

Responsible Data Science

Interpretability & Legal Frameworks

May 2 & 4, 2022

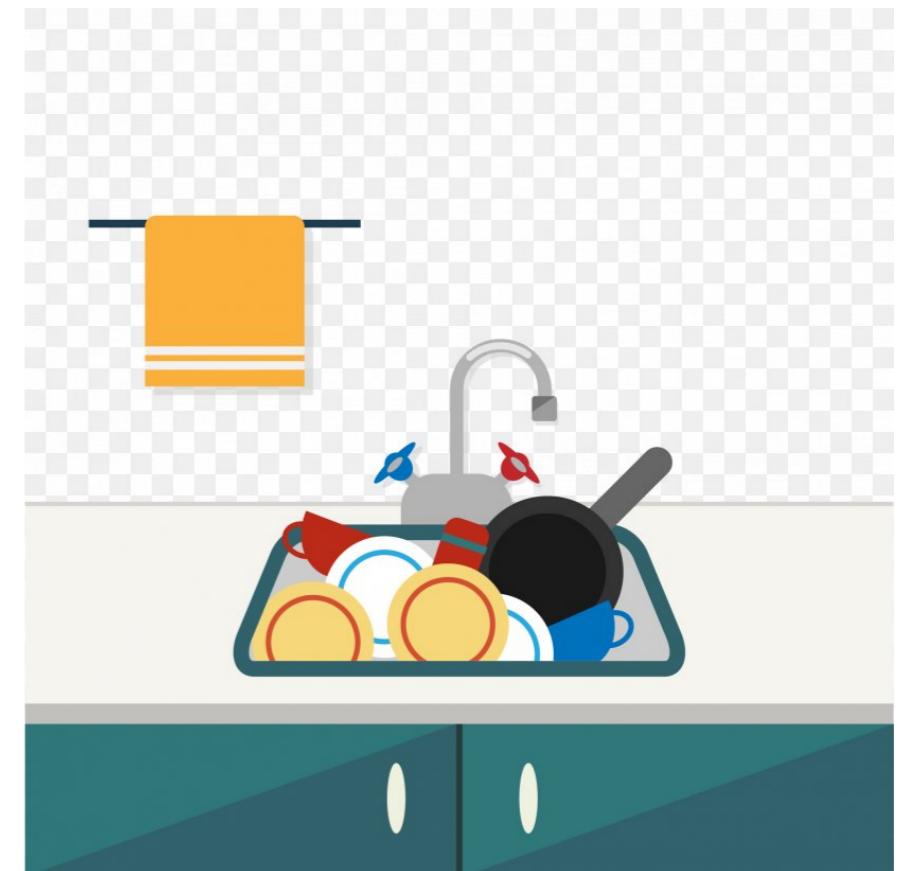
Prof. George Wood

Center for Data Science
New York University

What is interpretability?

- Explaining black-box models
- Online ad targeting
- Interpretability

A kitchen sink? Or a foundational concept
for responsible data science?



(Image source)

Algorithmic rankers

<https://freedom-to-tinker.com/2016/08/05/revealing-algorithmic-rankers/>

Input: database of items (individuals, colleges, cars, ...)

Score-based ranker: computes the score of each item using a known formula, often a monotone aggregation function, then sorts items on score

Output: permutation of the items, complete or top-k

Do we have transparency?

\mathcal{D}			f
id	x_1	x_2	$x_1 + x_2$
t_1	0.63	0.71	1.34
t_2	0.72	0.65	1.37
t_3	0.58	0.78	1.36
t_4	0.7	0.68	1.38
t_5	0.53	0.82	1.35
t_6	0.61	0.79	1.4

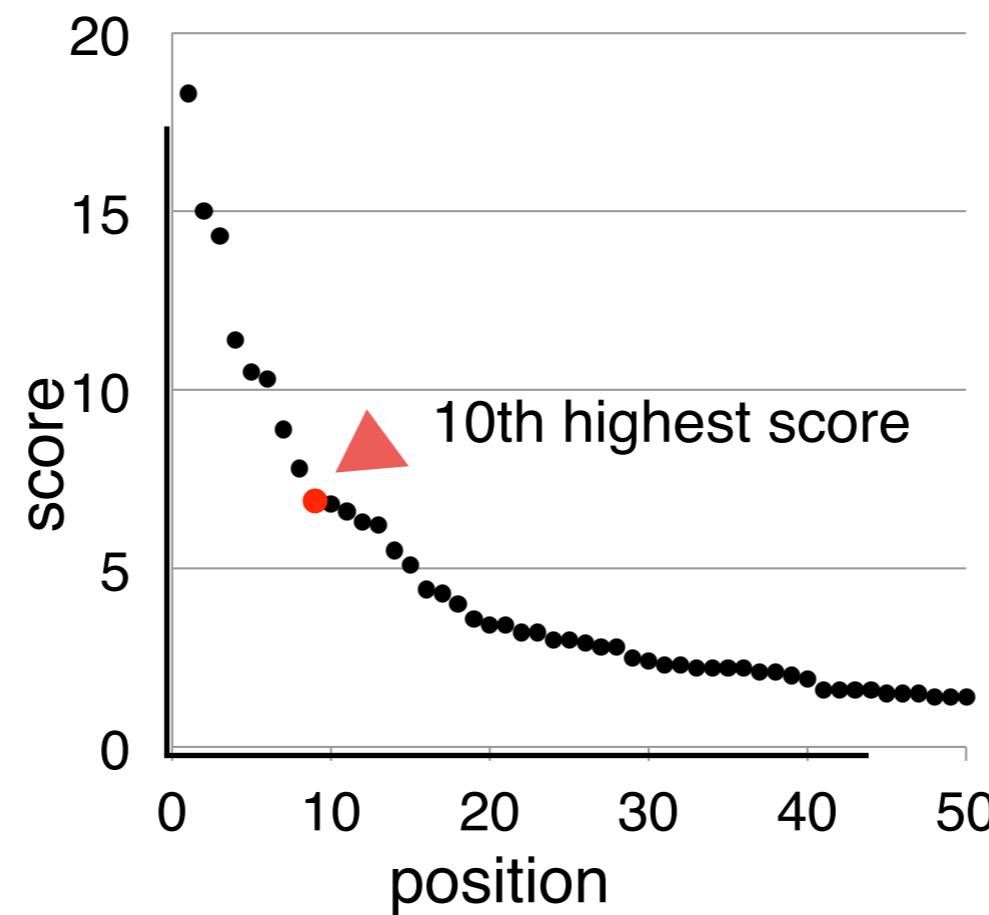
We have syntactic transparency, but lack interpretability!

Opacity in algorithmic rankers

<https://freedom-to-tinker.com/2016/08/05/revealing-algorithmic-rankers/>

Reason 1: The scoring formula alone does not indicate the relative rank of an item.

Scores are absolute, rankings are relative. Is 5 a good score? What about 10? 15?



Opacity in algorithmic rankers

<https://freedom-to-tinker.com/2016/08/05/revealing-algorithmic-rankers/>

Reason 2: A ranking may be unstable if there are tied or nearly-tied items.

Rank	Institution	Average Count	Faculty
1	► Carnegie Mellon University	18.4	123
2	► Massachusetts Institute of Technology	15.6	64
3	► Stanford University	14.8	56
4	► University of California - Berkeley	11.5	50
5	► University of Illinois at Urbana-Champaign	10.6	56
6	► University of Washington	10.3	50
7	► Georgia Institute of Technology	8.9	81
8	► University of California - San Diego	8	51
9	► Cornell University	7	45
10	► University of Michigan	6.8	63
11	► University of Texas - Austin	6.6	43
12	► University of Massachusetts - Amherst	6.4	47

Opacity in algorithmic rankers

<https://freedom-to-tinker.com/2016/08/05/revealing-algorithmic-rankers/>

Reason 3: A ranking methodology may be unstable:
small changes in weights can trigger significant re-shuffling.

THE NEW YORKER

DEPT. OF EDUCATION FEBRUARY 14 & 21, 2011 ISSUE

THE ORDER OF THINGS

What college rankings really tell us.



By Malcolm Gladwell

- | | | |
|---------------------------|---------------------------|---------------------------|
| 1. Porsche Cayman 193 | 2. Chevrolet Corvette 186 | 1. Chevrolet Corvette 205 |
| 3. Lotus Evora 182 | 2. Lotus Evora 195 | 3. Porsche Cayman 195 |
| 1. Lotus Evora 205 | 2. Porsche Cayman 198 | |
| 3. Chevrolet Corvette 192 | | |

<https://www.newyorker.com/magazine/2011/02/14/the-order-of-things>

Opacity in algorithmic rankers

<https://freedom-to-tinker.com/2016/08/05/revealing-algorithmic-rankers/>

Reason 4: The weight of an attribute in the scoring formula does not determine its impact on the outcome.

Rank	Name	Avg Count	Faculty	Pubs	GRE
1	CMU	18.3	122	2	791
2	MIT	15	64	3	772
3	Stanford	14.3	55	5	800
4	UC Berkeley	11.4	50	3	789
5	UIUC	10.5	55	3	772
6	UW	10.3	50	2	796
39	U Chicago	2	• • •	28	779
40	UC Irvine	1.9	28	2	787
41	BU	1.6	15	2	783
41	U Colorado Boulder	1.6	32	1	761
41	UNC Chapel Hill	1.6	22	2	794
41	Dartmouth	1.6	18	2	794

Given a score function:
$$0.2 * faculty +$$
$$0.3 * avg\ cnt +$$
$$0.5 * gre$$

Rankings are not benign!

THE NEW YORKER

DEPT. OF EDUCATION FEBRUARY 14 & 21, 2011 ISSUE

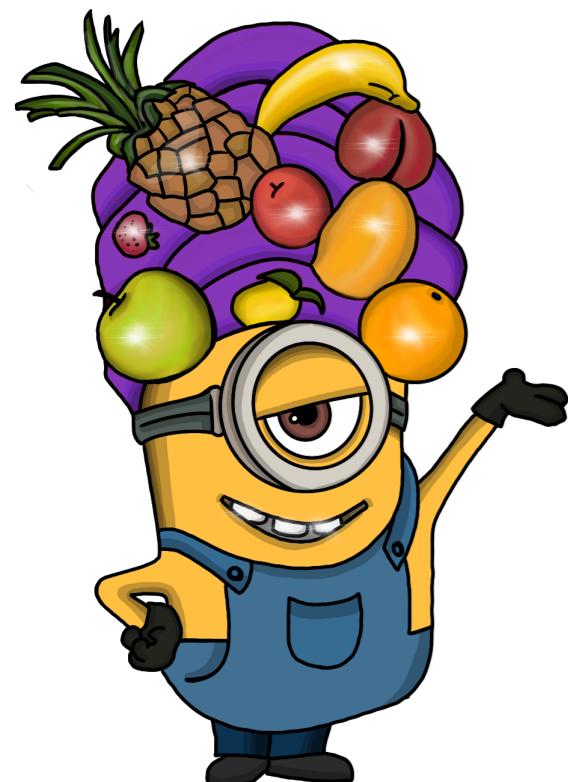
THE ORDER OF THINGS

What college rankings really tell us.



By Malcolm Gladwell

Rankings are not benign. They enshrine very particular ideologies, and, at a time when American higher education is facing a crisis of accessibility and affordability, we have adopted **a de-facto standard of college quality** that is uninterested in both of those factors. And why? Because a group of magazine analysts in an office building in Washington, D.C., decided twenty years ago to **value selectivity over efficacy**, to **use proxies** that scarcely relate to what they're meant to be proxies for, and to **pretend that they can compare** a large, diverse, low-cost land-grant university in rural Pennsylvania with a small, expensive, private Jewish university on two campuses in Manhattan.



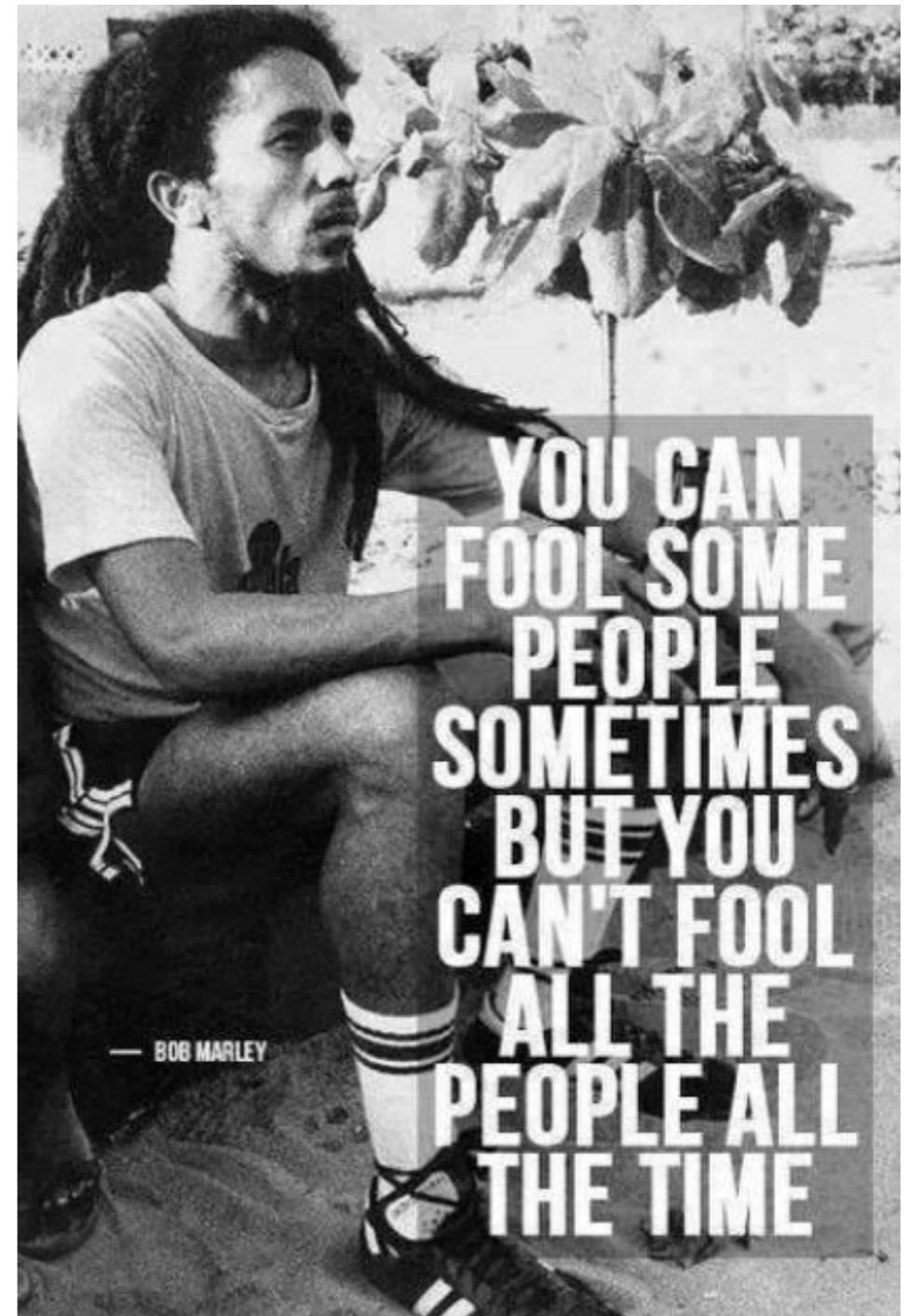
Interpretability in the service of trust!

Gladwell makes the point that rankings are claiming objectivity, yet are comparing apples and oranges.

In that sense, **a score-based ranker is a quintessential “black box” of data science**, and perhaps the simplest possible such black box.

AI is a red herring, privacy / IP / gaming arguments are overused. The truly difficult issues are that:

- 1) using math to pretend that we are correct when making intrinsically subjective decisions reinforcing the balance of power in society
- 2) that math / objectivity is used as a substitute for trust, but **trust must run deeper than math!**
- 3) need to find a kind of an interpretability that will enable trust!



The fairness you asked for is inside this box



[Arif Khan, Manis & Stoyanovich, 2021]

from transparency to
interpretability

New York City Local Law 49

January 11, 2018

Local Law 49 of 2018 in relation to automated decision systems used by agencies

 THE NEW YORK CITY COUNCIL Sign In

Corey Johnson, Speaker LEGISLATIVE RESEARCH CENTER

Council Home Legislation Calendar City Council Committees RSS Alerts

Details Reports

File #: Int 1696-2017 Version: A A Name: Automated decision systems used by agencies.

Type: Introduction Status: Enacted Committee: [Committee on Technology](#)

On agenda: 8/24/2017

Enactment date: 1/11/2018 Law number: 2018/049

Title: A Local Law in relation to automated decision systems used by agencies

Sponsors: [James Vacca](#), [Helen K. Rosenthal](#), [Corey D. Johnson](#), [Rafael Salamanca, Jr.](#), [Vincent J. Gentile](#), [Robert E. Cornegy, Jr.](#), [Jumaane D. Williams](#), [Ben Kallos](#), [Carlos Menchaca](#)

Council Member Sponsors: 9

Summary: This bill would require the creation of a task force that provides recommendations on how information on agency automated decision systems may be shared with the public and how agencies may address instances where people are harmed by agency automated decision systems.

Indexes: Oversight

Attachments: 1. [Summary of Int. No. 1696-A](#), 2. [Summary of Int. No. 1696](#), 3. [Int. No. 1696](#), 4. [August 24, 2017 - Stated Meeting Agenda with Links to Files](#), 5. [Committee Report 10/16/17](#), 6. [Hearing Testimony 10/16/17](#), 7. [Hearing Transcript 10/16/17](#), 8. [Proposed Int. No. 1696-A - 12/12/17](#), 9. [Committee Report 12/7/17](#), 10. [Hearing Transcript 12/7/17](#), 11. [December 11, 2017 - Stated Meeting Agenda with Links to Files](#), 12. [Hearing Transcript - Stated Meeting 12-11-17](#), 13. [Int. No. 1696-A \(FINAL\)](#), 14. [Fiscal Impact Statement](#), 15. [Legislative Documents - Letter to the Mayor](#), 16. [Local Law 49](#), 17. [Minutes of the Stated Meeting - December 11, 2017](#)

The original draft

Int. No. 1696

August 16, 2017

By Council Member Vacca

A Local Law to amend the administrative code of the city of New York, in relation to automated processing of **data** for the purposes of targeting services, penalties, or policing to persons

Be it enacted by the Council as follows:

- 1 Section 1. Section 23-502 of the administrative code of the city of New York is amended
- 2 to add a new subdivision g to read as follows:
 - 3 g. Each agency that uses, for the purposes of targeting services to persons, imposing
 - 4 penalties upon persons or policing, an algorithm or any other method of automated processing
 - 5 system of **data** shall:
 - 6 1. Publish on such agency's website, the source code of such system; and
 - 7 2. Permit a user to (i) submit **data** into such system for self-testing and (ii) receive the
 - 8 results of having such **data** processed by such system.
- 9 § 2. This local law takes effect 120 days after it becomes law.

MAJ
LS# 10948
8/16/17 2:13 PM

Point 1

algorithmic transparency is not synonymous with releasing the source code

publishing source code helps, but it is sometimes unnecessary and often insufficient

Point 2

algorithmic transparency requires data transparency

data is used in training, validation, deployment

validity, accuracy, applicability can only be understood in the data context

data transparency is necessary for all ADS, not only for ML-based systems

Point 3

data transparency is not synonymous
with making all data public

release data whenever possible;

also release:

data selection, collection and pre-processing
methodologies; data provenance and quality
information; known sources of bias; privacy-
preserving statistical summaries of the data

Data Synthesizer

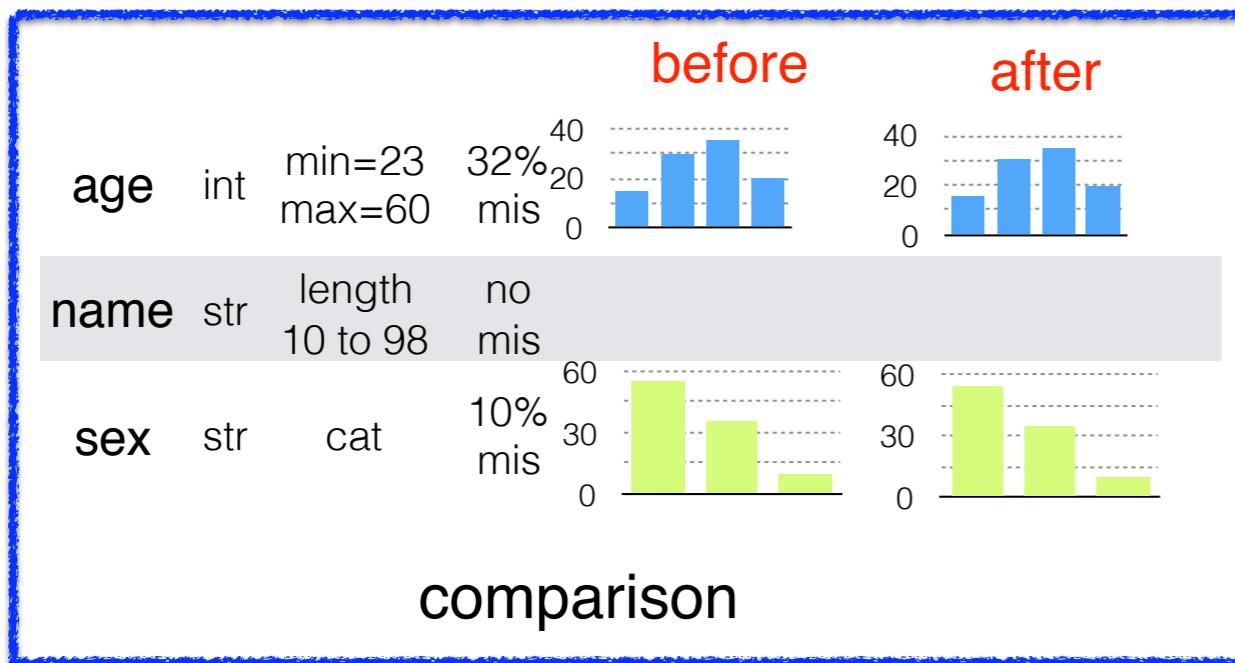
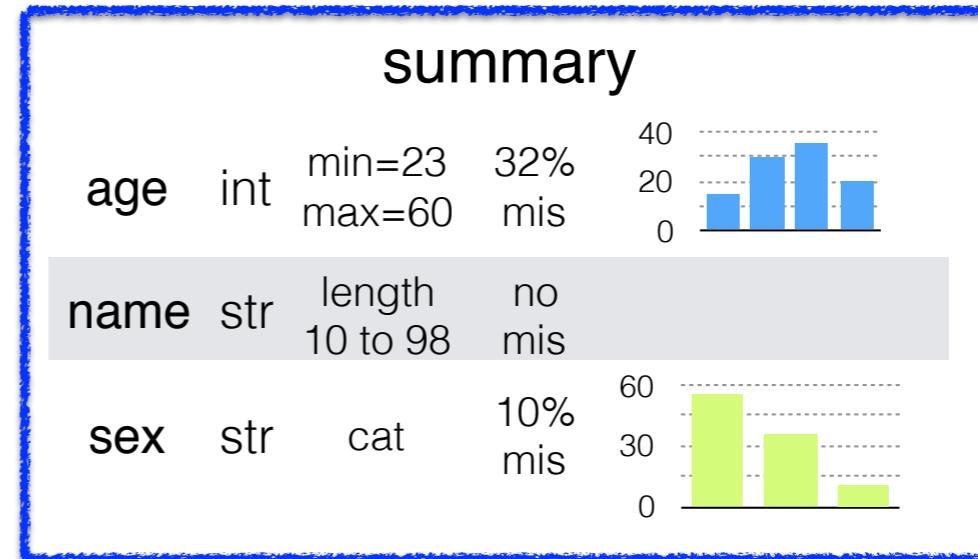
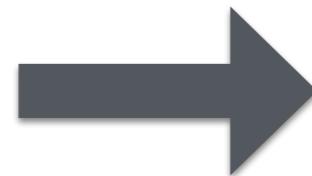
[Ping, Stoyanovich, Howe **SSDBM 2017**]

<http://demo.dataresponsibly.com/synthesizer/>

1	A	B	C	D	E	F	G	H		
1	UID	sex	race	MarriageStat	DateOfBirth	age	juv_fel	cour	decile	score
2	1	0	1	1	4/18/47	69	0	1		
3	2	0	2	1	1/22/82	34	0	3		
4	3	0	2	1	5/14/91	24	0	4		
5	4	0	2	1	1/21/93	23	0	8		
6	5	0	1	2	1/22/73	43	0	1		
7	6	0	1	3	8/22/71	44	0	1		
8	7	0	3	1	7/25/74	41	0	6		
9	8	0	1	2	2/25/73	49	0	4		
10	9	0	3	1	6/10/94	21	0	3		
11	10	0	3	1	6/10/88	27	0	4		
12	11	1	3	2	8/22/78	37	0	1		
13	12	0	2	1	12/2/74	41	0	4		
14	13	1	3	1	6/14/68	47	0	1		
15	14	0	2	1	3/25/85	31	0	3		
16	15	0	4	4	1/25/79	37	0	1		
17	16	0	2	1	6/22/90	25	0	10		
18	17	0	3	1	12/2/81	31	0	5		
19	18	0	3	1	5/8/85	31	0	3		
20	19	0	2	3	6/28/51	64	0	6		
21	20	0	2	1	11/29/94	21	0	9		
22	21	0	3	1	8/6/88	27	0	2		
23	22	1	3	1	3/22/95	21	0	4		
24	23	0	4	1	1/23/92	24	0	4		
25	24	0	3	3	1/10/73	43	0	1		
26	25	0	1	1	8/24/83	32	0	3		
27	26	0	2	1	2/8/89	27	0	3		
28	27	1	3	1	9/5/79	36	0	3		

input

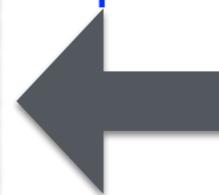
Data
Describer



comparison

Data
Generator

Model
Inspector



1	A	B	C	D	E	F	G	H		
1	UID	sex	race	MarriageStat	DateOfBirth	age	juv_fel	cour	decile	score
2	1	0	1	1	4/18/47	69	0	1		
3	2	0	2	1	1/22/82	34	0	3		
4	3	0	2	1	5/14/91	24	0	4		
5	4	0	2	1	1/21/93	23	0	8		
6	5	0	1	2	1/22/73	43	0	1		
7	6	0	1	3	8/22/71	44	0	1		
8	7	0	3	1	7/25/74	41	0	6		
9	8	0	1	2	2/25/73	49	0	4		
10	9	0	3	1	6/10/94	21	0	3		
11	10	0	3	1	6/10/88	27	0	4		
12	11	1	3	2	8/22/78	37	0	1		
13	12	0	2	1	12/2/74	41	0	4		
14	13	1	3	1	6/14/68	47	0	1		
15	14	0	2	1	3/25/85	31	0	3		
16	15	0	4	4	1/25/79	37	0	1		
17	16	0	2	1	6/22/90	25	0	10		
18	17	0	3	1	12/2/81	31	0	5		
19	18	0	3	1	5/8/85	31	0	3		
20	19	0	2	3	6/28/51	64	0	6		
21	20	0	2	1	11/29/94	21	0	9		
22	21	0	3	1	8/6/88	27	0	2		
23	22	1	3	1	3/22/95	21	0	4		
24	23	0	4	1	1/23/92	24	0	4		
25	24	0	3	3	1/10/73	43	0	1		
26	25	0	1	1	8/24/83	32	0	3		
27	26	0	2	1	2/8/89	27	0	3		
28	27	1	3	1	9/5/79	36	0	3		

output

Point 4

actionable transparency requires
interpretability

explain assumptions and effects, not details of
operation

engage the public - technical and non-technical

“Nutritional labels” for data and models

[K. Yang, J. Stoyanovich, A. Asudeh, B. Howe, HV Jagadish, G. Miklau; SIGMOD 2018]

Ranking Facts

Recipe			
Top 10:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	9.6	6.2
Faculty	122	52.5	45
GRE	800.0	796.3	771.9
Overall:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	2.9	1.4
Faculty	122	32.0	14
GRE	800.0	790.0	757.8

Stability



Slope at top-10: -6.91. Slope overall: -1.61.
Unstable when absolute value of slope of fit line in scatter plot <= 0.25 (slope threshold). Otherwise it is stable.

← Recipe

Attribute	Weight
PubCount	1.0
Faculty	1.0
GRE	1.0

Ingredients

Attribute	Correlation
PubCount	1.0
CSRankingAllArea	0.24
Faculty	0.12

Correlation strength is based on its absolute value. Correlation over 0.75 is high, between 0.25 and 0.75 is medium, under 0.25 is low.

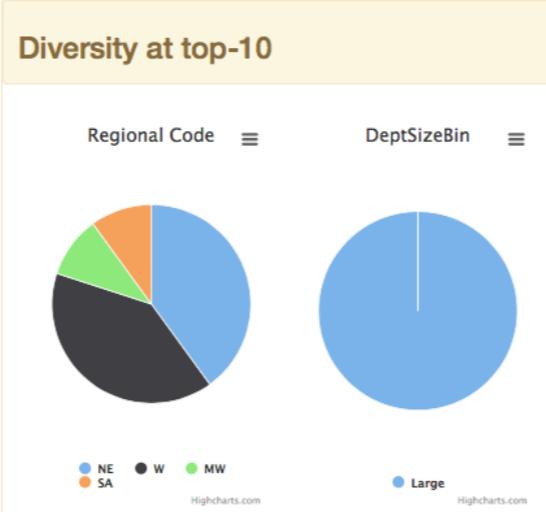
← Ingredients

Top 10:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	9.6	6.2
CSRankingAllArea	13	6.5	1
Faculty	122	52.5	45

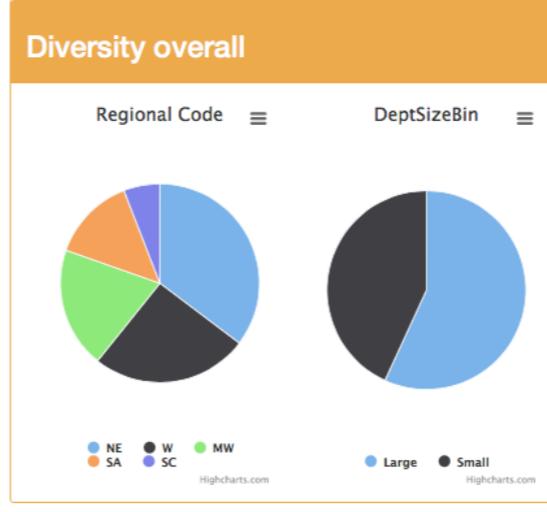
Overall:

Attribute	Maximum	Median	Minimum
PubCount	18.3	2.9	1.4
CSRankingAllArea	48	26.0	1
Faculty	122	32.0	14

Diversity at top-10



Diversity overall



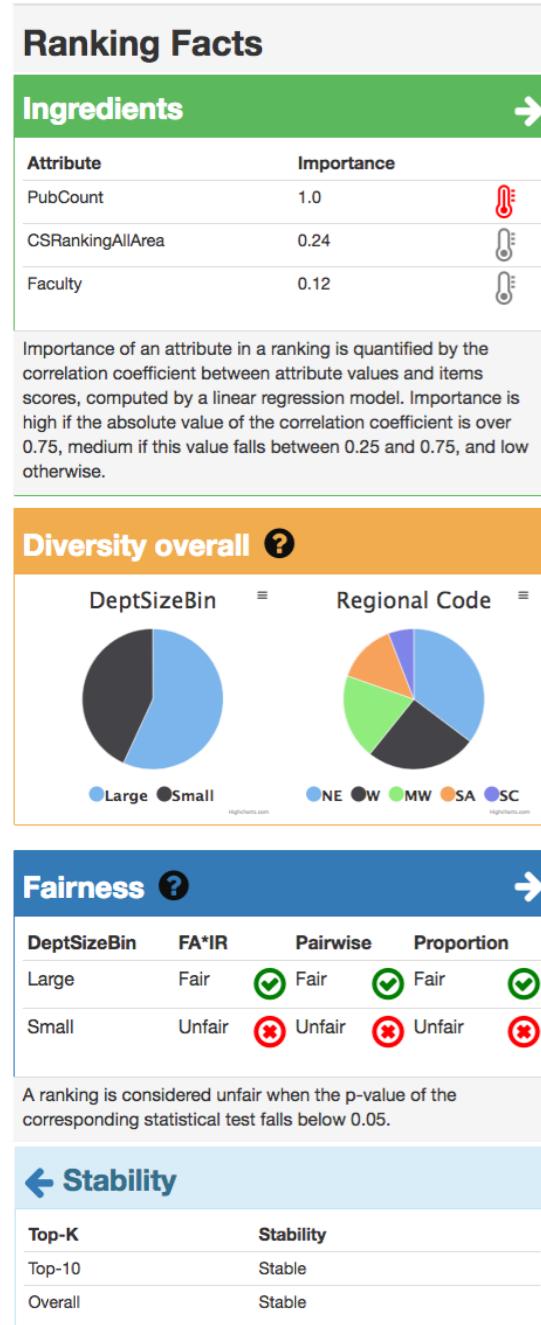
← Fairness

DeptSizeBin	FA*IR		Pairwise		Proportion	
	p-value	adjusted α	p-value	α	p-value	α
Large	1.0	0.87	0.99	0.05	1.0	0.05
Small	0.0	0.71	0.0	0.05	0.0	0.05

Top K = 26 in FA*IR and Proportion oracles. Setting of top K: In FA*IR and Proportion oracle, if N > 200, set top K = 100. Otherwise set top K = 50%N. Pairwise oracle takes whole ranking as input. FA*IR is computed as using code in [FA*IR codes](#). Proportion is implemented as statistical test 4.1.3 in [Proportion paper](#).

<http://demo.dataresponsibly.com/rankingfacts/>

Properties of a nutritional label



comprehensible: short, simple, clear

consultative: provide actionable info

comparable: implying a standard

concrete: helps determine a dataset's fitness for use for a given task

computable: produced as a “by-product” of computation - interpretability-by-design

Point 5

transparency / interpretability by design,
not as an afterthought

provision for transparency and interpretability at
every stage of the data lifecycle

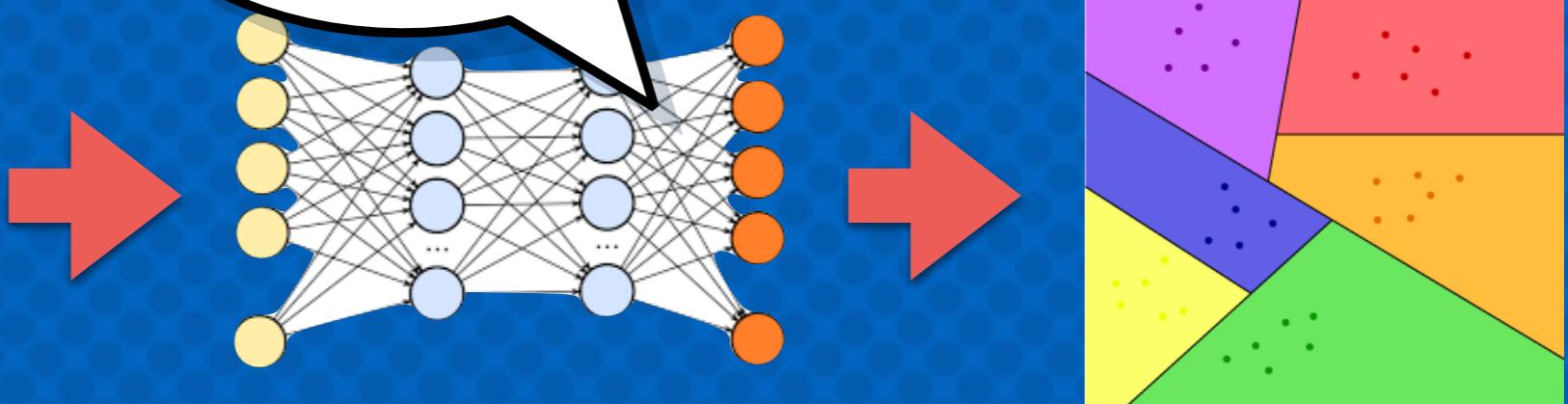
useful internally during development, for
communication and coordination between
agencies, and for accountability to the public

Frog's eye view

where did the data come from?

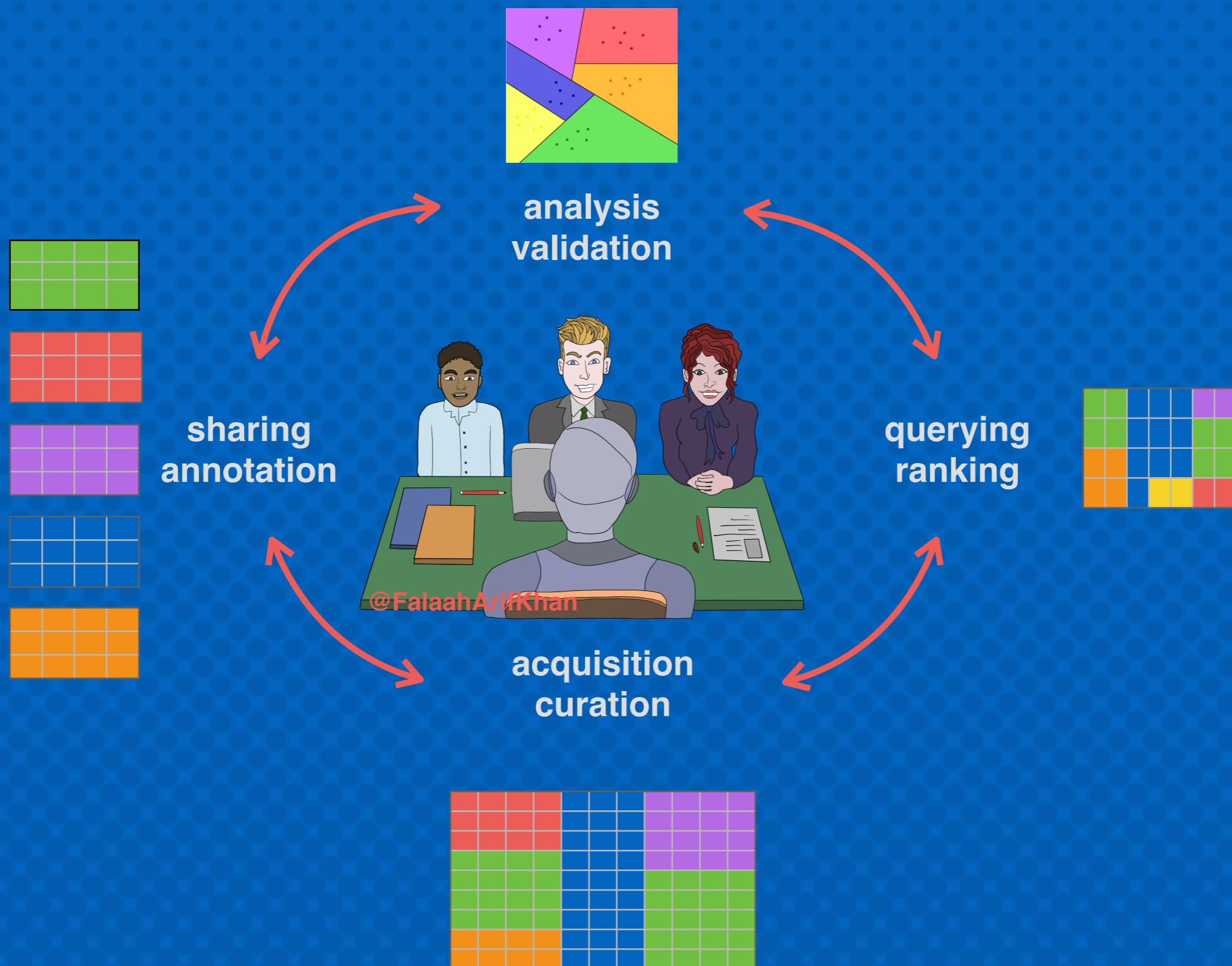
1	A	B	C	D	E	G	H
UID	sex	race	MarriageStat	DateOfBirth	age	iuv	fel
2	1	0	1	1	4/18/47	69	0
3	2	0	2	1	1/22/82	34	0
4	3	0	2	1	5/14/91	24	0
5	4	0	2	1	1/21/93	23	0
6	5	0	1	2	1/22/73	43	0
7	6	0	1	3	8/22/71	44	0
8	7	0	3	1	7/23/74	41	0
9	8	0	1	2	2/25/73	43	0
10	9	0	3	1	6/10/94	21	0
11	10	0	3	1	6/1/88	27	0
12	11	1	3	2	8/22/78	37	0
13	12	0	2	1	12/2/74	41	0
14	13	1	3	1	6/14/68	47	0
15	14	0	2	1	3/25/85	31	0
16	15	0	4	4	1/25/79	37	0
17	16	0	2	1	6/22/90	25	0
18	17	0	3	1	12/24/84	31	0
19	18	0	3	1	1/8/85	31	0
20	19	0	2	3	6/28/51	64	0
21	20	0	2	1	11/29/94	21	0
22	21	0	3	1	8/6/88	27	0
23	22	1	3	1	3/22/95	21	0
24	23	0	4	1	1/23/92	24	0
25	24	0	3	3	1/10/73	43	0
26	25	0	1	1	8/24/83	32	0
27	26	0	2	1	2/8/89	27	0
28	27	1	3	1	9/3/79	36	0
29	28	0	2	1	4/27/80	26	0

what happens inside the box?



how are results used?

Data lifecycle of an ADS



interpretability: in the
eye of the beholder

What are we explaining?

[J. Stoyanovich, J. Van Bavel, T. West; *NMI 2020*]

process (same for everyone? **why** is this the process?) vs. outcome

procedural justice aims to ensure that algorithms are perceived as fair and legitimate

data transparency is unique to algorithm-assisted decision-making, relates to the justification dimension of interpretability

To whom are we explaining and why?

[J. Stoyanovich, J. Van Bavel, T. West; *NMI 2020*]

accounting for the needs of different stakeholders

social identity - people trust their in-group members more

moral cognition - is a decision or outcome morally right or wrong?

How do we know that we explained well?

[J. Stoyanovich, J. Van Bavel, T. West; *NMI 2020*]

nutritional labels! :)

... but do they work?