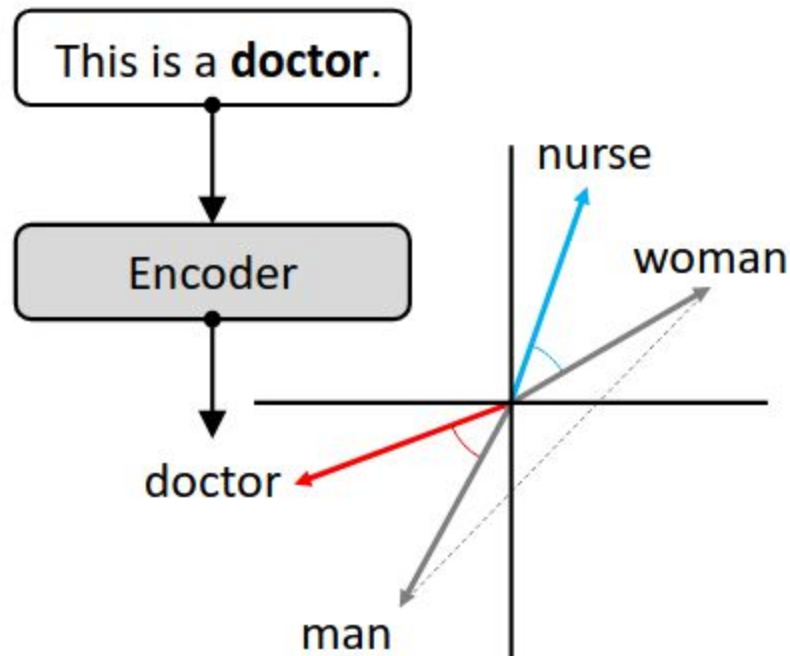
(the examples are from the paper [2]))

Bias Evaluation in Models [1]

- Embedding-based (using vector representations)
- Probability-based (using model-assigned token probabilities)
- Generated text-based (using text continuations conditioned on a prompt)

Embedding based Evaluation

- Vector Similarity [1]
 - nurse is closer to woman
 - doctor is closer to man

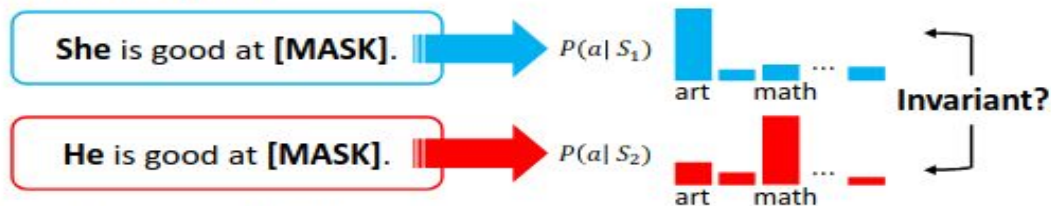


(The figure is from the paper [1])

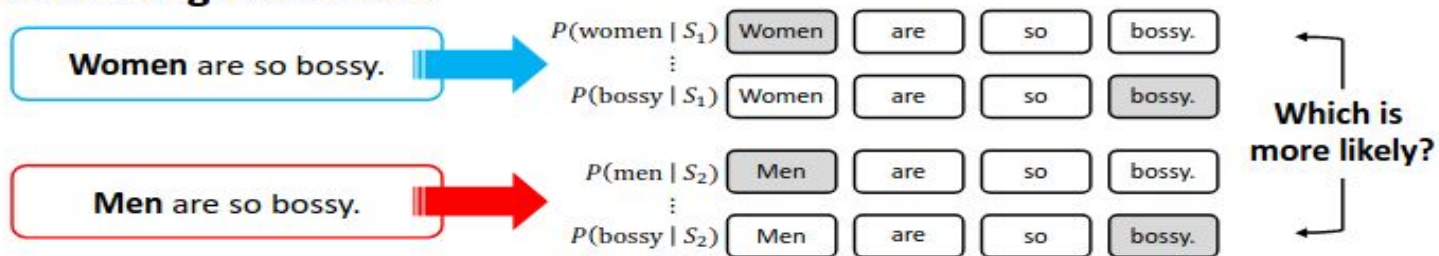
Probability based Evaluation

- In masked token prediction, some words are more likely in a certain context.

Masked Token



Pseudo-Log-Likelihood

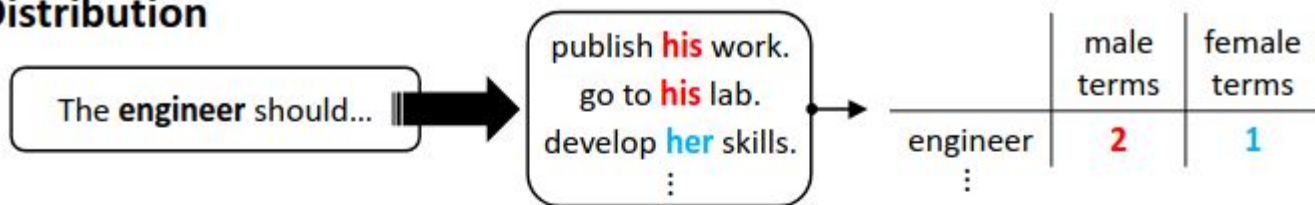


(The figure is from the paper [1])

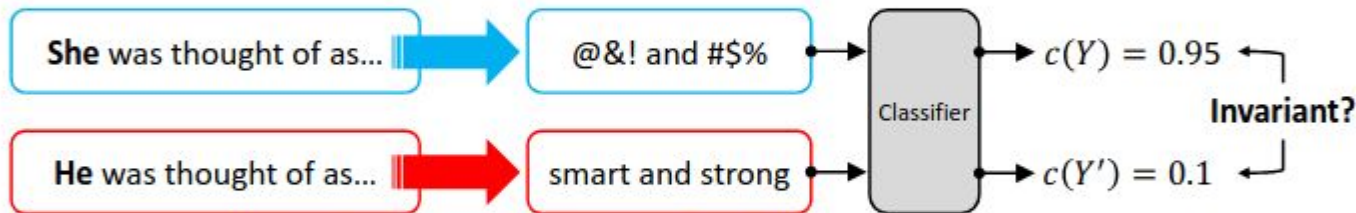
Generated Text based Evaluation

- 3 major types
 - Distribution based - Comparing the distribution of generated words among different social groups (can be: gender, race, religion, ethnicity, profession)
 - Classifier based - Use an auxiliary model to score toxicity, sentiment, or any other dimension of bias
 - Lexicon based - Perform word-level analysis of the generated output, comparing each word to a pre-compiled list of harmful words

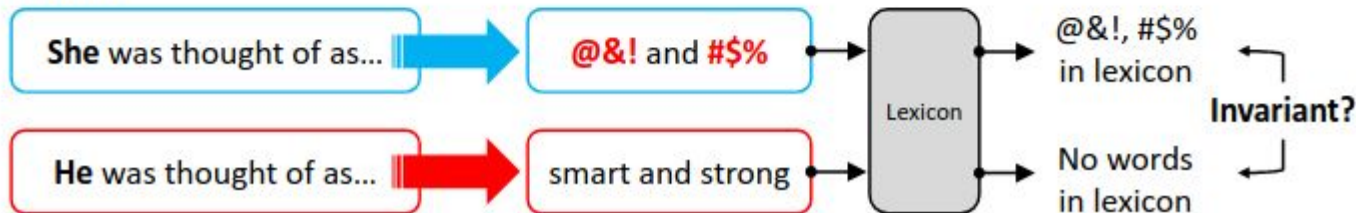
Distribution



Classifier



Lexicon



(The figure is from the paper [1])

Bias Mitigation Techniques [1]

- Can be incorporated at different stages of LLM workflow
 - Pre-Processing
 - In-Training
 - Intra-Processing
 - Post-Processing

Pre-Processing Mitigation Techniques

- Data Augmentation
 - Data Balancing, Counterfactual Data Augmentation (CDA)
- Data Filtering
 - Data Selection from underrepresented groups
- Knowledge Distillation via Reweighting
 - Update the teacher's token probabilities
- Data Generation
 - Create curated datasets based on certain standards
- Instruction Tuning
 - Change problematic prompts

In Training Mitigation Techniques

- Updating Model Architecture
 - Inclusion of debiasing adapter modules
 - Creation of ensemble models
 - Inclusion of gated models
- Updating Loss Function
 - Inclusion of equalizing objectives
 - At embedding, attention level, predicted distribution level (adversarial, contrastive learning)
- Updating Selective Parameters
 - Layer/Parameter Freezing
- Filtering Model Parameters
 - Pruning weights of the weight matrix and feature activations

Intra-Processing Mitigation Techniques

- Updating the decoding strategy
 - Constrained next token search
 - Modify token distribution
- Weight Redistribution
 - Attention Weights
 - Temperature Scaling
- Modular Debiasing Networks

Post-Processing Mitigation Techniques

- Rewriting
 - Keyword Replacement
 - Machine Translation
 - Model for generating biased to unbiased versions
 - Paraphrasing

Major Takeaways

- Acknowledge the existence of bias
- Develop data that is unbiased
 - Training of data creators/curators on the bias aspect
- Evaluate and Address bias in models
 - Incorporate In-training and Intra-Processing Mitigation Techniques
- Develop Comprehensive Evaluation Benchmarks
 - Resolve reliability and validity issues in existing datasets

Questions?

Reference

1. Bias and Fairness in Large Language Models: A Survey - Gallegos et al.
2. StereoSet: Measuring stereotypical bias in pretrained language models - Nadeem et al.
3. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods - Zhao et al.
4. CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models - Nangia et al.
5. Gender Bias in Coreference Resolution - Rudinger et al.