# Unveiling Bias: Analyzing Artificial Intelligence and Machine Learning's Impact on Fairness in the Criminal Justice System

## Sukanya Konatam[1],  Venkata Naga Murali Konatam[2]

[1]Senior Manager of Enterprise Data Governance and Data Science, IT Department, Vialto Partners, Texas, USA
[2]Data Architect, IT Department, Capital One, Texas, USA

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract -** *This paper examines the integration of artificial intelligence (AI) and machine learning (ML) in the criminal justice system, highlighting both the potential benefits and significant concerns related to bias. AI and ML are being employed to enhance various aspects of criminal justice, including crime prediction, tracking, and judicial decision-making. However, these technologies are susceptible to biases inherent in the historical data they rely on, which can perpetuate and amplify existing disparities within the justice system. The article delves into the historical context of AI advancements in criminal law, from early digitization efforts to the current deployment of sophisticated AI*

*applications. It explores notable AI and ML technologies used in criminal justice, such as risk assessment tools, predictive policing algorithms, and facial recognition systems. The discussion emphasizes the ethical implications of AI bias, particularly its impact on marginalized communities. To address these issues, the article proposes various strategies for detecting and mitigating biases in AI/ML systems, including bias detection tools, data pre-processing techniques, and the importance of transparency and accountability. By scrutinizing these technologies and their applications, the article aims to contribute to the development of fairer and more equitable AI systems in criminal justice.*

*Key Words***:  Bias in Judicial AI, Ethical AI in Criminal Law, Bias in Artificial Intelligence, Machine Learning Accountability, ML Bias in Criminal Justice

## 1.INTRODUCTION - ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING IN THE CRIMINAL JUSTICE SYSTEM

Artificial intelligence (AI) and machine learning (ML) algorithms are increasingly being adopted within the criminal justice system to tackle various challenges, such as predicting and tracking crimes and criminals and assisting in criminal court proceedings. [1] These technologies hold the potential to enhance efficiency and effectiveness in crime prevention and investigation efforts. However, the use of AI and ML in criminal justice also raises significant concerns regarding potential biases and privacy infringements. These models are built on historical data, which often reflect and perpetuate existing societal biases. [1] Biases in AI and ML systems can influence various stages of the criminal justice process, from arrest and bail decisions to sentencing and

parole. [1] Addressing these biases is crucial to ensure fairness and justice in AI and ML applications in criminal justice. This involves scrutinizing the data used to train these models, implementing robust bias detection and mitigation techniques, and ensuring transparency and accountability in AI decision-making processes. [1] As we continue to integrate AI and ML into the criminal justice system, it is essential to remain vigilant about the potential for bias and to develop strategies to counteract it, ensuring these technologies contribute positively to the criminal justice system rather than reinforcing existing disparities. [1]

## 2. EVOLUTION AND CURRENT APPLICATIONS OF ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING IN CRIMINAL JUSTICE

### A.     History of AI Advancements in Criminal Law

The integration of Artificial Intelligence (AI) and Machine Learning (ML) into the criminal justice system has been a gradual process, marked by key advancements over the past few decades. [3] In the early 2000s, the digitization of court filings and processes set the stage for the initial adoption of AI-powered tools to assist in partial criminal court decisions. [4] This was followed by the emergence of online dispute resolution (ODR) as an alternative to traditional in-person court proceedings, further showcasing the potential of AI to transform the judicial landscape. [4] The 2010s saw a significant surge in the development and deployment of more sophisticated AI applications within the criminal justice domain. [3] Algorithmic risk assessment tools were introduced to inform bail, sentencing, and parole decisions, leveraging predictive analytics to assist human decision-makers. [4] Jurisdictions such as the United States and the European Union began exploring the use of AI-powered digital tools to streamline various legal processes and support the work of judges and law enforcement. [4] More recently, in the late 2010s and early 2020s, the legal sector has witnessed the emergence of AI-driven solutions that can potentially automate certain adjudicative functions, raising complex ethical and legal questions about the role of technology in judicial decision-making. [2] As these advancements continue to unfold, the criminal justice system is grappling with the challenges and opportunities presented by the integration of AI and ML, seeking to harness the benefits while ensuring the protection of fundamental rights and the rule of law. [5]

(1)

**Chart -1**:This image represents the flow of online dispute resolution (ODR) as a transformative alternative to traditional court proceedings. ODR leverages technology to facilitate judicial processes, allowing for efficient, accessible, and often more cost-effective resolutions. The integration of AI within ODR platforms exemplifies the ongoing evolution of the judicial landscape, demonstrating how digital tools can enhance the efficiency and accessibility of justice in the mid to late 2000s and beyond.

*B.* **Notable AI/ML Technologies in the System of Criminal Justice**

Several AI and machine learning algorithms and models have been integrated into the criminal justice system to aid various processes. One prominent example is the use of risk assessment tools, such as the Correctional Offender Management Profiling for Alternative Solutions (COMPAS) and the Public Safety Assessment (PSA), which leverage algorithms to predict an individual's potential for future misconduct, informing crucial decisions like pretrial incarceration. [6] These tools assess factors like age, criminal history, and past misconduct to generate risk scores that judges utilize to set release conditions. [6] Additionally, predictive policing algorithms analyze historical crime data to identify areas and times where crimes are most likely to occur, enabling law enforcement agencies to allocate resources more effectively.

[7] Facial recognition technology is another AI application in the criminal justice sector, used to assist intelligence analysts in establishing an individual's identity and whereabouts by automating the analysis of large volumes of visual data. [8] These AI-driven tools promise enhanced consistency, accuracy, and transparency in judicial decision-making and law enforcement operations. [6] [7] However, concerns have been raised about the potential for bias, lack of individualization, and issues of transparency in these AI systems. [6]

## 3. Artificial Intelligence Biases in THE CRIMINAL Justice Industry

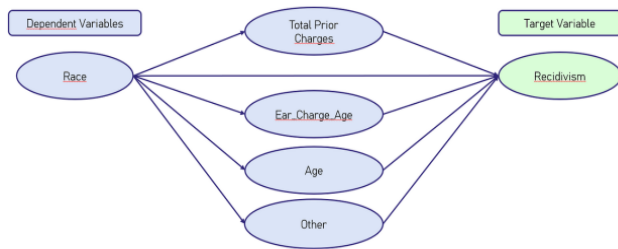The integration of AI and machine learning technologies into the criminal justice system has raised significant concerns regarding the potential for bias and discrimination. [9] These advanced algorithmic systems are increasingly being employed for critical decision-making processes, such as risk assessment, predictive policing, and sentencing recommendations. [9] [10] However, the reliance on historical data, which often reflects the biases inherent in the criminal justice system, can lead to the perpetuation and amplification of these biases within the algorithms. [9] [10] The lack of transparency and accountability surrounding the development and deployment of these AI systems further exacerbates the risk of biased outcomes, undermining the principles of fairness and justice. [10]

**A. Popular Types of AI Biases found in Criminal Law Industry**

• **Racial Bias**- Racial Bias in the Criminal Justice System Racial bias is a prevalent issue that permeates the criminal justice system in the United States. Studies have shown that African Americans are disproportionately more likely to be arrested, convicted, and given harsher sentences compared to their white counterparts, even for similar crimes. [12] This disparity can be attributed to factors such as discriminatory policing practices, biased decision-making by prosecutors and judges, and the legacy of systemic racism. [12] For example, the COMPAS algorithm, widely used in the U.S. criminal justice system, was shown to falsely flag Black defendants as being at a higher risk of recidivism compared to their white counterparts. [13] Such biases have the potential to perpetuate structural inequalities and exacerbate existing racial disparities within the criminal justice system. These findings highlight the deep-rooted racial biases that continue to plague the criminal justice system, with severe consequences for individuals and communities of color.

• **Socioeconomic Bias- Socioeconomic** Bias in the Criminal Justice System Socioeconomic bias is another pervasive issue in the criminal justice system, where individuals from lower socioeconomic backgrounds are more likely to face harsher treatment and outcomes compared to those from more affluent backgrounds. [11] [12] This bias manifests in various ways, such as the inability to afford adequate legal representation, the increased likelihood of being detained pre-trial due to an inability to post bail, and the imposition of fines and fees that perpetuate the cycle of poverty. [11] [12] For Instance, research on The Public Safety Assessment (PSA) has shown that even ostensibly "objective" factors like prior arrests and convictions can reflect underlying biases and inequities in policing and prosecution, which disproportionately target individuals from disadvantaged socioeconomic backgrounds. [15] As a result, the PSA's risk assessments may still perpetuate socioeconomic disparities, with defendants from poorer neighborhoods or lower income levels more likely to be labeled as high-risk, even when controlling for their actual criminal histories. [14] This

demonstrates how the use of historical data in AI/ML criminal justice tools can entrench systemic inequities, underscoring the importance of carefully examining the data sources and design of these algorithms. These socioeconomic disparities within the criminal justice system further exacerbate the challenges faced by marginalized communities, perpetuating a cycle of injustice and inequality.
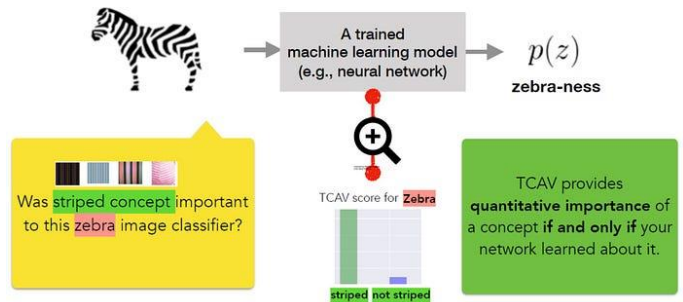


(2)

**Chart -2**:This diagram illustrates the dependent variables used by the COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) algorithm to predict the likelihood of recidivism. The variables include race, age, total prior charges, age at first charge, and other relevant factors. The model shows how these variables influence the target variable, recidivism. Studies and critiques of the COMPAS algorithm highlight concerns about biases, particularly the disproportionate impact of race on the predicted outcomes, raising questions about fairness and accuracy in the criminal justice system.

### B. Ethical implications of bias in Criminal Justice Artificial Intelligence & Machine Learning

The use of AI and machine learning (ML) systems in the criminal justice system raises significant ethical concerns regarding bias and fairness. [16] These algorithms are only as unbiased as the data used to train them, and if the training data contains inherent biases, the algorithms can perpetuate or amplify those biases, leading to discriminatory outcomes for marginalized groups. [16] This can result in unfair sentencing, reinforce systemic racism, and perpetuate unjust practices within the criminal justice system. [17] Furthermore, the complexity of these algorithms can make it challenging to explain how decisions are reached, undermining transparency and accountability. [16] Protecting the privacy and security of the personal data used by these systems is also crucial to maintain public trust and prevent misuse. [16] Addressing the ethical implications of AI/ML in criminal justice requires prioritizing fairness, transparency, and accountability in the development and deployment of these technologies. [17]



(3)

**Chart -3**:This diagram demonstrates the application of TCAV (Testing with Concept Activation Vectors) to determine the significance of specific concepts, such as "striped" for recognizing zebras, within machine learning models. TCAV allows researchers to quantitatively assess the importance of a concept, revealing whether the model has learned about it. Such tools are crucial in the criminal justice sector to detect and mitigate biases in AI/ML systems. By visualizing and understanding the internal workings of these models, stakeholders can ensure fair and accurate predictions, promoting justice and equity in criminal justice systems.

### C. Strategies for Detecting Biases in the System of Criminal Justice AI/ML.

Various approaches and techniques are employed to identify and detect biases in AI/ML systems within the criminal justice sector. These include the use of tools like TCAV and What-If Tool to visualize the internal representations of machine learning models and explore how different inputs affect the model's predictions. [18] Additionally, techniques such as data pre-processing, fairness-aware machine learning algorithms, and bias-correction algorithms are utilized to mitigate biases. [18] However, effectively detecting biases in these systems remains challenging due to issues like biased training data, lack of transparency in AI decision-making, and the potential for the AI systems themselves to perpetuate and amplify biases. [18]

### D. D. Detection of Famous Forms of Biases found in Criminal Justice AI

- **Discovering Racial Bias**- Racial bias in the criminal justice sector is perpetuated by algorithms that replicate or amplify historical biases, leading to discriminatory practices in areas such as predictive policing and sentencing. [19] To detect racial bias, researchers have employed methods like algorithmic bias detection and mitigation, implicit bias testing, and the examination of racial and ethnic disparities in discretionary criminal justice decisions. [19] [20] For example, a landmark investigation by ProPublica in 2016 revealed the racial biases present in an algorithmic risk assessment tool used in criminal justice systems across the United States. [23]

The researchers obtained data on risk scores assigned to over 10,000 defendants in Broward County, Florida, and tracked which of those individuals were charged with new crimes within two years. [23] Through detailed analysis, they found that Black defendants were twice as likely as white defendants to be incorrectly labeled as higher risk, even when controlling for factors like criminal history, age, and gender. [23] ProPublica employed a range of analytical techniques, including subgroup comparisons, logistic regression modeling, and visualization methods, to rigorously detect and quantify the algorithmic bias. [22] This groundbreaking investigation demonstrated the critical importance of comprehensive, data-driven evaluations to uncover and address racial disparities embedded within criminal justice technologies. In summary, ProPublica's work highlighted the need for robust bias detection approaches to ensure the fairness and accountability of AI/ML systems deployed in high-stakes domains like criminal justice. One key approach to detecting racial bias is to scrutinize the training data used to develop risk assessment algorithms, as biases in historical crime data can lead to algorithms disproportionately flagging minority groups as high-risk. [19] [20] Techniques such as undersampling crime data involving certain demographics and oversampling data related to others can help address this issue. [20] Additionally, ensuring transparency around the algorithms' inner workings and enabling independent audits are crucial for identifying and mitigating racial biases. [20] Meaningful community engagement in the design and deployment of these tools is also important to build trust and accountability. [20] [21]

- **Identifying Socioeconomic Bias- Socioeconomic** bias in the criminal justice sector refers to the unfair advantage or disadvantage faced by individuals based on their socioeconomic status, which can be perpetuated by AI/ML systems. [19] [20] To detect socioeconomic bias, researchers have explored methods like processing training data to achieve balanced representation, adjusting algorithms to assign less weight to data points related to certain demographics, and implementing continuous monitoring and evaluation of model outputs. [20] For example, A study by researchers at the University of Chicago developed a crime prediction algorithm that revealed significant biases in police enforcement. [24] The study found that while the algorithm was able to predict crimes one week in advance with 90% accuracy, it also highlighted disparities in police response. [24] Crimes in wealthier areas resulted in more arrests, while arrests in disadvantaged neighborhoods dropped, suggesting inherent bias in how law enforcement resources were allocated and utilized. [24] This investigation demonstrates how advanced analytical techniques and tools can be employed to uncover systemic biases in the criminal justice system, particularly around socioeconomic factors that disproportionately impact marginalized communities. [25]

Ultimately, these findings underscore the critical need to address the societal biases that become embedded within the data and algorithms used in predictive policing applications. The lack of transparency from companies developing risk assessment tools is a significant challenge in detecting socioeconomic bias. [20] Algorithmic hygiene, which involves surfacing and responding to algorithmic bias upfront, as well as proactive addressing of factors contributing to bias, is crucial for detecting and mitigating socioeconomic bias in the criminal justice sector. [20]

## 4. Conclusion

The integration of artificial intelligence (AI) and machine learning (ML) into the criminal justice system represents a transformative shift with the potential to enhance efficiency and effectiveness across various processes, from crime prediction to judicial decision-making. Despite the promising advancements, the deployment of these technologies has exposed critical concerns about bias and fairness. AI and ML systems, often reliant on historical data, can perpetuate and even exacerbate existing inequalities within the justice system. The biases embedded in these models pose significant challenges, influencing outcomes in arrest decisions, risk assessments, and sentencing, which can undermine the principles of justice and equity.

This article has explored the evolution of AI and ML in criminal justice, identifying key technologies and their applications. It has also highlighted the ethical implications of AI bias, emphasizing the need for vigilance and proactive measures to address these issues. Effective strategies for detecting and mitigating bias are crucial to ensure that AI and ML systems do not reinforce systemic disparities but instead contribute to a more equitable and just legal framework.

One promising approach to addressing these challenges is the use of Large Language Models (LLMs). LLMs can enhance transparency and accountability in AI decision-making by providing more nuanced and interpretable explanations of algorithmic processes and outputs. They can be employed to analyze and interpret complex data, identify patterns of bias, and suggest adjustments to algorithms. Additionally, LLMs can facilitate the development of more inclusive and representative datasets, improving the fairness of AI systems.

As the criminal justice system continues to integrate AI and ML technologies, it is imperative that stakeholders prioritize transparency, accountability, and fairness. By leveraging LLMs and other advanced tools, and fostering ongoing dialogue about the ethical implications of AI, we can work towards harnessing the benefits of these technologies while safeguarding the fundamental principles of justice. Future research and policy development should focus on refining these technologies and establishing robust frameworks to ensure their responsible and equitable use in the criminal justice domain.

## References

1.  AI Applications in the Criminal Justice System: The Next Logical Step or Violation of Human Rights, https://jolets.org/ojs/index.php/jolets/article/download/124/67

2.  Justice Augmented: Navigating the Ethical and Legal Terrains of AI Integration in International Criminal Proceedings, https://www.dmejournals.com/index.php/DMEJL/article/download/300/147

3.  Algorithmic Decision Making: Can Artificial Intelligence and the Metaverse Provide Technological Solutions to Modernise the United Kingdom's Legal Services and Criminal Justice?, https://lifescienceglobal.com/pms/index.php/FIA/article/download/9582/5010

4.  Exploring the impact of reporting medium on online crime reporting experiences: comparing live chat with human and AI operators, https://www.semanticscholar.org/paper/47cae47d70e91465409d35eaec6c7e85c1457abc

5.  Technical Solutions for Legal Challenges: Equality of Arms in Criminal Proceedings, https://www.semanticscholar.org/paper/ca6744ca37b940fa742ee6d40130f6f9c6548b4d

6.  AI in Criminal Justice: Shaping the Future of Law and Order - Plat.AI, https://plat.ai/blog/ai-in-criminal-justice/

7.  Predictive Policing using Machine Learning (With Examples), https://www.cogentinfo.com/resources/predictive-policing-using-machine-learning-with-examples

8.  Using Artificial Intelligence to Address Criminal Justice Needs, https://nij.ojp.gov/topics/articles/using-artificial-intelligence-address-criminal-justice-needs

9.  AI in Criminal Justice: How it Can Become Biased - Inspirit AI, https://www.inspiritai.com/blogs/ai-student-blog/ai-in-criminal-justice-how-it-can-become-biased

10. Gender and Sexuality Bias in the Application of Artificial Intelligence ..., https://www.linkedin.com/pulse/gender-sexuality-bias-application-artificial-systems-justice-angiey-6fouf

11. 5 steps to help eliminate socio-economic bias, https://www.americanbar.org/news/abanews/publications/youraba/2019/march-2019/5-steps-to-help-eliminate-socio-economic-bias/

12. Report to the United Nations on Racial Disparities in the U.S. ..., https://www.sentencingproject.org/reports/report-to-the-united-nations-on-racial-disparities-in-the-u-s-criminal-justice-system/

13. Fairness Deconstructed: A Sociotechnical View of 'Fair' Algorithms in Criminal Justice, https://www.semanticscholar.org/paper/525a9e9688eb734e6c1bea2afe60eb367f79b77c

14. Artificial Intelligence (AI) & Criminal Justice System - PixelPlex, https://pixelplex.io/blog/artificial-intelligence-criminal-justice-system/

15. AI in the Criminal Justice System - Epic.org, https://epic.org/issues/ai/ai-in-the-criminal-justice-system/

16. The Ethical Implications of AI in Criminal Justice - Alexi, https://www.alexi.com/post/the-ethical-implications-of-ai-in-criminal-justice

17. The Ethical Implications of AI-Powered Criminal Justice - LinkedIn, https://www.linkedin.com/pulse/ethical-implications-ai-powered-criminal-justice-can-biased-pathak

18. Navigating the AI Bias: Exploring Tools and Techniques, https://arunapattam.medium.com/navigating-the-ai-bias-exploring-tools-and-techniques-c42b0f26fd29

19. Algorithmic bias detection and mitigation: Best practices and policies ..., https://www.brookings.edu/articles/algorithmic-bias-detection-and-mitigation-best-practices-and-policies-to-reduce-consumer-harms/

20. Fairness and Bias in Artificial Intelligence: A Brief Survey of Sources ..., https://www.mdpi.com/2413-4155/6/1/3

21. Racial bias in AI: unpacking the consequences in criminal justice ..., https://www.irissd.org/post/racial-bias-in-ai-unpacking-the-consequences-in-criminal-justice-systems

22. Bias in Criminal Risk Scores Is Mathematically Inevitable ..., https://www.propublica.org/article/bias-in-criminal-risk-scores-is-mathematically-inevitable-researchers-say

23. ProPublica Responds to Company's Critique of Machine Bias Story, https://www.propublica.org/article/propublica-responds-to-companys-critique-of-machine-bias-story

24. [PDF] Predicting Police Misconduct - The University of Chicago, https://bfi.uchicago.edu/wp-content/uploads/2024/05/BFI_WP_2024-62.pdf

25. Predictive Policing using Machine Learning (With Examples), https://www.cogentinfo.com/resources/predictive-policing-using-machine-learning-with-examples

26. (1) **[Chart -1]** *Online Dispute Resolution Diagram. Making Justice Effortless and Accessible Online Dispute Resolution*, Imagesoft, https://imagesoftinc.com/courts/online-dispute-resolution/.

27. (2) **[Chart -2]** Eger, Felix. *COMPAS Directed graph displaying potential causal effects of the features.* 20 Apr. 2022. *Fairness in American Courts: An Exploration of the COMPAS Algorithm*, Medium, https://medium.com/@felix.eger17/fairness-in-american-courts-an-exploration-of-the-compas-algorithm-415a23affb39. Accessed 22 June 2024.

28. (3) **[Chart -3]** Rahimi, Reza. TCAV Tool Diagram. 29 July 2022. Google AI Open-Sourced a New ML Tool for Conceptual and Subjective Queries over Images, Infoq, https://www.infoq.com/news/2022/07/google-ai-tcav/. Accessed 22 July 2024.

## BIOGRAPHIES



Sukanya Konatam is the principal author of the research, leading the idea of writing a comprehensive literature survey journal on AI and ML biases. With over 18+ years of experience in the IT industry, enriched by five years dedicated to specializing in AI and AI Governance. Her expertise spans a comprehensive range of IT disciplines, underscored by a profound depth of knowledge in the governance of artificial intelligence. This unique combination of experience positions her as a leading figure in the field, adept at navigating the intricacies of AI technology with a strategic and ethical approach. She has implemented data centric solutions for several industries including banking & financials, telecom, health care, automobile, criminal justice and many more.

Her significant contributions to the paper included an in-depth analysis of biases in the Criminal Justice system, as well as authoring the abstract, conclusion, and ensuring overall grammatical accuracy. Sukanya's robust background in data governance, data warehousing, machine learning, and AI, combined with her proficiency in AI Governance made her instrumental in identifying and proposing mitigation strategies for AI biases. Additionally, she holds a postgraduate degree in Data Science, Machine Learning, and Artificial Intelligence.



Venkat Konatam, a secondary author of this literature survey, played a key role in investigating artificial intelligence and machine learning biases within the

Criminal Justice system. With over 25 years of experience in IT, Venkat has excelled in constructing resilient data platforms and delivering impactful data products. His expertise lies in data architecture, database design, and performance optimization, with a strong focus on AI and ML applications.

Venkat is currently pursuing a postgraduate program in data science, machine learning, and artificial intelligence.