

Chapter 9

Optimizing Interpretability and Dataset Bias in Modern AI Systems

L. K. Hema

*Department of ECE, Vinayaka Mission.s
Research Foundation (DU), Aarupadai Veedu
Institute of Technology, India*

Rajat Kumar Dwibedi

*Department of ECE, Vinayaka Mission.s
Research Foundation (DU), Aarupadai Veedu
Institute of Technology, India*

Muppala Deepak Varma

*Department of ECE, Vinayaka Mission.s
Research Foundation (DU), Aarupadai Veedu
Institute of Technology, India*

Anamika Reang

*Department of ECE, Vinayaka Mission.s
Research Foundation (DU), Aarupadai Veedu
Institute of Technology, India*

S. Silvia Priscila

*Bharath Institute of Higher Education and
Research, India*

A. Chitra

*Dharmamurthi Rao Bahadur Calavala Cunnan
Chetty's Hindu College, India*

ABSTRACT

As AI systems become deeply ingrained in societal infrastructures, the need to comprehend their decision-making processes and address potential biases becomes increasingly urgent. This chapter takes a critical approach to the issues of interpretability and dataset bias in contemporary AI systems. The authors thoroughly dissect the implications of these issues and their potential impact on end-users. The chapter presents mitigative strategies, informed by extensive research, to build AI systems that are not only fairer but also more transparent, ensuring equitable service for diverse populations. Interpretability and dataset bias are critical aspects of AI systems, particularly in high-stakes applications like healthcare, criminal justice, and finance. In the study, the authors delve deep into the challenges associated with interpreting the decisions made by complex AI models.

INTRODUCTION

The widespread adoption of artificial intelligence (AI) systems across various industries has ushered in a new era of technological advancement, offering promises of increased efficiency, automation, and

DOI: 10.4018/979-8-3693-5951-8.ch009

data-driven decision-making (Le & Viviani, 2018; Ahmed Chhipa et al., 2021). However, as these AI systems become increasingly integrated into our daily lives, there is a growing concern regarding their interpretability and dataset bias (Bose et al., 2023). These issues are of paramount importance, especially in the context of OneWebbie's focus on catering to diverse client interests and demographics in the unique market landscape of Australia (Angeline et al., 2023; Saxena & Chaudhary, 2023).

Interpretability, in the area of AI, refers to the capacity to comprehend and trust the decisions made by AI systems (Chakrabarti & Goswami, 2008). It is, without a doubt, one of the fundamental prerequisites for the acceptance and widespread adoption of AI technologies. When AI systems generate decisions or recommendations, users and stakeholders need to understand how those conclusions were reached (Senapati & Rawal, 2023a). In an agency like OneWebbie, which specializes in web development, digital marketing, SEO, software and app development, e-commerce, content creation, social media, and a wide array of services, ensuring interpretability is paramount (Senapati et al., 2024). Clients must have confidence in the AI-driven strategies and solutions being employed on their behalf (Senapati & Rawal, 2023b).

However, achieving interpretability in AI is not always straightforward. Many AI models, particularly deep learning models, are often considered "black boxes" because they operate on complex mathematical computations that are difficult for humans to interpret (Cristian Laverde Albarracín et al., 2023). For OneWebbie, which offers services such as PPC, SEO, content creation, and social media management, understanding how AI algorithms make decisions is vital to tailoring strategies effectively (Sharma et al., 2021). This requires the development of interpretable AI models and the implementation of transparent processes that demystify the decision-making process (Vignesh Raja et al., 2023). In the Australian market, where clients come from diverse backgrounds and industries, ensuring that AI-driven solutions are understandable and trustworthy becomes even more critical (Shah et al., 2020).

Dataset bias, on the other hand, is another pressing concern in the AI landscape. It arises when the data used to train AI models is not representative of the real-world scenarios the AI system will encounter (Jasper et al., 2023). This bias can lead to unfair or inaccurate predictions, which could have detrimental consequences for OneWebbie's clients (Haro-Sosa & Venkatesan, 2023). In a country as culturally diverse as Australia, where the audience target location is specified, the importance of dataset diversity cannot be overstated (Sharma et al., 2022). Ensuring that the data used to train AI models reflects the nuances and diversity of the Australian market is essential for delivering accurate and fair results (Gaayathri et al., 2023).

Addressing dataset bias requires meticulous attention to data collection, curation, and ongoing monitoring (Jeba et al., 2023). OneWebbie's focus on services like ORM (Online Reputation Management), CRM (Customer Relationship Management), and social media engagement means that the data used to train AI systems often comes from various sources, including social media, customer interactions, and website traffic. It's crucial to identify and mitigate biases present in these datasets to avoid reinforcing existing stereotypes or making decisions that unintentionally discriminate against certain demographics (Karn et al., 2022a; Sivapriya et al., 2023). In the Australian market, understanding and addressing dataset bias is a multifaceted challenge. Australia is known for its cultural diversity, with a wide range of languages, cultures, and social contexts (Kumar et al., 2023). OneWebbie's commitment to tailoring services to diverse client interests and demographics in this context means that dataset bias can manifest in many subtle ways. It's not just about avoiding explicit bias but also about ensuring that AI systems are sensitive to the cultural nuances and diversity of the Australian audience (Karn et al., 2022b).

As OneWebbie continues to serve its diverse clientele in Australia across a wide spectrum of services, understanding and addressing the challenges of AI interpretability and dataset bias is paramount. Achiev-

17 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/optimizing-interpretability-and-dataset-bias-in-modern-ai-systems/349525

Related Content

Impact of Artificial Intelligence on Marketing Research: Challenges and Ethical Considerations

Laura Sáez-Ortuño, Javier Sanchez-Garcia, Santiago Forgas-Coll, Rubén Huertas-García and Eloi Puertas-Prat (2023). *Philosophy of Artificial Intelligence and Its Place in Society* (pp. 18-42).

www.irma-international.org/chapter/impact-of-artificial-intelligence-on-marketing-research/332598

Extracting Functional Dependencies in Large Datasets Using MapReduce Model

K. Amshakala, R. Nedunchezian and M. Rajalakshmi (2014). *International Journal of Intelligent Information Technologies* (pp. 19-35).

www.irma-international.org/article/extracting-functional-dependencies-in-large-datasets-using-mapreduce-model/116741

Lightweight ConvNet Model for American Sign Language Hand Gesture Recognition

Shamik Tiwari (2022). *Challenges and Applications for Hand Gesture Recognition* (pp. 175-193).

www.irma-international.org/chapter/lightweight-convnet-model-for-american-sign-language-hand-gesture-recognition/301062

The Concept of Exaptation Between Biology and Semiotics

Davide Weible (2012). *International Journal of Signs and Semiotic Systems* (pp. 72-87).

www.irma-international.org/article/concept-exaptation-between-biology-semiotics/64639

Bridge Crack Recognition Method Based on Yolov5 Neural Network Fused With Attention Mechanism

Yingjun Wu, Junfeng Shi, Wenxue Ma and Bin Liu (2024). *International Journal of Intelligent Information Technologies* (pp. 1-25).

www.irma-international.org/article/bridge-crack-recognition-method-based-on-yolov5-neural-network-fused-with-attention-mechanism/361575