



# Biases in chatbots

Abhilash Chutia(CSM23010)

Bishal Sharma(CSM23024)



# What is Biases in chatbot

Bias in chatbots refers to the presence of unfair or prejudiced outcomes in their responses or decision-making processes. These biases can stem from various sources, including the data used to train the chatbot, the algorithms employed, or the design choices made during development.

# Different types of Biases

- 1.Data Bias
- 2.User Interaction Bias
- 3.Algorithmic Bias
- 4.Presentation and ranking Biases
5. Gender Bias

# Data Bias

Data bias in chatbots occurs when the training data used to develop the chatbot contains biases, leading the chatbot to generate responses that reflect those biases. These biases can arise from various sources:

1. Data collection bias: The way data is collected can introduce biases. For example, if data is collected from a specific population or region, it may not be representative of the larger population.
2. Data representation bias: The way data is represented can also introduce biases. For example, if data is represented in a way that is biased towards a particular group or characteristic, it can affect the model's performance.

3. Historical Bias: It refers to biases present in data or systems that arise from the historical context in which that data was collected or developed. This type of bias reflects the norms, beliefs, and practices that were prevalent during a specific period in history, which may no longer be valid, equitable, or socially acceptable.

4. Selection Bias: The way data is selected or filtered during the training process can introduce bias. If developers only include certain types of data while excluding others, this can skew the chatbot's understanding of various topics.

# User Interaction Bias

In the chatbot development phase, biases mainly originate from developers, training data, and algorithms. When the service is launched, the chatbot interacts with users. It gets prompts and feedback from users. It learns from interactions, which makes it more capable but introduces bias. This is a significant difference between chatbot systems and other traditional ML systems. Users first give the chatbot a prompt to start a conversation. The chatbot will return with a response. Users can grade the response, and ask for regeneration or give a new prompt. In this way, users and the chatbot can exchange information with their biases being propagated mutually.

# Algorithmic Bias

Algorithmic bias results in unfair outcomes due to skewed or limited input data, unfair algorithms, or exclusionary practices during AI development. Algorithmic bias happens due to the biased training data.

Different Types of Algorithmic biases are gender bias, racial bias, age bias, cultural bias etc.

# Gender Bias

Some chatbots may exhibit gender bias by assuming stereotypical gender roles, responding differently to male and female users.

we can classify gender biases into different categories such as:

**Occupational Stereotypes:** Chatbots might suggest certain careers or activities based on gender, such as suggesting nursing for women and engineering for men.

**Naming Bias:** Assuming gender based on first names, which can be incorrect and exclusionary.



# Presentation and Ranking Bias

Since there is too much information on the Internet, it is difficult for chatbots to present all of them. Which information to present and whether the presented information is balanced are determined by the algorithms. The information not presented cannot be received by users. All of these lead to presentation bias. The presented information by chatbots may be ranked or with a certain focus, causing ranking bias.

# Impact of biases in chatbot

Chatbots that perpetuate stereotypes can reinforce harmful societal beliefs and contribute to the marginalization of certain groups.

- Biases can lead to legal and ethical issues, particularly if they result in discrimination or the denial of services.
- Users may lose trust in the technology if they feel they are being treated unfairly or disrespected. Once biases are detected in a chatbot, users may disengage or lose confidence in the system, which can harm an organization's reputation.
- Biases in chatbots can lead to distortion in the information provided, especially in domains where accuracy is critical, such as education or medical advice. This can have serious real-world consequences, as users may act on incorrect or biased information.

# Strategies for mitigating Biases in chatbot

1. Diverse and Representative Training Data: Use datasets that represent diverse populations, cultures, and perspectives.

Ensure balanced representation of different groups in the training data.

Regularly update and expand datasets to include new and diverse information.

2. Bias Detection and Auditing Tools: Design chatbot models with bias-mitigating techniques during the development phase. For example remove sensitive attributes (e.g., gender, race) from training data to avoid bias based on those characteristics.

3. User Feedback and Continuous Learning: Allow users to report biased or inappropriate chatbot behavior, and implement mechanisms to incorporate this feedback into future model improvements.

4. Frequent Model Retraining: Regularly retrain the chatbot model with updated and balanced datasets to ensure it stays relevant and unbiased.

Thank you