

AI

Introducing AI Fairness 360

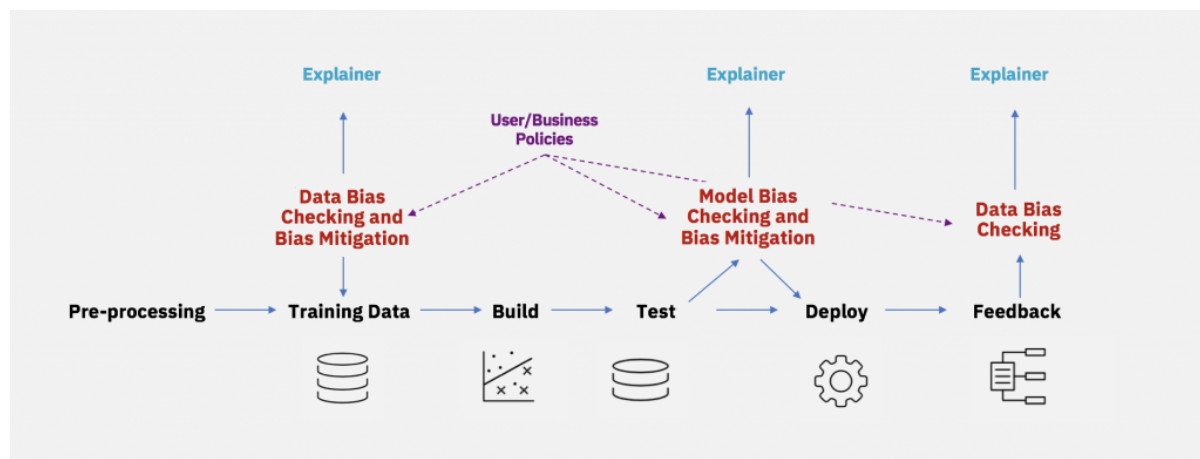
- September 19, 2018 | Written by: [Kush R. Varshney](#)

Categorized: [AI](#)

Share this post:



We are pleased to announce [AI Fairness 360](#) (AIF360), a comprehensive open-source toolkit of datasets and machine learning models, and state-of-the-art algorithms to mitigate such bias. We hope it to help engender [trust in AI](#) and make the world more equitable for all.



Mitigating bias throughout the AI lifecycle

Machine learning models are increasingly used to inform high-stakes decisions about people. A machine learning model, by its nature, is always a form of statistical discrimination, the discrimination becomes objectionable when it disadvantages certain groups at systematic advantage and certain unprivileged groups at systematic disadvantage. Bias in data, such as prejudice in labels or under-/over-sampling, [yields](#) models with unwanted bias.

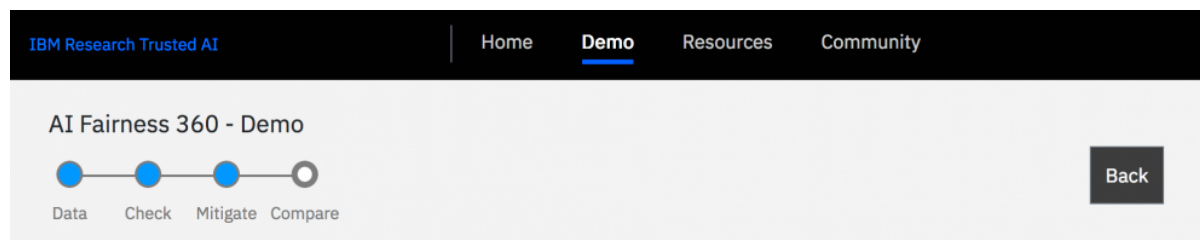
This initial release of the [AIF360 Python package](#) contains nine different algorithms, developed

research community, to mitigate that unwanted bias. They can all be called in a standard way, v paradigm. In this way, we hope that the package is not only a way to bring all of us researchers

IBM Research Blog Topics ▾ Labs ▾ About

focus on industrial usability, and its software engineering.

AIF360 contains three [tutorials](#) (with more to come soon) on credit scoring, predicting medical [images by gender](#). I would like to highlight the medical expenditure example; we've worked in t health insurance clients (without explicit fairness considerations), but it has not been considere before. (For background, here are [two papers](#) describing our earlier applied data science work i



4. Compare original vs. mitigated results

Dataset: Adult census income

Mitigation: [Optimized Pre-processing algorithm applied](#)

Protected Attribute: Race

Privileged Group: **White**, Unprivileged Group: **Non-white**

Accuracy after mitigation changed from 82% to 74%

Bias against unprivileged group was reduced to acceptable levels* for 1 of 2 previously biased metrics
(1 of 5 metrics still indicate bias for unprivileged group)



AI Fairness 360 interactive experience

AIF360 is not just a Python package. It is also an [interactive experience](#) that provides a gentle i capabilities of the toolkit. Being a comprehensive set of capabilities, it may be confusing to figu are most appropriate for a given use case. To help, we have created some [guidance material](#) tha

Our team includes members from the IBM India Research Lab and the T. J. Watson Research Center. We are a diverse lot in terms of national origin, scientific background, and expertise. We are using the toolkit as a summer project this year. We are a diverse lot in terms of national origin, scientific background, and expertise.

IBM Research Blog Topics ▾ Labs ▾ About

One of the reasons we decided to make AIF360 an open source project as a companion to the toolkit is to encourage the contribution of researchers from around the world to add their metrics and algorithms. AIF360 becomes the hub of a flourishing [community](#).

The currently implemented set of metrics and algorithms are described in the following list of papers:

[Flavio P. Calmon](#), [Dennis Wei](#), [Bhanukiran Vinzamuri](#), [Karthikeyan Natesan Ramamurthy](#), and [Michael Feldman](#), “[Processing for Discrimination Prevention](#),” Conference on Neural Information Processing Systems, 2016.

[Michael Feldman](#), [Sorelle A. Friedler](#), [John Moeller](#), [Carlos Scheidegger](#), and [Suresh Venkatasubramanian](#), “[Removing Disparate Impact](#),” ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016.

[Moritz Hardt](#), [Eric Price](#), and [Nathan Srebro](#), “[Equality of Opportunity in Supervised Learning](#),” Conference on Neural Information Processing Systems, 2016.

[Faisal Kamiran](#) and [Toon Calders](#), “[Data Preprocessing Techniques for Classification without Fairness](#),” Conference on Artificial Intelligence and Law, 2012.

[Faisal Kamiran](#), [Asim Karim](#), and [Xiangliang Zhang](#), “[Decision Theory for Discrimination-Aware Classification](#),” Conference on Data Mining, 2012.

[Toshihiro Kamishima](#), [Shotaro Akaho](#), [Hideki Asoh](#), and [Jun Sakuma](#), “[Fairness-Aware Classifiers with Counterfactual Fairness](#),” Joint European Conference on Machine Learning and Knowledge Discovery in Artificial Intelligence, 2012.

[Geoff Pleiss](#), [Manish Raghavan](#), [Felix Wu](#), [Jon Kleinberg](#), and [Kilian Q. Weinberger](#), “[On Fairness and Accuracy in Machine Learning](#),” Conference on Neural Information Processing Systems, 2017.

[Till Speicher](#), [Hoda Heidari](#), [Nina Grgic-Hlaca](#), [Krishna P. Gummadi](#), [Adish Singla](#), [Adrian Weller](#), and [David Rosenberg](#), “[A Unified Approach to Quantifying Algorithmic Unfairness: Measuring Individual & Group Unfairness](#),” ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2018.

[Richard Zemel](#), [Yu \(Ledell\) Wu](#), [Kevin Swersky](#), [Toniann Pitassi](#), and [Cynthia Dwork](#), “[Learning without Fairness](#),” Conference on Machine Learning, 2013.

[Brian Hu Zhang](#), [Blake Lemoine](#), and [Margaret Mitchell](#), “[Mitigating Unwanted Biases with Adversarial Debiasing](#),” Conference on Artificial Intelligence, Ethics, and Society, 2018.

¹Some of the excellent repositories are [Aequitas](#), [Audit-AI](#), [FairML](#), [Fairness Comparison](#), [Fairness Indicators](#), [Themis-ML](#).

²AIF360 team members are Rachel Bellamy, Kuntal Dey, Mike Hind, Sam Hoffman, Stephanie Hoyer, Jacquelyn Martino, Sameep Mehta, Saška Mojsilović, Seema Nagar, Karthi Natesan Ramamurthy, and David Rosenberg.

Sattigeri, Moninder Singh, Kush Varshney, Dakuo Wang, and Yunfeng Zhang.

IBM Research Blog

Topics ▾

Labs ▾

About



Principal Research Staff Member and Manager, IBM Research

AI

Artificial Intelligence

trusted AI

[← Previous Post](#)

[Helping to Improve Medical Image Analysis with Deep Learning](#)

[Trust an](#)

More AI stories

IBM Research

The world is our lab

[Learn more](#)

Connect with us

IBM Research Blog

Topics ▾

Labs ▾

About

