

Chapter 17

AI, Ethics, and Hate Speech: A Collaborative Approach to Social Media Influence

Chalamalla Venkateshwarlu

 <https://orcid.org/0000-0003-3559-2475>

Osmania University, India

ABSTRACT

Hate speech is a serious social problem, disrupts social order, impairs individual well-being, and threatens democracy, demonize individuals or groups based on attributes like race, religion, caste ethnicity, gender or sexual orientation, appears in many forms: in spoken language, written materials, symbols and online posts. With the advent of social media, the reach and impact of hate speech has increased exponentially, making it an ever-pressing issue of the digital age. Legislative measures, education awareness campaigns, the increasing prevalence of harmful content online has led to the exploration of artificial intelligence (AI) as the tool to be used to detect and mitigate hate speech online. The paper concludes that preserving ethical concepts in the design and implementation of the AI-driven systems. It points to policies that empower responsible innovation without undermining fundamental democratic values. Through creating inclusive digital spaces and advocating for rights and freedoms for developing socially responsible, inclusivity and respect for individual liberties

DOI: 10.4018/979-8-3693-9904-0.ch017

INTRODUCTION

Artificial intelligence (AI) solutions have also been developed for the rising problem of social media hate speech, which is primarily designed to pick out and filter against such content. Yet the whole matter poses massive ethical concerns when AI is deployed in this way. This paper explores the ethical minefield that AI interventions in hate-speech management create, and highlights considerations such as free speech issues, privacy concerns and second-order effects like self-censorship or chilling effects on open discourse, all of which have ample precedence. Next, it examines the struggles of defining hate speech in general when having to address it from a global platform and introduces a necessity that a specific practice is necessary. And finally, the chapter suggests ethical responsibilities that transparency and accountability should impose on AI developers, platform operators and regulators, for fair play in these systems. It suggests a collaborative effort for ethicists with technologists, lawyers and sociologists to see that the AI solutions are made ethically applicable and socially responsible. As social media and anti-social attacks continue to play a growing role in our lives, there must be some way of drawing the line between the scientifically destructive that needs curbing and fundamental rights, the protections for which are irreplaceable. Significance and Contribution

This chapter will also add significantly to the discussion on what is ethically acceptable use of AI in content moderation. Several AI solutions targeting hate speech have already been developed but can face ethical dilemmas that challenge core human freedoms of freedom of speech and privacy. Identifying these troubling issues, the chapter will provide a nuanced response to how AI can be responsibly used to prevent hate speech on the one hand while safeguarding the rights and freedoms of individuals on the other. The chapter will also present tangible guidance on how improving AI-based systems to bring ethics into future technologies.

The rapid rise of social media, and its impact, has been the greatest challenge in policing dangerous speech, including expressions of hate. It was only a matter of time before the use of more AI-based models to screen and remove automatic hate speech shifted to these platforms that wield the stronghold on public opinion. But AI dispositions in them all bring their own sets of challenging ethical issues, the combination of free speech with privacy, and underlying bias coded into algorithms.

In this chapter, four core ethical challenges that are at stake in the development and application of solutions for AI-driven management of hate speech are focused on balancing technology potential with rights and responsibilities.

Hate speech is any expression, whether spoken, written or symbolic, that denigrates a person or group based on (specific attributes), vilifies them or incites violence against them. These characteristics generally include race, ethnicity, nationality, religion, gender, sexual orientation, disability and/or other status. Utterances aimed

22 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/ai-ethics-and-hate-speech/371746

Related Content

Assessing Factors Affecting the Blockchain Adoption in Public Procurement Delivery in Ghana: A Correlational Study Using UTAUT2 Theoretical Framework

David King Boison, Ebenezer Malcalm, Ahmed Antwi-Boampong, Musah Osumanu Doumbia and Kamal Kant Hiran (2022). *International Journal of Ambient Computing and Intelligence* (pp. 1-13).

www.irma-international.org/article/assessing-factors-affecting-the-blockchain-adoption-in-public-procurement-delivery-in-ghana/314568

Evolutionary Algorithms in Discredibility Detection

Bohumil Sulc and David Klimanek (2009). *Encyclopedia of Artificial Intelligence* (pp. 567-574).

www.irma-international.org/chapter/evolutionary-algorithms-discredibility-detection/10304

Unleashing the Power of Customer Personalization in the Digital Age With Artificial Intelligence

Premendra Sahu and Pinaki Mandal (2024). *Improving Service Quality and Customer Engagement With Marketing Intelligence* (pp. 97-113).

www.irma-international.org/chapter/unleashing-the-power-of-customer-personalization-in-the-digital-age-with-artificial-intelligence/350878

Reasoning Temporally Attributed Spatial Entity Knowledge Towards Qualitative Inference of Geographic Process

Jayanthi Ganapathy and Uma V. (2019). *International Journal of Intelligent Information Technologies* (pp. 32-53).

www.irma-international.org/article/reasoning-temporally-attributed-spatial-entity-knowledge-towards-qualitative-inference-of-geographic-process/225068

Examining the Effect of Implementation Obstacles on Artificial Intelligence Use in the Insurance Sector

Ankit Garg, T. R. Pandey, Laxmi Pandey, Ajay Kumar Varshney and Neha Verma (2026). *AI-Driven Innovations in the Insurance Sector* (pp. 197-226).

www.irma-international.org/chapter/examining-the-effect-of-implementation-obstacles-on-artificial-intelligence-use-in-the-insurance-sector/391888