

# RSA® Conference 2022

San Francisco & Digital | February 7 – 10

## TRANSFORM

SESSION ID: MLAI-M06

# It's Not Fair! Detecting Algorithmic Bias with Open Source Tools

**Mike Kiser**

Director of Strategy and Standards  
SailPoint  
 @\_MikeKiser

**Mo Badawy**

Principal Data Scientist  
SailPoint



# Disclaimer

Presentations are intended for educational purposes only and do not replace independent professional judgment. Statements of fact and opinions expressed are those of the presenters individually and, unless expressly stated to the contrary, are not the opinion or position of RSA Conference LLC or any other co-sponsors. RSA Conference does not endorse or approve, and assumes no responsibility for, the content, accuracy or completeness of the information presented.

Attendees should note that sessions may be audio- or video-recorded and may be published in various media, including print, audio and video formats without further notice. The presentation template and any media capture are subject to copyright protection.

©2022 RSA Conference LLC or its affiliates. The RSA Conference logo and other trademarks are proprietary. All rights reserved.

TEMPLE OF JUSTICE







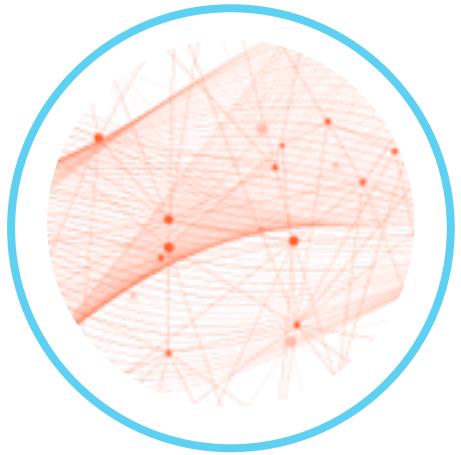








**Transparency**



**Fairness**



**Transparency**



**Fairness**

# Why are bubbles round?



Consider a smooth surface in  $\mathbb{R}^{n+1}$  representing the graph of a function  $x_{n+1} = u(x_1, \dots, x_n)$  defined on a bounded open set  $\Omega$  in  $\mathbb{R}^n$ . Assuming that  $u$  is sufficiently smooth, the area of the surface is given by the nonlinear functional

$$\mathcal{A}(u) = \int_{\Omega} \left(1 + |\nabla u|^2\right)^{1/2} dx_1 \dots dx_n, \quad (1)$$

where  $\nabla u$  is the gradient vector  $(\partial u / \partial x_1, \dots, \partial u / \partial x_n)$  and  $|\nabla u|^2 = (\nabla u) \cdot (\nabla u)$ .



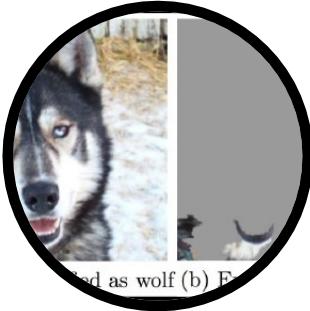
# Impact of Transparency



Safety

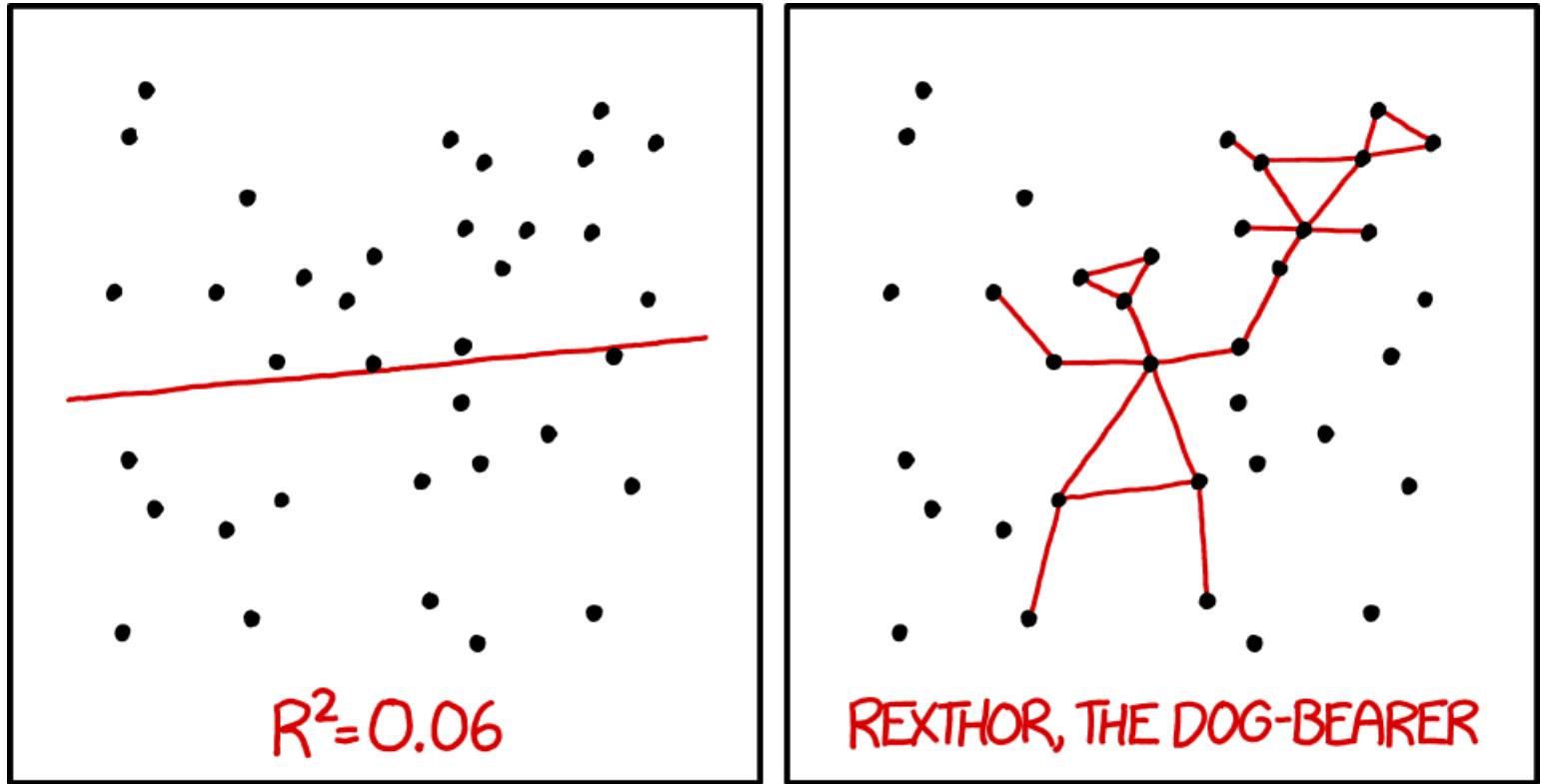


Model Adjustment



Objective Assessment

# Model Interpretability

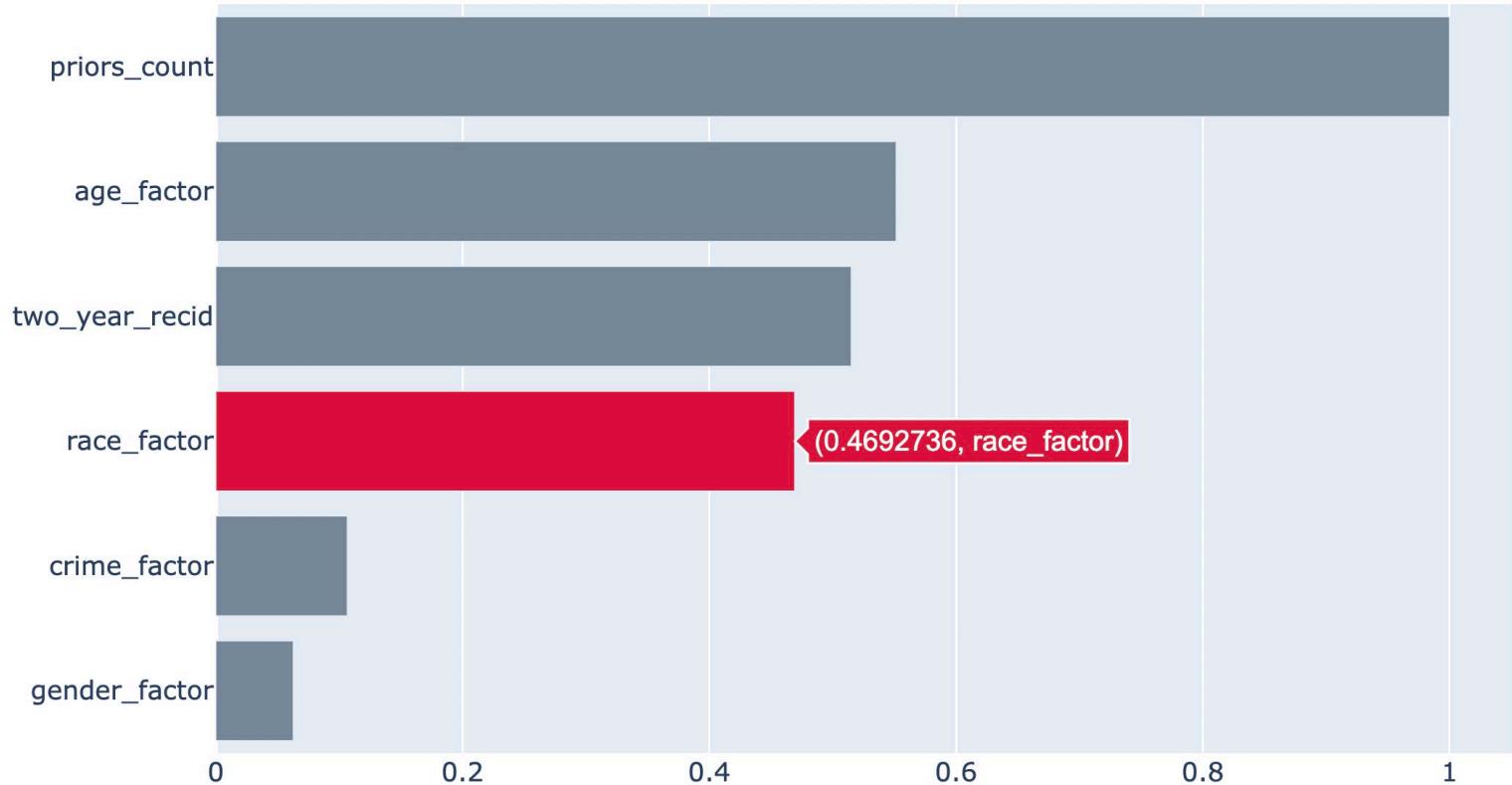


I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

<https://xkcd.com/1725/>

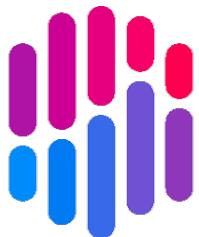
# COMPAS Model Interpretability

COMPAS non-violent 2 yr redicivism - Variable importance



# Model Agnostic Methods

- **Global vs Local**
- **LIME - Local Interpretable Model-Agnostic Explanations**
- **SHAP - SHapley Additive exPlanations**



# Comparison of Transparency Tools

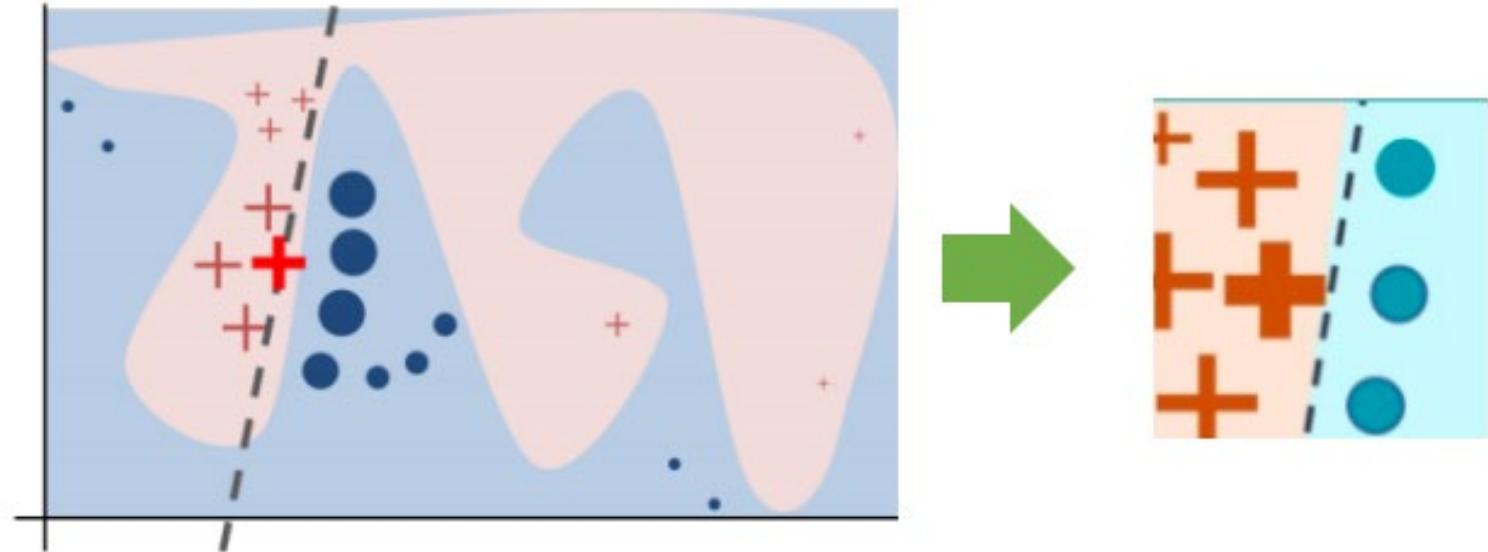
## Local Interpretable Model-agnostic Explanations (LIME)

- Local Model Approximation
- Less Intuitive to Human Reasoning
- Relatively Fast



- Game Theory
- More Intuitive to Human Reasoning
- Relatively Slow  
*(Faster on tree structures,  
slow on k-nearest neighbor)*

# Global vs Local Interpretation



- **Global interpretation:** general landscape of features/attributes
- **Local interpretation:** a particular instance or a small set of instances.

# Diabetes Progression Dataset

**Prediction Target:**

Quantitative measure of disease progress

**Instances:**

442

**Attributes:**

Age

Sex

Body Mass Index (BMI)

Blood Pressure (BP)

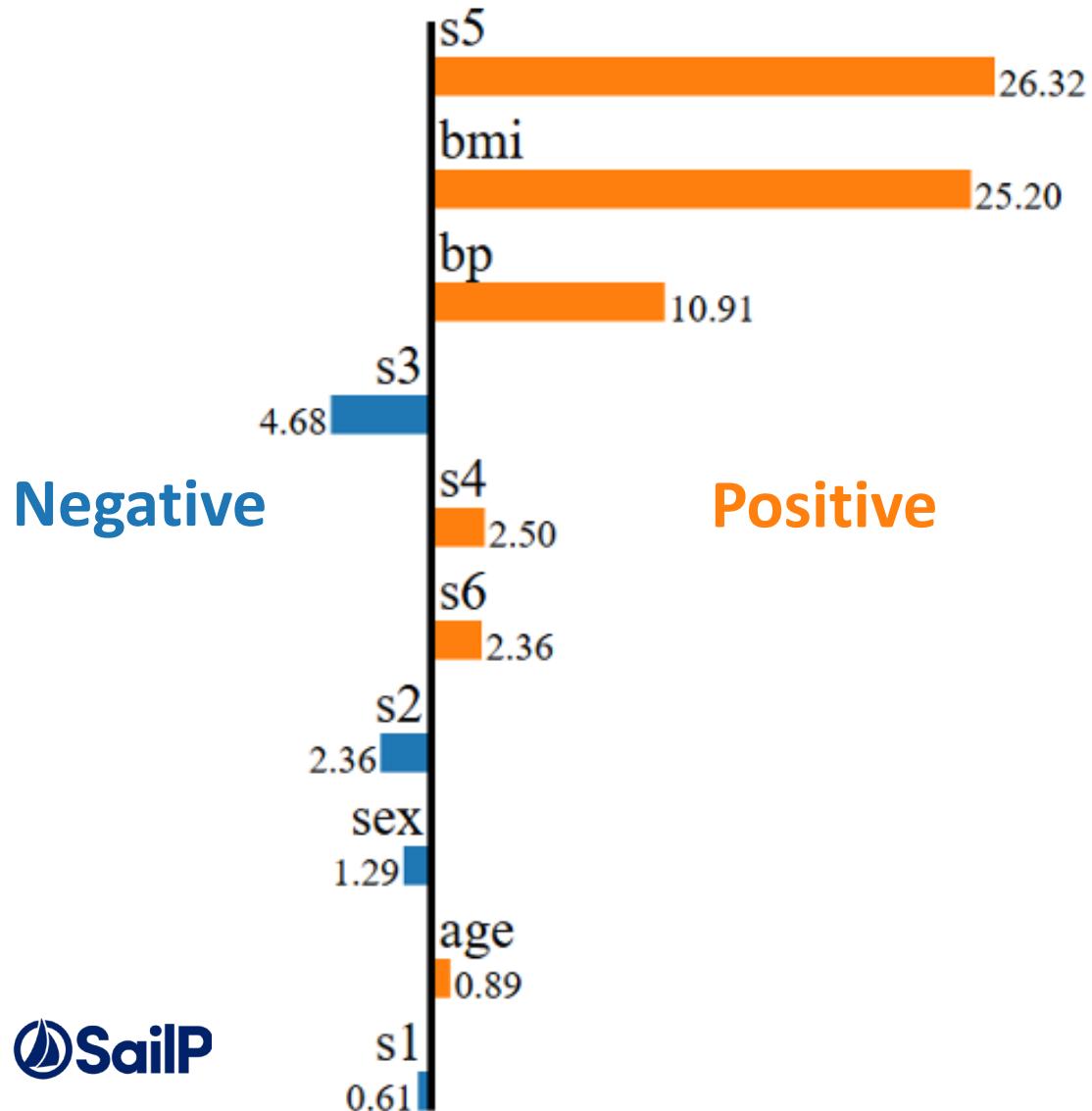
LDL

HDL

...

[https://web.stanford.edu/~hastie/StatLearnSparsity\\_files/DATA/diabetes.html](https://web.stanford.edu/~hastie/StatLearnSparsity_files/DATA/diabetes.html)

# LIME Results: Diabetes Dataset



Predicted value

68.69 (min) 175.45 280.27 (max)

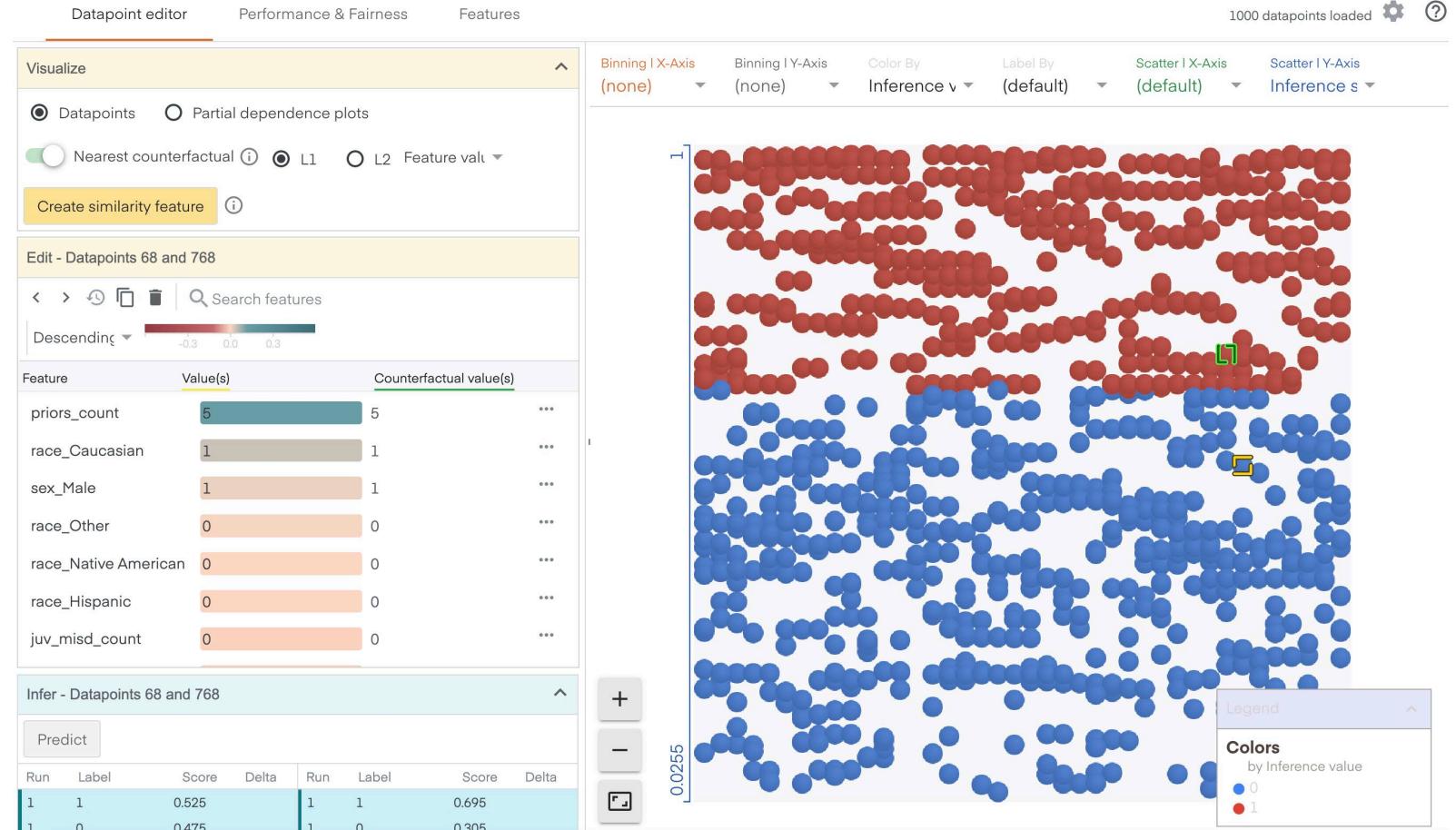
- Predicted value is a representation of diabetes progression
- Features are identified by their relative impact

# SHAP Results: Diabetes Dataset

- Shapley value in game theory
- Shows change in prediction value by feature



# Google's What-If Tool



## Data Exploration

**“Counter Factuals”**

**Fairness**

**Supported Models:**

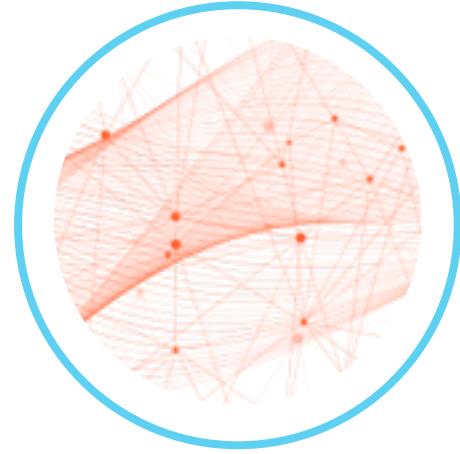
Binary classification

Mutli-class classification

Regression tasks



Transparency



Fairness

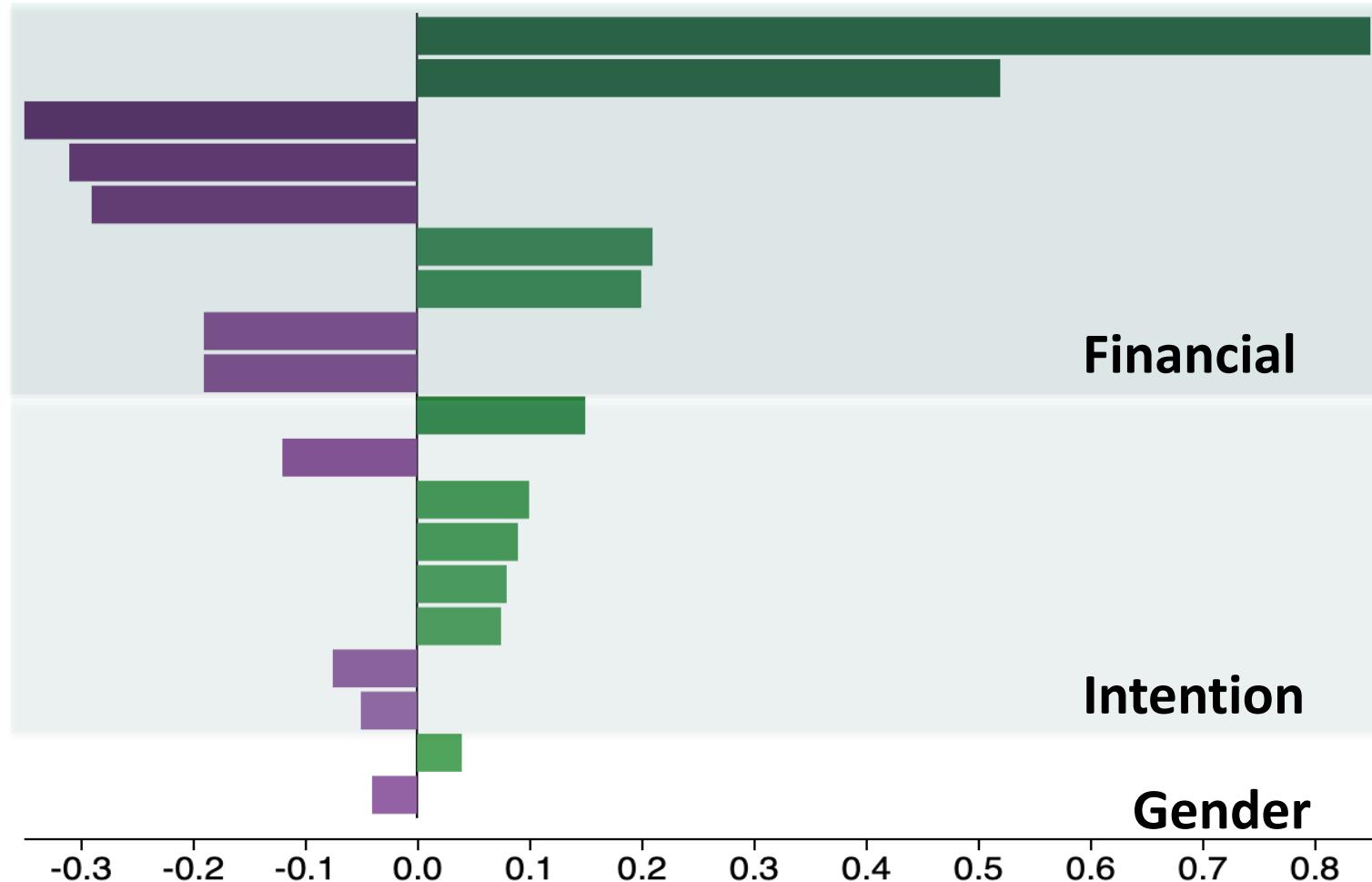
# Bias

# Shortcut

# Shortcut | Impartiality

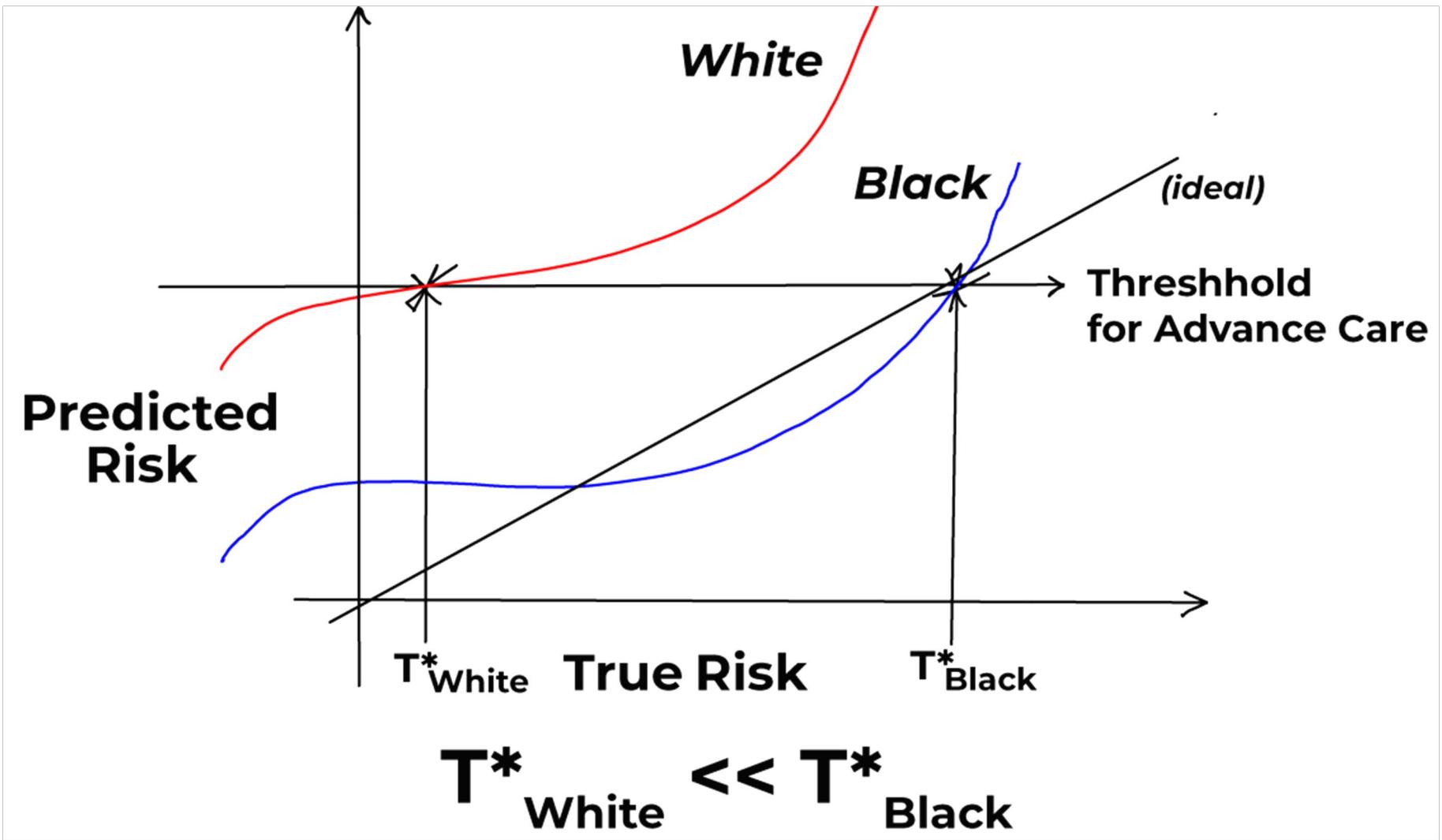
**Shortcut | Impartiality | Self-Interest**

# Gender Bias in German Credit Dataset



- Gender part of original data
- Checking for bias is relatively straightforward
- Gender is relatively low impact

# Hidden Bias: Renal Failure Dataset



- Racial bias not immediately apparent
- Algorithm was optimized for cost

# Aequitas Fairness Toolkit

## Bias and Fairness Audit Toolkit

The Bias Report is powered by [Aequitas](#), an open-source bias audit toolkit for machine learning developers, analysts, and policymakers to audit machine learning models for discrimination and bias, and make informed and equitable decisions around developing and deploying predictive risk-assessment tools.



See an [example report](#) on COMPAS risk assessment scores.

Or try out the audit tool using your own data or one of our sample data sets.

[Get Started!](#)

# Aequitas Result Set



## Equal Parity

Also known as  
Demographic or Statistical  
Parity



## Proportional Parity

Also known as Impact Parity  
or Minimizing Disparate  
Impact



## False Positive Parity

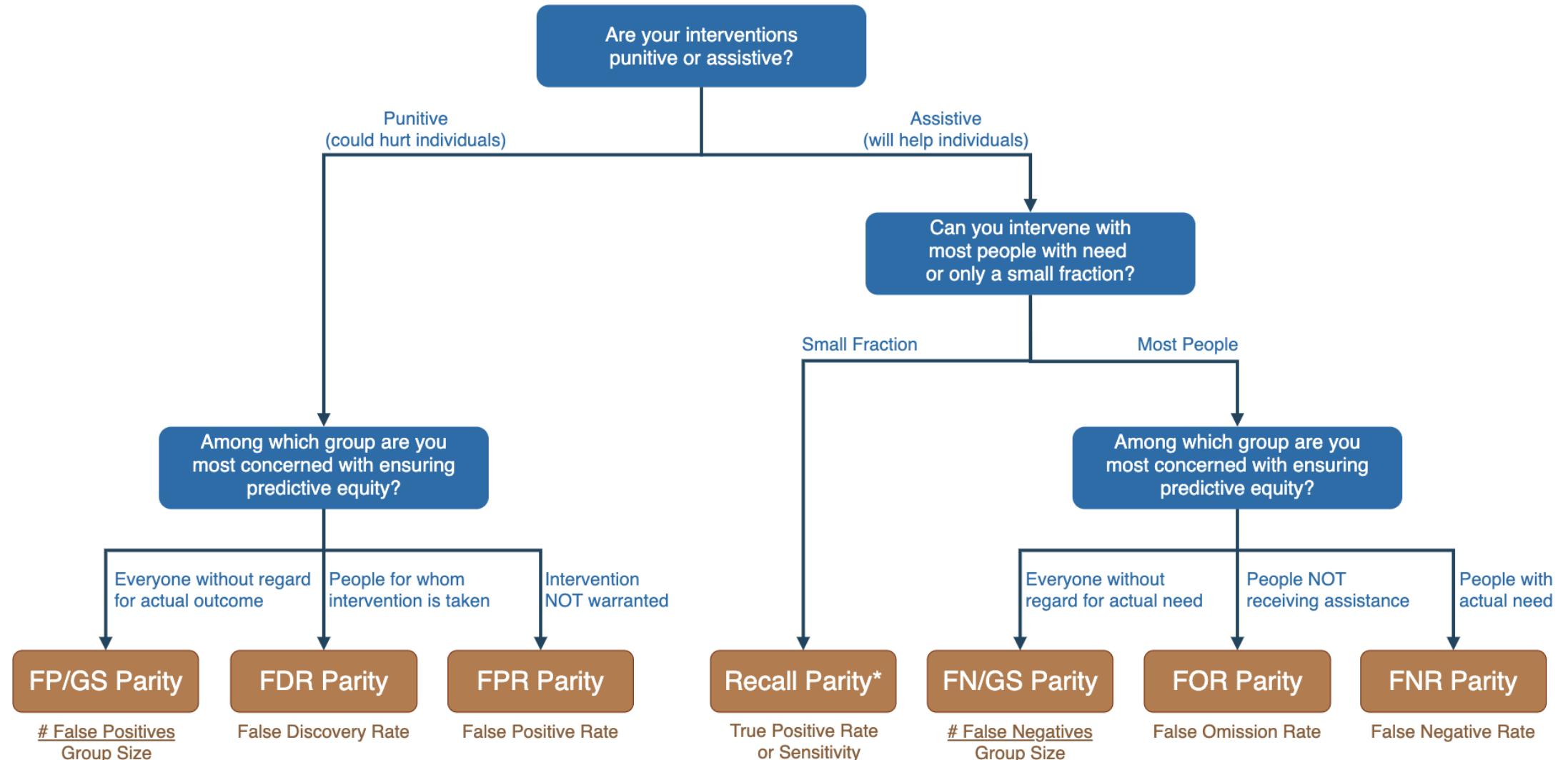
Desirable when your  
interventions are punitive



## False Negative Parity

Desirable when your  
interventions are  
assistive/preventative

# Aequitas Fairness Decision Tree



# Aequitas: Formatting the Data Set

**score**

The predicted outcome

**label\_value**

Ground truth (the actual result)

**attribute\_n**

Attribute within the model

	<b>score</b>	<b>label_value</b>	<b>attribute_1</b>	...	<b>attribute_n</b>
<b>score</b> The predicted outcome	0	1	A		40
	1	0	B		30
<b>label_value</b> Ground truth (the actual result)	1	1	A		32
	0	1	B		50
<b>attribute_n</b> Attribute within the model					
	↑	↑	↑		↑
	binary	binary	categorical or continuous		categorical or continuous (will be discretized into quartiles)



# Comparison of Fairness Tools



<b>Fairness Comparison</b>	Benchmarks fairness of algorithms
<b>AI Fairness 360</b>	Embeddable code for bias/mitigation Developer-focused Input cleansing
<b>Audit-AI</b>	Equal Employment Opportunity Commission Standards Fairness in hiring driven
<b>Thesis-ML</b>	Fairness-aware ML algorithms
<b>Aequitas</b>	Toolkit (libraries, command line, web interface) Usable by policymakers and auditors
<b>FairML</b>	Calculates attribute significance for black-box models
<b>FairTest</b>	Subgroup / Decision Tree created from user population Reports on subgroup / feature association for potential bias
<b>Themis</b>	Python module for statistical analysis of bias
<b>Fairness Measures</b>	Top K ranking algorithms

# So, What Now?

- Next week:
  - Inventory and consider the use cases of machine learning within your organization (both internal and outsourced)
  - Discuss with IRB, etc about potential issues that might need to be addressed
- In three months:
  - Investigate the tools that might be applicable for your datasets and machine learning usage:
    - Transparency: LIME / SHAP
    - Fairness: (Various tools)
- Within six months:
  - Begin to implement these tools and seek to rectify issues that may present themselves

TEMPLE OF JUSTICE

# Links to Various Tools



## Transparency

Lime	<a href="https://github.com/marcotcr/lime">https://github.com/marcotcr/lime</a>
SHAP	<a href="https://github.com/slundberg/shap">https://github.com/slundberg/shap</a>
Google What-if?	<a href="https://pair-code.github.io/what-if-tool/">https://pair-code.github.io/what-if-tool/</a>

## Fairness

Aequitas	<a href="https://github.com/dssg/aequitas">https://github.com/dssg/aequitas</a>
Fairness Comparison	<a href="https://github.com/algofairness/fairness-comparison">https://github.com/algofairness/fairness-comparison</a>
AI Fairness 360	<a href="https://github.com/IBM/aif360">https://github.com/IBM/aif360</a>
Audit-AI	<a href="https://github.com/pymetrics/audit-ai">https://github.com/pymetrics/audit-ai</a>
Thesis-ML	<a href="https://github.com/cosmicBboy/themis-ml">https://github.com/cosmicBboy/themis-ml</a>
FairML	<a href="https://github.com/adebayoj/fairml">https://github.com/adebayoj/fairml</a>
FairTest	<a href="https://github.com/columbia/fairtest">https://github.com/columbia/fairtest</a>
Themis	<a href="https://themis-ml.readthedocs.io/en/latest/index.html">https://themis-ml.readthedocs.io/en/latest/index.html</a>
Fairness Measures	<a href="https://fairnessmeasures.github.io/">https://fairnessmeasures.github.io/</a>

***Explore, Explain and Examine Predictive Models:*** <https://pbiecek.github.io/ema/>