

Title

J.T. Cho

joncho@
seas.upenn.edu

Karinna Loo

kloo@
seas.upenn.edu

Veronica Wharton

whartonv@
seas.upenn.edu

Abstract

TODO: Abstract

1 Introduction

For our CIS 625 final project, our team — JT Cho, Karinna Loo, and Veronica Wharton — took a closer look at the topic of fairness in machine learning. The paper that piqued our interest was *Rawlsian Fairness for Machine learning* (Joseph et al., 2016), which discusses the procedures of studying the quantifying the discriminatory behaviors in automated decision-making procedures and proposing algorithms that both learn at a rate comparable to (but necessarily worse than) the best algorithms absent of a fairness constraint and also satisfy a specified fairness constraint. Specifically, our team was interested in exploring the paper’s results via implementation, as well as exploring what further applications for fairness analysis in real-world data might be possible.

TODO: Problem overview

TODO: Potential applications

TODO: Literature review, including overview of Joseph et al. (2016)

2 Project overview

Our project consisted of the following steps:

1. We read the paper *Rawlsian Fairness for Machine Learning* (Joseph et al., 2016).
2. We implemented the TopInterval, IntervalChaining, and RidgeFair algorithms from the paper in Python.
3. We ran our implementations on a Yahoo! dataset containing a fraction of the user click log for news articles displayed in the Featured Tab of the Today Module on the Yahoo! Front Page during the first ten days in May 2009, to see how well they performed on real data.

4. To empirically evaluate our implementations, we ran experiments similar to those in (Joseph et al., 2016) with randomly-drawn contexts.

5. We compiled our findings into a written report.

3 Implementation: IntervalChaining

4 Implementation: RidgeFair

5 Experimental results: generated data

TODO: Pretty figures

6 Experimental results: real data

TODO: Pretty figures

7 Conclusion

References

- [Joseph et al.2016] Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2016. Rawlsian fairness for machine learning. *CoRR*, abs/1610.09559.