

HAIID Exam assignment

Task 1: Paper review

For the first task I've chosen paper №9: Angwin et al.: Machine bias: There's software used across the country to predict future criminals, and it's biased against blacks. (ProPublica, 2016).

Motivation

The ProPublica article, "Machine Bias," authored by Julia Angwin and her colleagues, serves as a critical investigation into the deployment of predictive algorithms, specifically focusing on the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) system. Published in 2016, the motivation behind the article lies in uncovering potential biases within this widely-used software employed in the criminal justice system across the United States.

The primary motivation stems from the broader societal concern of fairness and equity in the criminal justice system, particularly in the context of pretrial risk assessments. COMPAS, designed to assist judges in evaluating the likelihood of an individual reoffending, raised questions about its impact on marginalized communities, specifically the alleged bias against black defendants.

In the context of the article, the authors aim to shed light on the potential racial disparities perpetuated by COMPAS. By scrutinizing the algorithm's predictions and their real-world consequences, the motivation is rooted in addressing the urgent need for transparency, accountability, and fairness in the use of technology within the criminal justice system.

The users targeted by this investigative piece are multifaceted, encompassing not only legal professionals such as judges and lawyers but also the general public. As algorithms increasingly play a role in shaping legal outcomes, it becomes imperative for those directly involved in the legal process and the public at large to understand the potential biases embedded in these systems. The users' need, in this case, is to comprehend the implications of algorithmic decision-making on fairness and justice, particularly with regard to racial disparities in the criminal justice system.

The research questions addressed in the article revolve around the reliability and fairness of the COMPAS system. The authors explore whether the algorithm exhibits bias in predicting recidivism, especially concerning racial groups. This question is of paramount importance due to the potential real-world impact on the lives of individuals, as algorithmic predictions can influence judicial decisions, potentially leading to unequal treatment based on race.

The significance of these research questions lies in the broader societal implications of algorithmic bias in the criminal justice system. If indeed COMPAS demonstrates bias, it raises concerns about the perpetuation and exacerbation of existing racial disparities in legal outcomes. Such biases could have far-reaching consequences, impacting not only individual lives but also perpetuating systemic inequalities within the criminal justice system.

Novelty and significance

A novel aspect of the paper is its targeted focus on predictive policing algorithms, addressing a burgeoning issue in law enforcement where automated systems are increasingly influencing critical decisions. By scrutinizing the COMPAS algorithm, a widely adopted tool for risk

assessment in pretrial and sentencing decisions, the study provided a concrete example of how such technologies can potentially perpetuate and exacerbate existing societal inequalities. In terms of methodology, Angwin et al. undertook a comprehensive approach combining quantitative analysis and investigative journalism techniques. The dataset comprised over 7,000 individuals in Broward County, Florida, subjected to the COMPAS algorithm. The researchers employed statistical analyses to assess the algorithm's fairness and accuracy, providing a robust foundation for their assertions. The inclusion of real-world narratives and specific cases added a qualitative dimension, offering a human perspective on the consequences of biased algorithmic predictions.

The results were striking, revealing a significant racial bias in the COMPAS algorithm. The erroneous classification of black defendants as high-risk, compared to their white counterparts with similar backgrounds, highlighted the potential for unjust outcomes within the criminal justice system. This disparity, as evidenced by a higher false positive rate for black individuals, underscored the urgent need for scrutiny and reevaluation of algorithmic decision-making processes.

The contribution of the paper is multifaceted. Firstly, it prompted a heightened awareness of the ethical dimensions surrounding the use of algorithms in the criminal justice system. The findings served as a catalyst for broader discussions on transparency, accountability, and the societal impact of relying on machine learning in legal decision-making. Secondly, the study underscored the importance of nuanced investigations into algorithmic biases, emphasizing the need for a multidimensional approach that combines quantitative rigor with real-world narratives.

Limitations

While Angwin et al.'s paper, "Machine Bias: There's software used across the country to predict future criminals, and it's biased against blacks," made significant strides in uncovering biases in the COMPAS algorithm and sparking crucial conversations on algorithmic fairness, it is essential to acknowledge some limitations and areas where further exploration might be warranted.

Single Algorithm Focus

The paper primarily centers around the COMPAS algorithm, providing a deep dive into its biases. However, this exclusive focus raises questions about the generalizability of the findings to other predictive policing algorithms. Different algorithms might have distinct mechanisms and sources of bias that merit separate investigations.

Limited Geographical Scope

The study concentrated on data from Broward County, Florida, potentially limiting the generalizability of the findings to other regions with different demographic compositions, legal systems, and law enforcement practices. A more extensive geographic sampling would enhance the study's external validity.

Algorithmic Transparency

While the paper sheds light on the biased outcomes of the COMPAS algorithm, it does not delve deeply into the algorithm's internal workings. Understanding the specific features or data points that contribute to biased predictions could offer valuable insights for mitigation strategies. A lack of transparency in proprietary algorithms is a broader issue that extends beyond the COMPAS case.

Societal Context

The study, while emphasizing the impact of technology on marginalized communities, does not extensively explore the broader societal context that contributes to biases in the criminal justice system. Factors such as historical prejudices, systemic racism, and socioeconomic disparities play a role and should be considered in a more comprehensive analysis.

Temporal Dynamics

The paper represents a snapshot in time, and the algorithms and practices it critiques may have evolved since its publication in 2016. Given the rapid pace of technological advancements, a longitudinal analysis could reveal whether improvements or exacerbations of biases have occurred over time.

Ethical Considerations

The paper touches on the ethical implications of algorithmic bias but does not extensively explore the ethical considerations surrounding the use of predictive policing algorithms. Ethical frameworks, biases in data collection, and the potential consequences of relying on algorithms for high-stakes decisions warrant further examination.

User Feedback and Agency Response

The study does not extensively incorporate user feedback, particularly from law enforcement agencies employing the COMPAS algorithm. Understanding how end-users interpret and respond to algorithmic predictions is crucial for a holistic evaluation of its impact on decision-making.

Intersectionality

The paper primarily addresses racial bias but does not extensively explore how other intersectional factors, such as gender, age, or socio-economic status, might compound or intersect with racial biases in algorithmic predictions.

Future work

Angwin et al.'s groundbreaking work on algorithmic bias in predictive policing, particularly in the COMPAS algorithm, opens avenues for further research and highlights crucial areas that demand exploration. The conclusions drawn from this study beckon researchers to delve deeper into the intricacies of algorithmic decision-making and its implications for fairness and justice in the criminal justice system.

Algorithmic Transparency and Explainability

One of the key takeaways is the need for increased transparency in the design and operation of predictive policing algorithms. Future research should aim to unravel the intricacies of these algorithms, emphasizing the importance of making their decision-making processes more accessible to stakeholders, including the public, policymakers, and those directly affected by algorithmic predictions. Understanding the specific features contributing to biased outcomes could facilitate the development of more transparent and accountable algorithms.

Broader Societal Context and Systemic Biases

The study prompts a call for research that considers the broader societal context influencing algorithmic biases. Investigating how historical prejudices, systemic racism, and socioeconomic factors interact with predictive policing algorithms could provide a more holistic understanding of the root causes of bias. Future work should explore the complex interplay between algorithmic decision-making and the societal structures in which these systems are embedded.

Longitudinal Analysis of Algorithmic Evolution

Given the rapid evolution of technology, a longitudinal analysis of predictive policing algorithms is essential. Monitoring how these algorithms have evolved since the publication of the study in 2016 would provide insights into whether improvements have been made to address biases or if new challenges have emerged. This longitudinal perspective would contribute to ongoing discussions about the adaptability and responsiveness of these systems to societal demands for fairness.

User Feedback and Decision-Maker Perspectives

Future research should incorporate feedback from end-users, such as law enforcement agencies and criminal justice practitioners employing predictive policing algorithms. Understanding how these individuals interpret and respond to algorithmic predictions is critical for refining these tools

and ensuring that they align with the objectives of the criminal justice system. Investigating the perspectives, biases, and decision-making processes of those responsible for implementing algorithmic recommendations would enrich the understanding of the human-AI interaction in law enforcement.

References to materials

Mehrabi et al. meticulously examine the landscape of bias and fairness in machine learning systems. By synthesizing existing literature, the authors categorize and analyze different types of biases that can emerge in these systems. The survey delves into both algorithmic biases and biases introduced through data collection processes. Furthermore, it provides valuable insights into various fairness metrics and mitigation techniques employed across different domains. Incorporating findings from this survey enriches our understanding of the broader context in which the biases observed in the COMPAS algorithm operate, allowing for a more nuanced and informed discussion on algorithmic fairness in the criminal justice system.

Holstein et al. present a pragmatic exploration of the challenges and requirements for improving fairness in machine learning systems, focusing on the perspective of industry practitioners. Drawing on real-world experiences, the authors discuss the trade-offs involved in addressing fairness concerns and the practical considerations that guide decision-making in industry settings. This paper contributes valuable insights into the implementation of fairness-enhancing measures, shedding light on the gap between theoretical discussions on algorithmic fairness and the practical needs of those tasked with deploying and managing machine learning systems. Integrating the perspectives from this paper ensures a well-rounded understanding of the actionable steps needed to enhance fairness in the criminal justice context.

Wang, Zhang, and Zhu provide a focused examination of algorithmic fairness. The paper synthesizes key concepts, methodologies, and challenges associated with ensuring fairness in algorithmic decision-making. By surveying existing literature, the authors outline prevalent fairness definitions and metrics while highlighting the evolving landscape of research in this domain. The review delves into the trade-offs between different fairness notions and addresses the complexities of mitigating biases in algorithmic systems. Incorporating insights from this review enriches our understanding of algorithmic fairness by providing a succinct overview of foundational concepts. The authors offer a valuable perspective on the nuances of defining and measuring fairness, laying the groundwork for a more comprehensive discussion on how these principles apply to predictive policing algorithms. This review serves as a complementary piece, bridging theoretical considerations with practical implications, and further strengthens the foundation for addressing algorithmic bias in the criminal justice system.

Lee et al. delve into the critical intersection of procedural justice and algorithmic fairness. The paper emphasizes the importance of procedural fairness in the design and deployment of algorithms, focusing on transparency and outcome control as key components. By evaluating the procedural aspects of algorithmic decision-making, the authors contribute valuable insights into how fairness considerations can be integrated into the mediation process. The work advocates for a balanced approach, considering not only the outcomes produced by algorithms but also the fairness of the procedures governing their operation. This paper enriches the discourse on algorithmic fairness by highlighting the procedural dimensions often overlooked in discussions centered solely on outcome fairness. The authors propose practical strategies for incorporating transparency and outcome control into algorithmic mediation, offering a roadmap for developers, policymakers, and stakeholders involved in deploying algorithmic systems. Integrating findings from this study into discussions on predictive policing algorithms enhances our understanding of how procedural justice principles can be harnessed to address bias and promote fairness in algorithmic decision-making within the criminal justice context.

Intersectionality in Algorithmic Bias

The study primarily addresses racial bias, but the intersectionality of biases remains a crucial area for exploration. Future research should investigate how multiple demographic factors, such as gender, age, and socio-economic status, intersect with racial biases in predictive policing algorithms. This more nuanced analysis would contribute to a comprehensive understanding of the varied ways individuals may be disproportionately impacted by algorithmic predictions.

Inspired by these open lines of research, one might consider interdisciplinary collaborations between computer scientists, social scientists, ethicists, and legal scholars to tackle the multifaceted challenges posed by algorithmic decision-making in criminal justice. Developing frameworks that prioritize transparency, fairness, and accountability, while considering the broader societal context, could be a promising avenue for future research. Additionally, engaging with communities affected by these algorithms to incorporate their perspectives and concerns can contribute to the development of more just and equitable predictive policing systems.

References

"Improving Fairness in Machine Learning Systems: What Do Industry Practitioners Need?" by Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miroslav Dudík, Hanna Wallach

"A Survey on Bias and Fairness in Machine Learning» by Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, Aram Galstyan

A brief review on algorithmic fairness by Xiaomeng Wang, Yishi Zhang & Ruilin Zhu

"Procedural Justice in Algorithmic Fairness: Leveraging Transparency and Outcome Control for Fair Algorithmic Mediation» by Min Kyung Lee, Anuraag Jain, Hea Jin Cha, Shashank Ojha, Daniel Kusbit

Task 2: Explanations with LIME

Introduction

In this experiment, I explore the application of the LIME (Local Interpretable Model-agnostic Explanations) framework on a classification model trained on the well-known Iris dataset. The objective is to evaluate the interpretability of the model's predictions at the local level, considering principles of user-centric explainable AI.

LIME basics

Here are the key concepts and features of LIME:

Local Interpretability

LIME focuses on providing explanations for individual predictions rather than explaining the entire model. It aims to make predictions interpretable on a per-instance basis, offering insights into why a specific prediction was made.

Model Agnostic

LIME is designed to be model-agnostic, meaning it can be applied to any machine learning model, regardless of its underlying architecture or learning algorithm. This flexibility allows LIME to be used with a wide range of models, including decision trees, support vector machines, neural networks, and more.

Simplified Proxies (Interpretable Models)

LIME generates local approximations of the complex model by training simpler, interpretable models (such as linear models or decision trees) on perturbed versions of the original data. These simpler models serve as "proxies" for the black-box model in the local neighborhood around the instance being explained.

Perturbation of Instances

To understand how the model behaves in the vicinity of a specific instance, LIME introduces slight perturbations to that instance and observes the changes in predictions. The generated perturbed instances, along with their corresponding model outputs, are then used to train the interpretable proxy model.

Feature Importance

The interpretable proxy model provides insight into the importance of different features for the specific prediction. LIME assigns weights to each feature, indicating their contribution to the model's output for the given instance.

Explanation for Humans

The final output of LIME is a human-interpretable explanation that highlights the most influential features and their impact on the model's decision for a particular instance. This helps users, stakeholders, or domain experts understand the rationale behind a model's prediction at a local level.

User-Centric Explainability

LIME aligns with the user-centric perspective on explainable AI. It emphasizes providing explanations that are understandable and useful to end-users, facilitating trust and transparency in the decision-making process of complex models.

In summary, LIME is a powerful tool for providing local, interpretable explanations for machine learning models, making their predictions more transparent and understandable. Its model-agnostic nature and focus on user-centric interpretability contribute to its wide adoption in the field of explainable AI.

Experimental Setup

Data Selection and Model Training

I chose the Iris dataset due to its simplicity and widespread use in machine learning. A Support Vector Machine (SVM) classifier was employed for its capability in handling multi-class classification tasks. The dataset was split into training and testing sets, with the SVM model trained on the former.

LIME Explanations

After training the SVM model, I applied LIME to generate local explanations for predictions on three instances from the test set. The LIME explanations provide insights into the decision-making process of the SVM model for each instance. Here are the LIME values for 3 inputs from the test set:

```
[('4.25 < petal length (cm) <= 5.10', 0.38294024349342254), ('0.30 < petal width (cm) <= 1.30', 0.05887276673104905), ('5.75 < sepal length (cm) <= 6.40', 0.032385473843352074), ('sepal width (cm) <= 2.80', 0.0033773058578228776)]
```

```
[('1.50 < petal length (cm) <= 4.25', 0.30435743511335245), ('petal width (cm) <= 0.30', 0.08919074496401534), ('sepal width (cm) > 3.40', -0.02906399625721439), ('5.10 < sepal length (cm) <= 5.75', -0.0232053382941127)]
```

```
[('petal length (cm) > 5.10', -0.3370096128381187), ('petal width (cm) > 1.80', -0.13810151296236559), ('sepal length (cm) > 6.40', 0.08396938168033771), ('sepal width (cm) <= 2.80', -0.0019029395061967774)]
```

Let's analyze the LIME explanations for the three instances.

- Instance 1:
 - Positive Contributions:
 - Petal length between 4.25 and 5.10 is a significant positive factor.
 - Petal width between 0.30 and 1.30 contributes positively.
 - Negative Contributions:
 - Sepal width less than or equal to 2.80 has a minor negative impact.
- Instance 2:
 - Positive Contributions:
 - Petal length between 1.50 and 4.25 is a major positive factor.
 - Petal width less than or equal to 0.30 contributes positively.
 - Negative Contributions:
 - Sepal width greater than 3.40 has a minor negative impact.
 - Sepal length between 5.10 and 5.75 has a minor negative impact.
- Instance 3:
 - Positive Contributions:
 - Petal length greater than 5.10 has a major positive impact.
 - Sepal length greater than 6.40 contributes positively.
 - Negative Contributions:
 - Petal width greater than 1.80 has a negative impact.

Discussion

The LIME explanations unveil crucial insights into how the SVM model is making predictions for individual instances in the Iris dataset. By identifying the specific ranges of features that contribute positively or negatively to predictions, we gain a deeper understanding of the decision boundaries learned by the model. For example, the positive contribution of petal length between 4.25 and 5.10 in Instance 1 suggests the model's reliance on this feature for classifying certain iris species.

Moreover, LIME highlights the significance of each feature within a local context. In Instance 2, petal length between 1.50 and 4.25 emerges as a major positive factor, indicating the model's emphasis on this feature for the given instance. This local interpretability is essential for users and domain experts seeking to grasp how the model weighs different features in specific scenarios.

Bias and Problems

While LIME explanations provide valuable insights, they do not explicitly flag biases in the model. However, certain negative contributions, such as the impact of sepal length in Instance 2, warrant further investigation. Understanding why the model considers a particular range as a negative factor for a specific instance is crucial for uncovering potential biases or areas where the model may be making suboptimal decisions.

The negative contribution associated with sepal length in Instance 2 might suggest that, for this specific instance, a longer sepal length contradicts the expected characteristics of a certain iris species. This observation underscores the nuanced nature of model interpretations, requiring a domain-specific lens to discern whether such features align with botanical expectations.

Trust and Transparency

Explanations generated by LIME contribute to the overall trustworthiness of the classification model. The transparency offered by LIME aligns with the user-centric principles of explainable AI, providing users with clear, understandable insights into the factors influencing individual predictions. In scenarios where model predictions impact critical decisions, such as species identification in botanical research, these explanations become instrumental in justifying model outputs to end-users and stakeholders.

The simplicity of the explanations, expressed as rule-based statements on feature ranges, enhances their interpretability. Users, including those without a deep understanding of machine learning, can readily comprehend why a specific prediction was made for a given instance. This aligns with the overarching goal of user-centric explainability, where the explanations are not only accurate but also accessible to a broader audience.

Framework of Wang et al. (2019)

The framework proposed by Wang et al. emphasizes the need for explanations that are actionable, understandable, and trust-inducing. LIME adheres to these principles by providing actionable insights—users can act on the identified feature ranges, such as modifying environmental conditions for iris growth. The understandability of LIME explanations is evident through their rule-based format, enabling users to comprehend the rationale behind model predictions.

Moreover, the transparency offered by LIME contributes to trust-building. Users are more likely to trust a model when they can comprehend its decision-making process. LIME's local interpretability aligns well with the user-centric approach, promoting trust by demystifying complex model behavior.

Fairness and Abstraction in Sociotechnical Systems

As we navigate the landscape of machine learning interpretability, it is crucial to not only comprehend the inner workings of models but also to critically examine the societal implications of their decisions. Diakopoulos and Friedler (2016) shed light on the intricate relationship between fairness and abstraction in sociotechnical systems, urging a thorough consideration of ethical dimensions in machine learning applications. This aligns with our exploration into interpretability using LIME, a framework designed to provide transparent, human-understandable explanations for individual predictions. The intersection of fairness considerations, as highlighted by Diakopoulos and Friedler, with the local interpretability offered by LIME underscores the broader commitment to user-centric explainable AI.

Conclusion

In conclusion, the detailed insights provided by LIME contribute significantly to the interpretability and trustworthiness of the SVM classifier on the Iris dataset. The nuanced understanding of feature contributions, identification of potential biases, and the alignment with user-centric principles underscore LIME's effectiveness in making machine learning models more transparent and actionable. As the field of explainable AI advances, frameworks like that of Wang et al. will continue to guide the development of models that not only make accurate predictions but also empower users through understandable and trustworthy explanations.

References

- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should I trust you?": Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016.
- Wang, F., Liu, Y., Liu, Z., & Tang, J. (2019). Explainable recommendation: A survey and new perspectives. *Foundations and Trends® in Information Retrieval*, 13(3), 161-309.

- Diakopoulos, N., & Friedler, S. A. (2016). Fairness and Abstraction in Sociotechnical Systems. *Data and Discrimination: Converting Critical Concerns into Productive Inquiry*, 41.