

DS-GA 3001.009: Responsible Data Science

Algorithmic Fairness (continued)

Prof. Julia Stoyanovich
Center for Data Science
Computer Science and Engineering at Tandon

@stoyanoj

<http://stoyanovich.org/>
<https://dataresponsibly.github.io/>

Two notions of fairness

individual fairness



equality

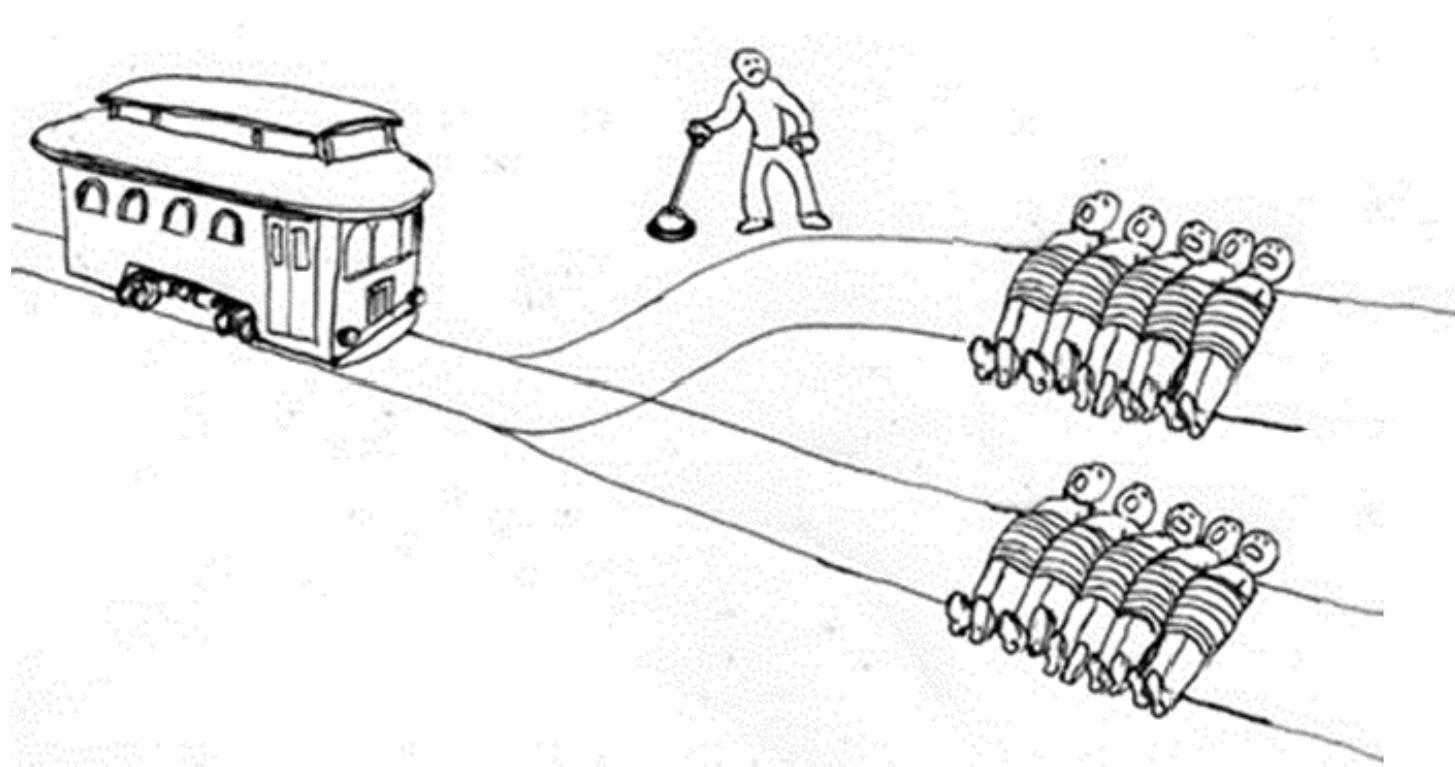
group fairness



equity

two intrinsically different world views

Fairness definitions as “trolley problems”



https://www.helpage.org/silo/images/blogs/16_1391611056.gif

Fairness in risk assessment

- A risk assessment tool **gives a probability estimate of a future outcome**
- Used in many domains:
 - insurance, criminal sentencing, medical testing, hiring, banking
 - also in less-obvious set-ups, like online advertising
- **Fairness** is concerned with **how different kinds of errors are distributed among sub-populations**
 - Recall our discussion on fairness in classification - similar?

Racial bias in criminal sentencing

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

A commercial tool **COMPAS** automatically predicts some categories of future crime to assist in bail and sentencing decisions. It is used in courts in the US.

Prediction Fails Differently for Black Defendants

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

Overall, Northpointe's assessment tool correctly predicts recidivism 61 percent of the time. But blacks are almost twice as likely as whites to be labeled a higher risk but not actually re-offend. It makes the opposite mistake among whites: They are much more likely than blacks to be labeled lower risk but go on to commit other crimes. (Source: ProPublica analysis of data from Broward County, Fla.)

<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Desirable properties of risk tools

[J. Kleinberg, S. Mullainathan, M. Raghavan; ITCS (2017)]

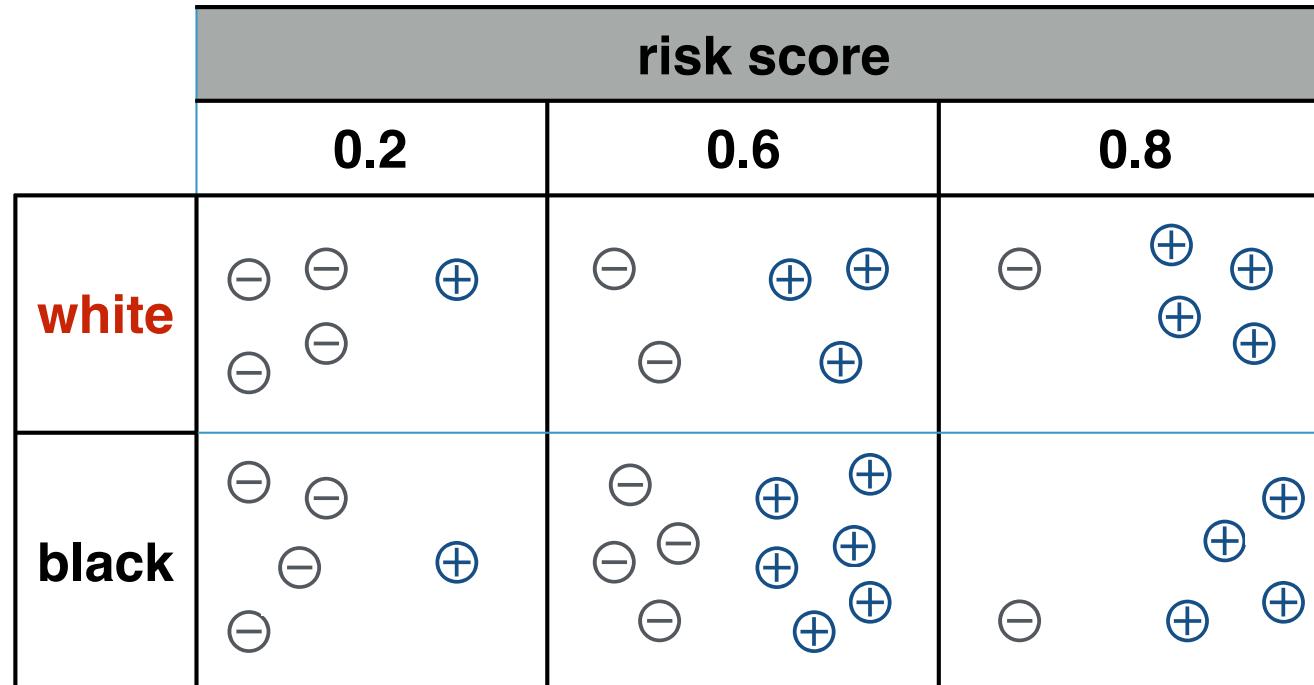
“risk assessment tool / instrument” = “**risk tool / instrument**”
for brevity in the rest of today’s slides

- Calibration
- Balance for the positive class
- Balance for the negative class

can we have all these properties?

Calibration

**positive
outcomes:
do recidivate**



given the output of a risk tool, likelihood of belonging to the positive class is independent of group membership

0.6 means 0.6 for any defendant - likelihood of recidivism

why do we want calibration?

Calibration in COMPAS

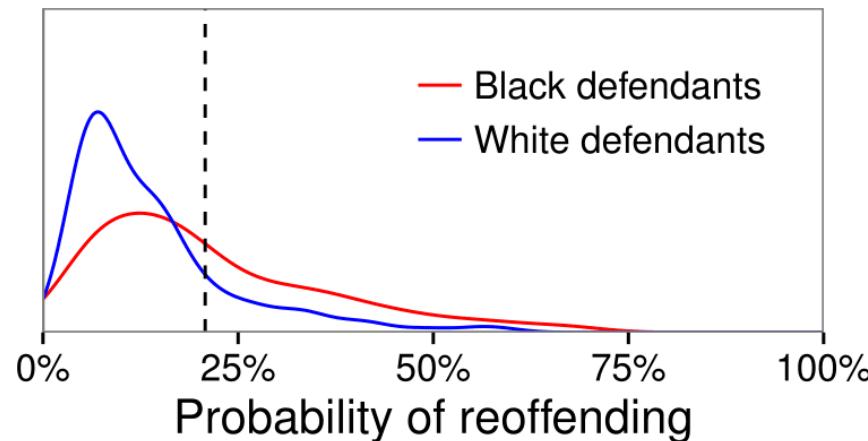
[J. Kleinberg, S. Mullainathan, M. Raghavan; *ITCS 2017*]

Predictive parity (also called **calibration**)

an risk tool identifies a set of instances as having probability x of constituting positive instances, then approximately an x fraction of this set are indeed positive instances, over-all and in sub-populations

COMPAS is **well-calibrated**: in the window around 40%, the fraction of defendants who were re-arrested is ~40%, both over-all and per group.

Broward County



[plot from Corbett-Davies et al.; *KDD 2017*]

Balance

[J. Kleinberg, S. Mullainathan, M. Raghavan; *ITCS 2017*]

- Balance for the positive class: Positive instances are those who go on to re-offend. The average score of positive instances should be the same across groups.
- Balance for the negative class: Negative instances are those who do not go on to re-offend. The average score of negative instances should be the same across groups.
- Generalization of: Both groups should have equal false positive rates and equal false negative rates.
- Different from statistical parity!

the chance of making a mistake does not depend on race

Desiderata, formally

[J. Kleinberg, S. Mullainathan, M. Raghavan; ITCS (2017)]

- For each group, a v_b fraction in each bin b is positive
- Average score of positive class same across groups
- Average score of negative class same across groups

can we have all these properties?

Achievable only in trivial cases

[J. Kleinberg, S. Mullainathan, M. Raghavan; ITCS (2017)]

- Perfect information: the tool knows who recidivates (score 1) and who does not (score 0)
- Equal base rates: the fraction of positive-class people is the same for both groups

cannot even find a good approximate solution

a negative result, need tradeoffs

proof sketched out in (starts 12 min in)

<https://www.youtube.com/watch?v=UUC8tMNxwV8>

Group fairness impossibility result

[A. Chouldechova; arXiv:1610.07524v1 (2017)]

If a predictive instrument **satisfies predictive parity**, but the **prevalence** of the phenomenon **differs between groups**, then the instrument **cannot achieve** equal false positive rates and equal false negative rates across these groups

Recidivism rates in the ProPublica dataset are higher for the black group than for the white group

<https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>

What is recidivism?: Northpointe [*the maker of COMPAS*] defined recidivism as “**a finger-printable arrest** involving a charge and a filing for any uniform crime reporting (UCR) code.”

Fairness for whom?

Decision-maker: of those I've labeled high-risk, how many will recidivate?

Defendant: how likely am I to be incorrectly classified high-risk?

Society: (think positive interventions) is the selected set demographically balanced?

based on a slide by Arvind Narayanan

	labeled low-risk	labeled high-risk
did not recidivate	TN	FP
recidivated	FN	TP

different metrics matter to different stakeholders

<https://www.propublica.org/article/propublica-responds-to-companys-critique-of-machine-bias-story>

Impossibility theorem

Metric	Equalized under
Selection probability	Demographic parity
Pos. predictive value	Predictive parity
Neg. predictive value	
False positive rate	Error rate balance
False negative rate	Error rate balance
Accuracy	Accuracy equity

based on a slide by Arvind Narayanan

Chouldechova
paper

All these metrics can be expressed in terms of FP, FN, TP, TN

If these metrics are equal for 2 groups, some trivial algebra shows that the prevalence (in the COMPAS example, of recidivism, as measured by re-arrest) is also the same for 2 groups

Nothing special about these metrics, can pick any 3!

Ways to evaluate binary classifiers

based on a slide by Arvind Narayanan

	True condition				
	Total population	Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	True positive , Power	False positive , Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative , Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
	True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$	$F_1 \text{ score} = \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$
	False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	True negative rate (TNR), Specificity (SPC) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$		

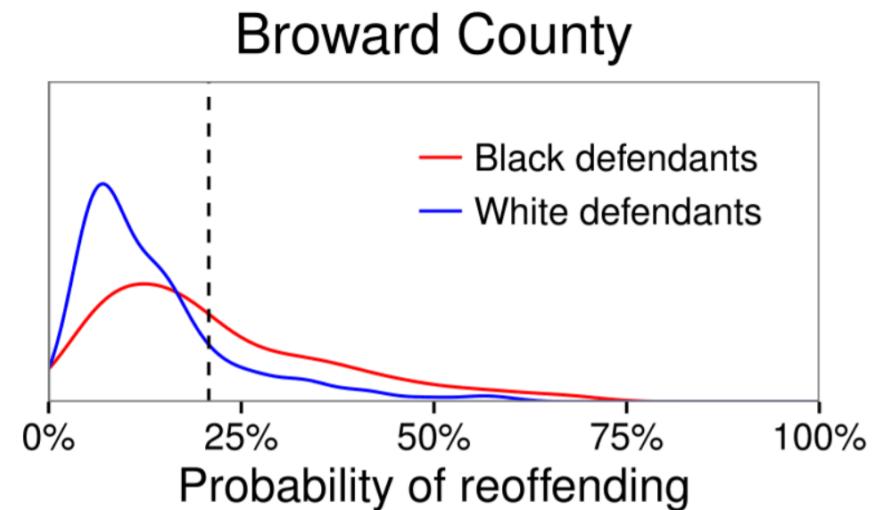
364 impossibility theorems :)

Individual fairness

based slides by Arvind Narayanan

Individual fairness:

assuming scores are calibrated, we cannot pick a single threshold for 2 groups that equalizes both the False Positives Rate and the False Negatives Rate



What's the right answer?

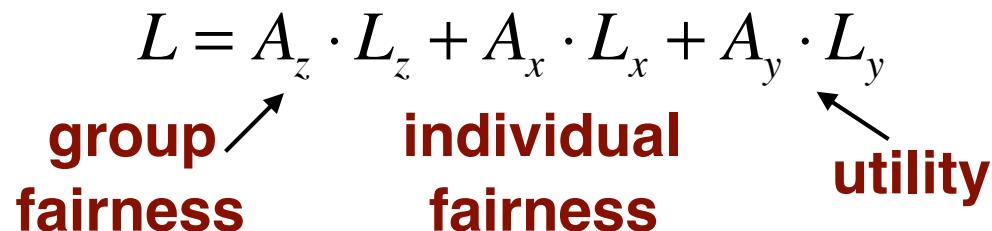
There is no single answer!

Need transparency and public debate

- Consider harms and benefits to different stakeholders
- Being transparent about which fairness criteria we use, how we trade them off
- Recall “Learning Fair Representations”: a typical ML approach

$$L = A_z \cdot L_z + A_x \cdot L_x + A_y \cdot L_y$$

group fairness individual fairness utility



apples + oranges + fairness = ?

Evaluating fairness-aware algorithms

[S. Friedler, C. Scheidegger, S. Venkatsubramanian, S. Chaudhary,
E. Hamilton, D. Roth; FAT* (2019)]

How do we know how to trade off different fairness objectives, and how to encode them in fairness-aware algorithms? - Societal context + experimental work!

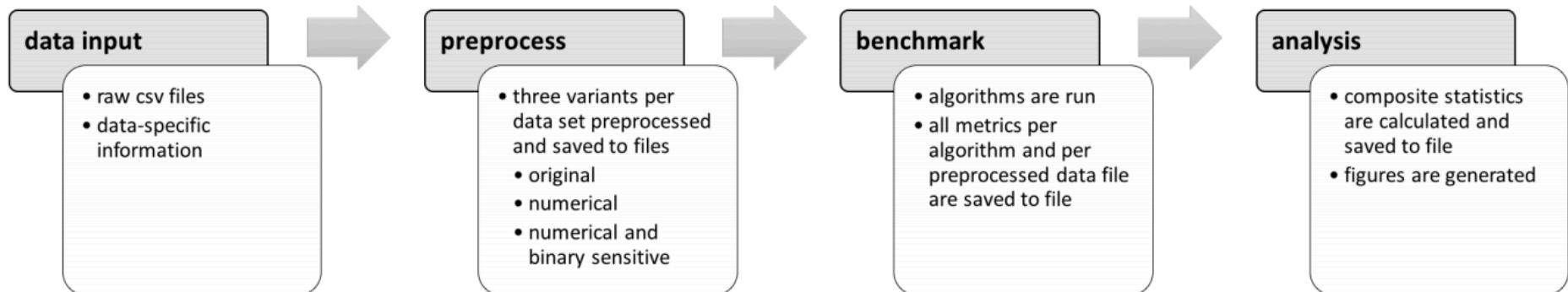


Figure 1: The stages of the fairness-aware benchmarking program: data input, preprocessing, benchmarking, and analysis. Intermediate files are saved at each stage of the pipeline to ensure reproducibility.

Insight 1: pre-processing matters

[S. Friedler, C. Scheidegger, S. Venkatsubramanian, S. Chaudhary, E. Hamilton, D. Roth; FAT* (2019)]

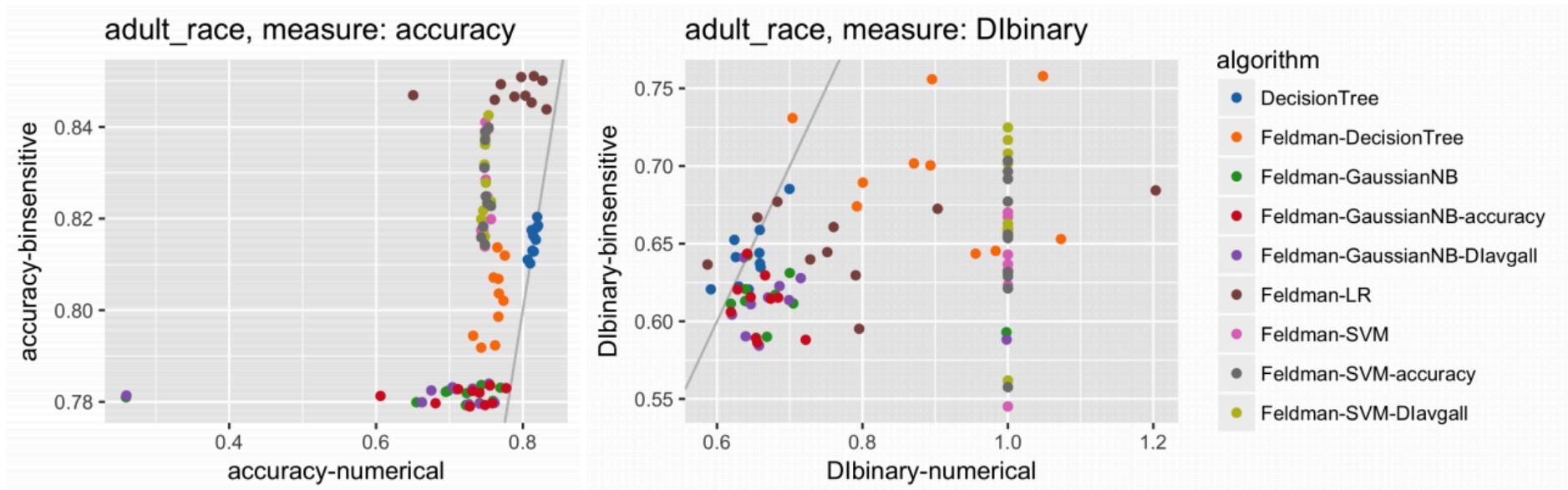
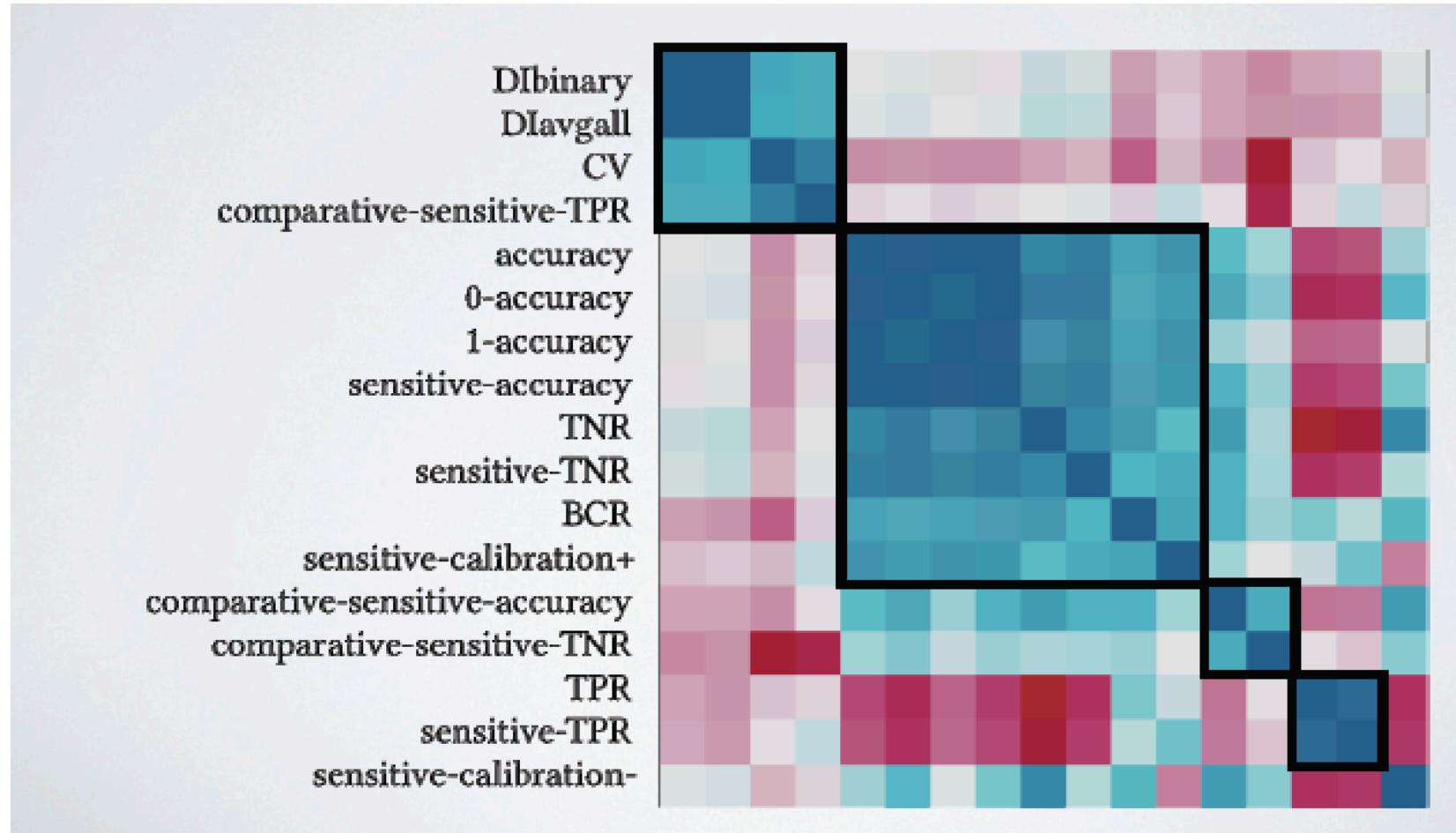


Figure 2: Examining the results of the Feldman et al. [10] algorithm under different preprocessing choices: numerical versus numerical+binary. Each dot plots the result of a single split of the data in terms of the labeled metric under both preprocessing choices. The gray line shows equality between the preprocessing choices. The model used within the Feldman algorithm is listed, and some variants of the algorithm had the tradeoff parameter optimized for either accuracy or disparate impact value.

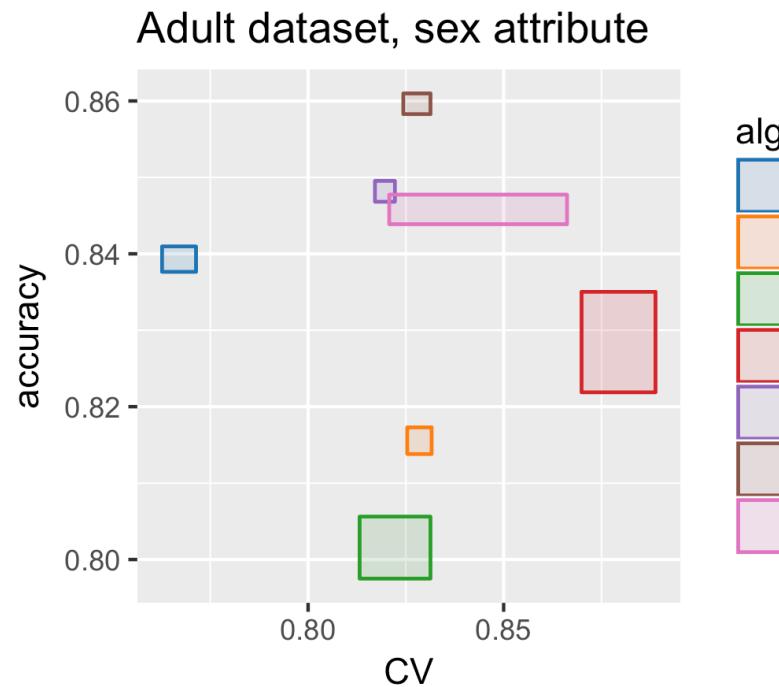
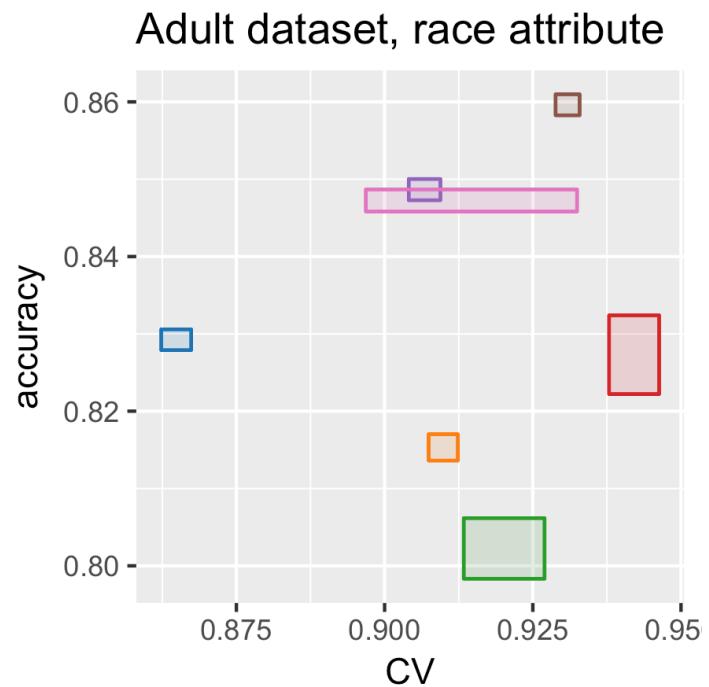
Insight 2: some measures correlate

[S. Friedler, C. Scheidegger, S. Venkatsubramanian, S. Chaudhary,
E. Hamilton, D. Roth; FAT* (2019)]



Insight 3: beware of variability

[S. Friedler, C. Scheidegger, S. Venkatsubramanian, S. Chaudhary, E. Hamilton, D. Roth; FAT* (2019)]



Feldman et al. varies in accuracy over splits while Zafar et al. varies in fairness.

Causal interpretations of fairness

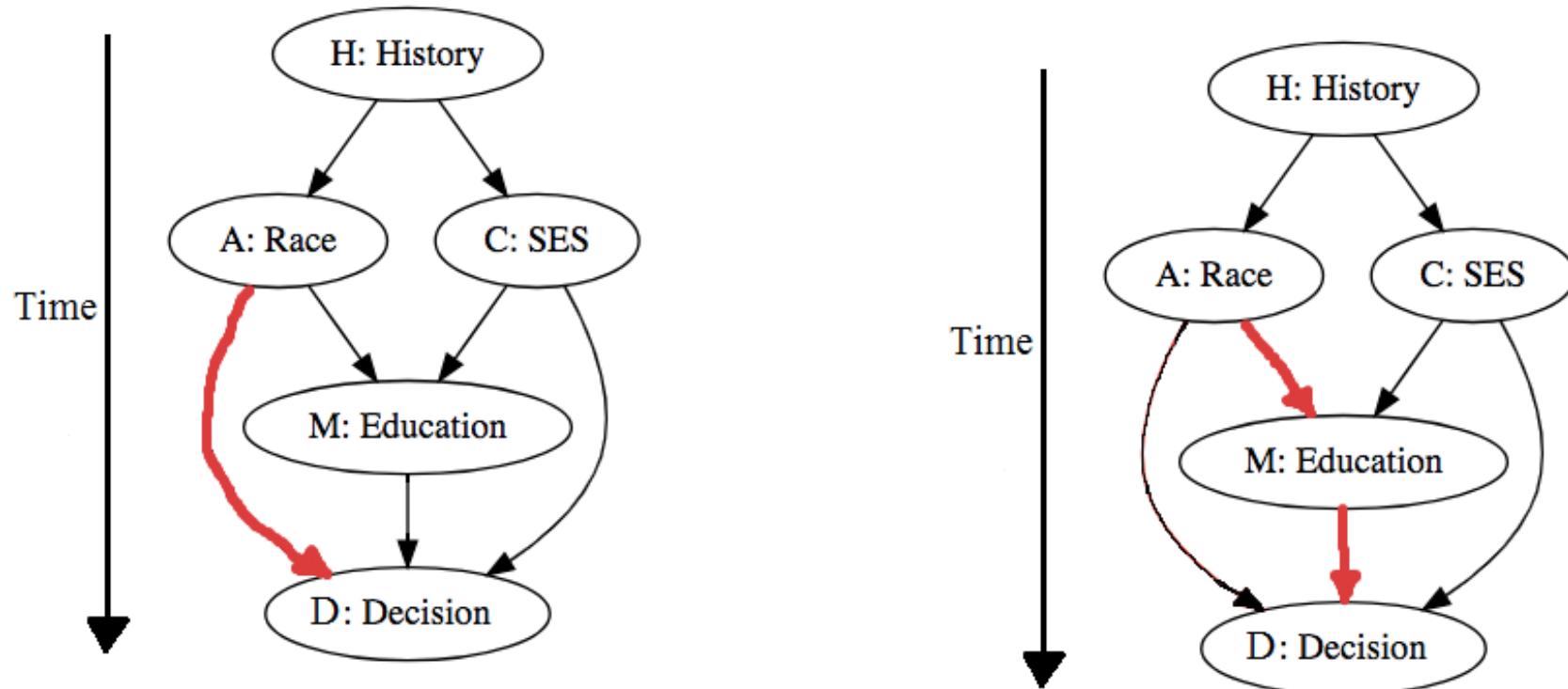
<https://shiraamitchell.github.io/fairness/>

- Will be covered in a **guest lecture by Shira Mitchell on Thursday, February 21** (no lecture that week, so we'll use lab time as lecture time)
- Starting with **counterfactuals**:
 - *Was I not hired because I was black? => Would I have been hired if I were non-black?*
 - *Is there an effect of race on hiring? => Would the rate of hiring be the same if everyone were black? If no-one were?*

Causal interpretations of fairness

[T.J. VanderWeele and W.R. Robinson; Epidemiology (2014)]

arrows represent possible causal relationships



we (society) decide which of these are “OK”

Fairness in ranking



Fairness in ranking

[K. Yang & J. Stoyanovich, FATML (2016)]

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

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

Output: permutation of the items (complete or top-k)

<u>id</u>	sex	race	age	cat
a	F	W	25	T
b	F	B	23	S
c	M	W	27	T
d	M	B	45	S
e	M	W	60	U



What is a positive outcome in a ranking?

Idea: Rankings are relative, fairness measures should be rank-aware

The order of things

THE NEW YORKER

1. Chevrolet Corvette 205

2. Lotus Evora 195

3. Porsche Cayman 195

1. Lotus Evora 205

2. Porsche Cayman 198

3. Chevrolet Corvette 192

1. Porsche Cayman 193

2. Chevrolet Corvette 186

3. Lotus Evora 182

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!

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.

Location-location-location

[K. Yang & J. Stoyanovich, FATML (2016)]

gender is the sensitive attribute, input is balanced

Algorithm 1 Ranking generator

Require: Ranking τ , fairness probability f .

{Initialize the output ranking σ }.

```
1:  $\sigma \leftarrow \emptyset$ 
2:  $\tau^+ = \tau \cap S^+$ 
3:  $\tau^- = \tau \cap S^-$ 
4: while  $(\tau^+ \neq \emptyset) \wedge (\tau^- \neq \emptyset)$  do
5:    $p = \text{random}([0, 1])$ 
6:   if  $p < f$  then
7:     Pop an item from the top of the list  $\tau^+$ .
8:      $\sigma \leftarrow \text{pop}(\tau^+)$ 
9:   else
10:    Pop an item from the top of the list  $\tau^-$ .
11:     $\sigma \leftarrow \text{pop}(\tau^-)$ 
12:  end if
13: end while
14:  $\sigma \leftarrow \tau^+$ 
15:  $\sigma \leftarrow \tau^-$ 
16: return  $\sigma$ 
```

rank	gender	rank	gender	rank	gender
1	M	1	M	1	M
2	M	2	M	2	M
3	M	3	F	3	F
4	M	4	M	4	F
5	M	5	M	5	M
6	F	6	F	6	F
7	F	7	M	7	M
8	F	8	F	8	F
9	F	9	F	9	M
10	F	10	F	10	F

$$f = 0$$

$$f = 0.3$$

$$f = 0.5$$

Rank-aware fairness

[K. Yang & J. Stoyanovich, FATML (2016)]

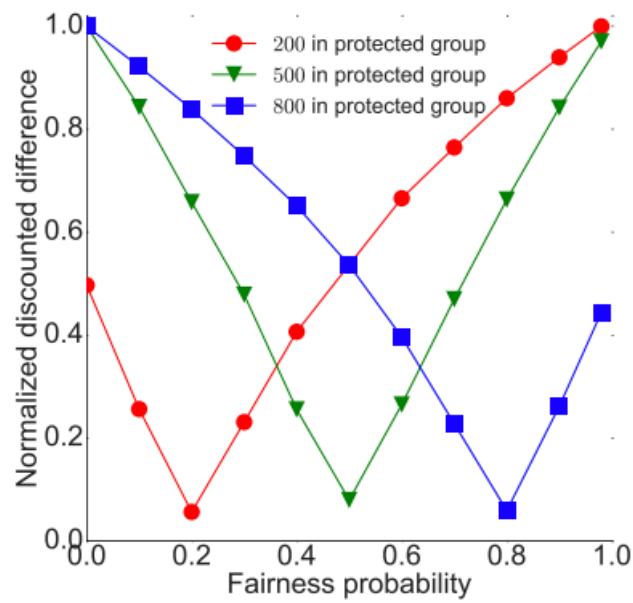


Figure 3: rND on 1,000 items

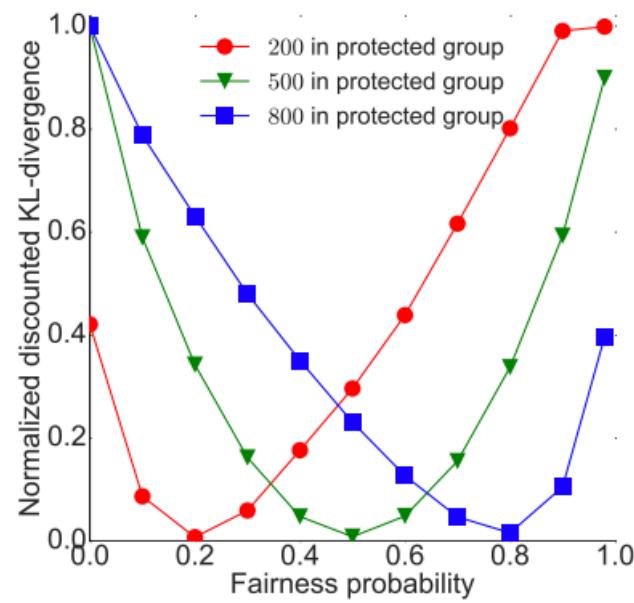


Figure 4: rKL on 1,000 items

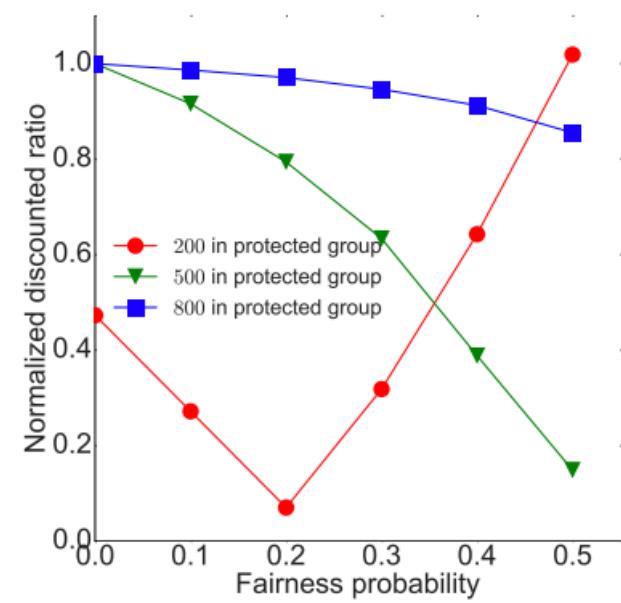


Figure 5: rRD on 1,000 items

In an optimization framework

[K. Yang & J. Stoyanovich, FATML (2016)]

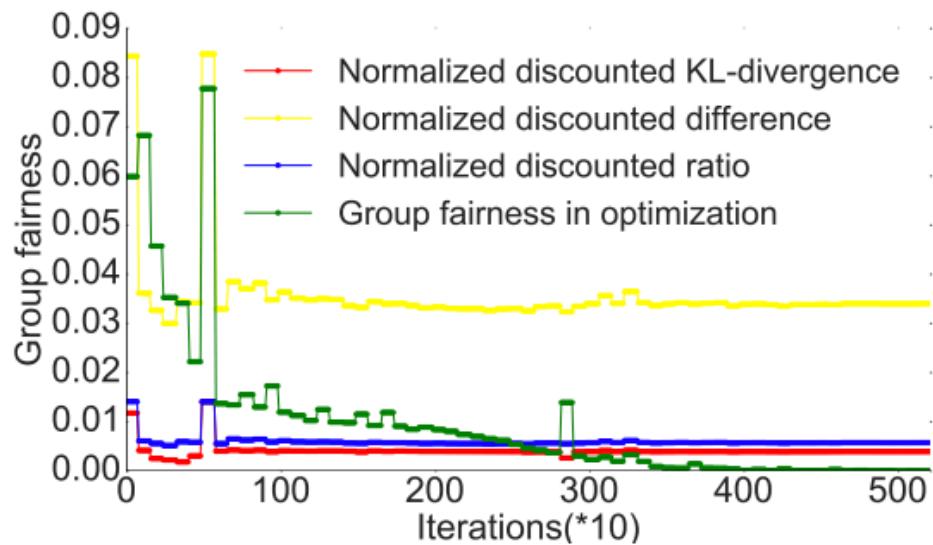
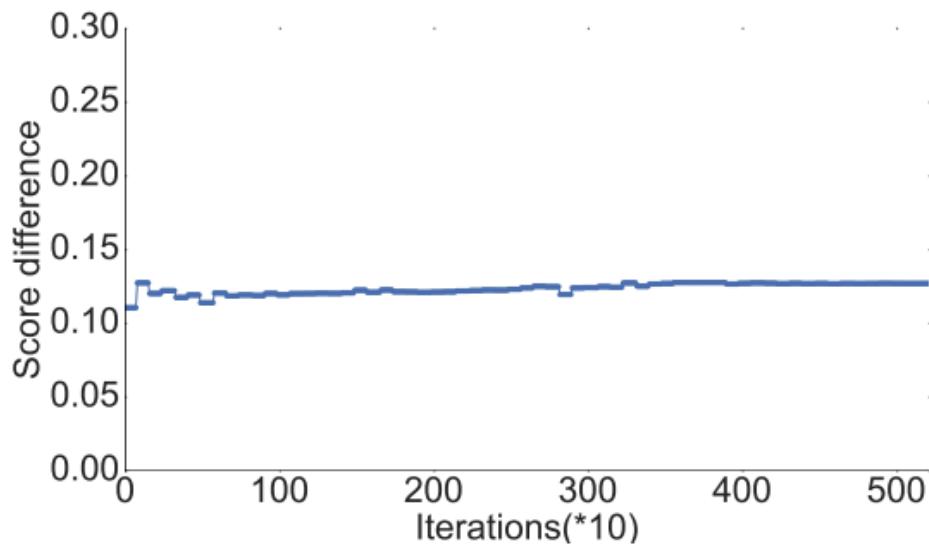
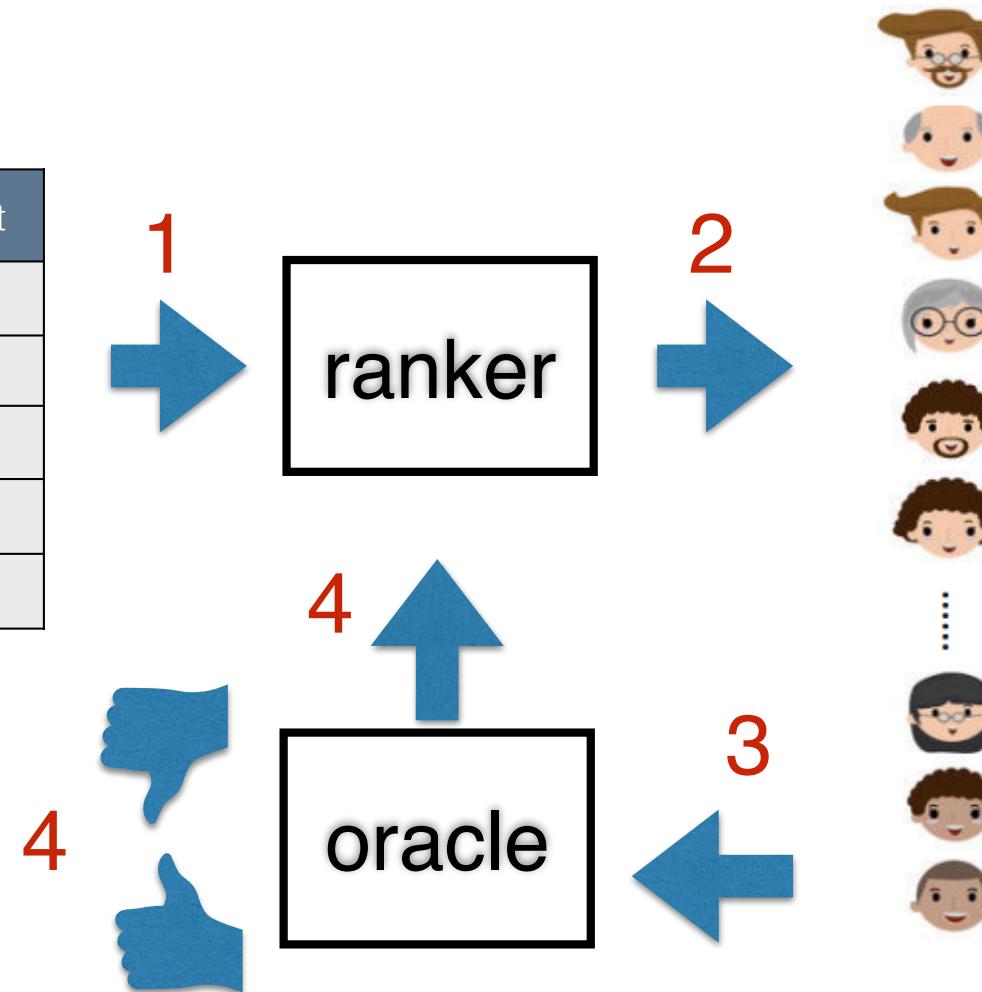


Figure 6: Accuracy and fairness on German Credit, ranked by sum of normalized attribute values, with $k = 10$.

Designing fair rankers

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

<u>id</u>	sex	race	age	cat
a	F	W	25	T
b	F	B	23	S
c	M	W	27	T
d	M	B	45	S
e	M	W	60	U



Score-based rankers

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

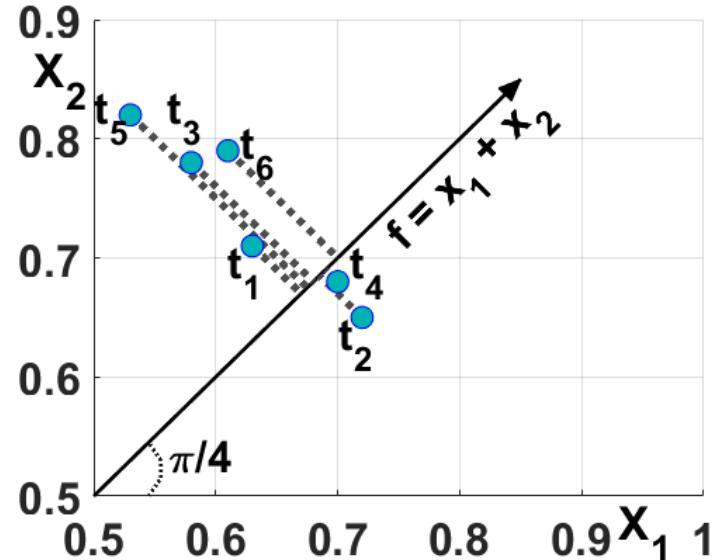
- tuple x in D ; $\text{score}(x)$: sum of attribute values, with non-negative weight (a common special case of **monotone aggregation**)
- weights **subjectively chosen by a user**: $0.5g + 0.5s$, where g - normalized GRE, s - normalized SAT; why not $0.45g + 0.55s$?

\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

Geometry of a (2D) ranker

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

\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



- tuples are points in 2D, scoring functions are rays starting from the origin
- to determine a ranking of the points, we read it off, walking the ray of the scoring function towards the origin
- examples: $f(x) = x_1 + x_2$; $f(x) = x_1$; $f(x) = x_2$

Goal: find a satisfactory function

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

Closest Satisfactory Function: *Given a dataset \mathcal{D} with n items over d scalar scoring attributes, a fairness oracle $O : \nabla_f(\mathcal{D}) \rightarrow \{\top, \perp\}$, and a linear scoring function f with the weight vector $\vec{w} = \langle w_1, w_2, \dots, w_d \rangle$, find the function f' with the weight vector \vec{w}' such that $O(\nabla_{f'}(\mathcal{D})) = \top$ and the angular distance between \vec{w} and \vec{w}' is minimized.*

How might we approach this? Why is this difficult?

Ordering exchange

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

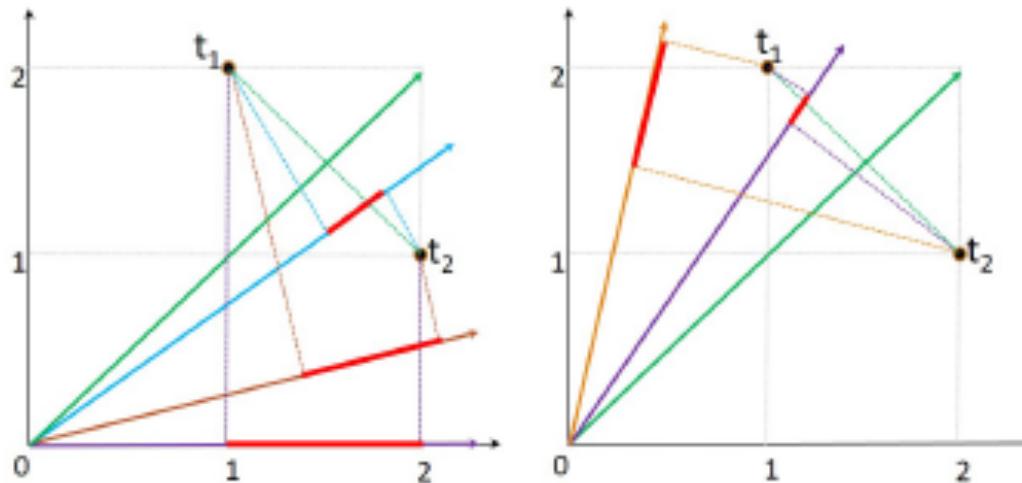
Key idea: only look at scoring functions that change the relative order between some pair of tuples (points). These are the only points where the fairness oracle may change its mind!

$$t_1 \langle 1, 2 \rangle \quad t_2 \langle 2, 1 \rangle$$

$$t_2 \succ_x t_1$$

$$t_2 =_{x+y} t_1$$

$$t_2 \prec_y t_1$$



An **ordering exchange** is a set of functions that score a pair of points equally. In 2D, it corresponds to a single function.

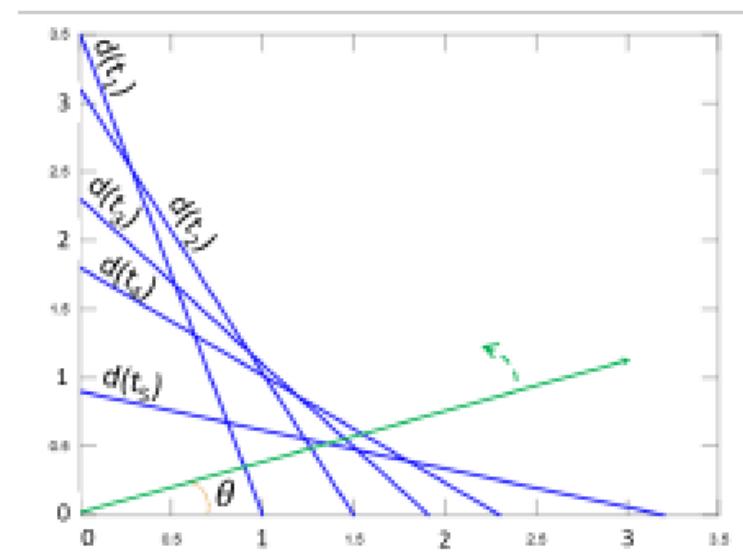
Outline of approach

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

Pre-processing

- Transform original space into dual space (in 2D, points become lines)
- Order the items on $f=x$; compute ordering exchanges between adjacent pairs of items
- Sweep the space with a ray from x-axis to y-axis, find satisfactory regions

t_1	1	3.5
t_2	1.5	3.1
t_3	1.91	2.3
t_4	2.3	1.8
t_5	3.2	0.9

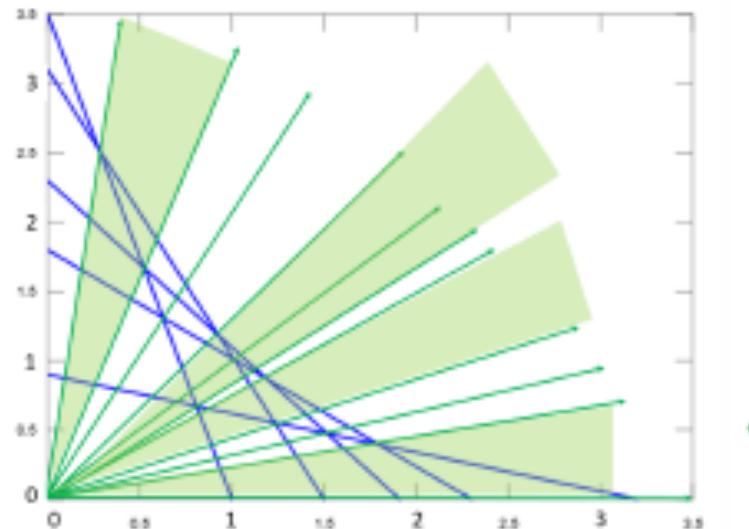
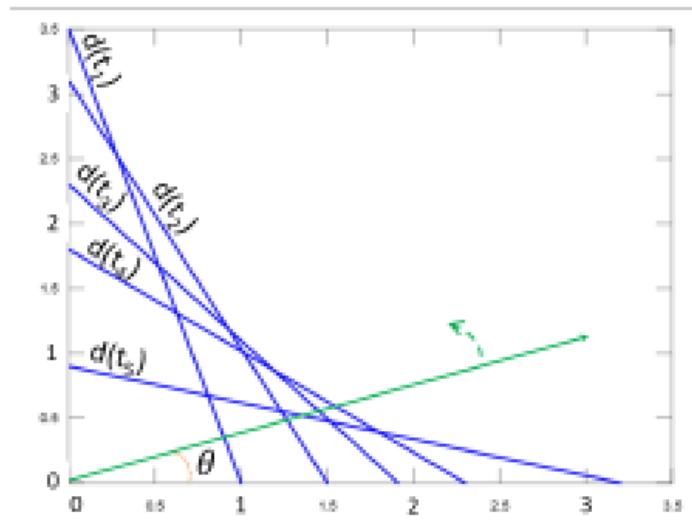
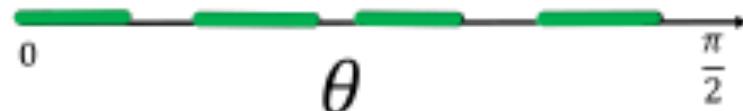


Outline of approach

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

At query time

- Look for a satisfactory region closest to the query function
- In 2D, this is simply binary search



And lots more algorithmic + systems work

[A. Asudeh, HV Jagadish, J. Stoyanovich, G. Das; ACM SIGMOD (2019)]

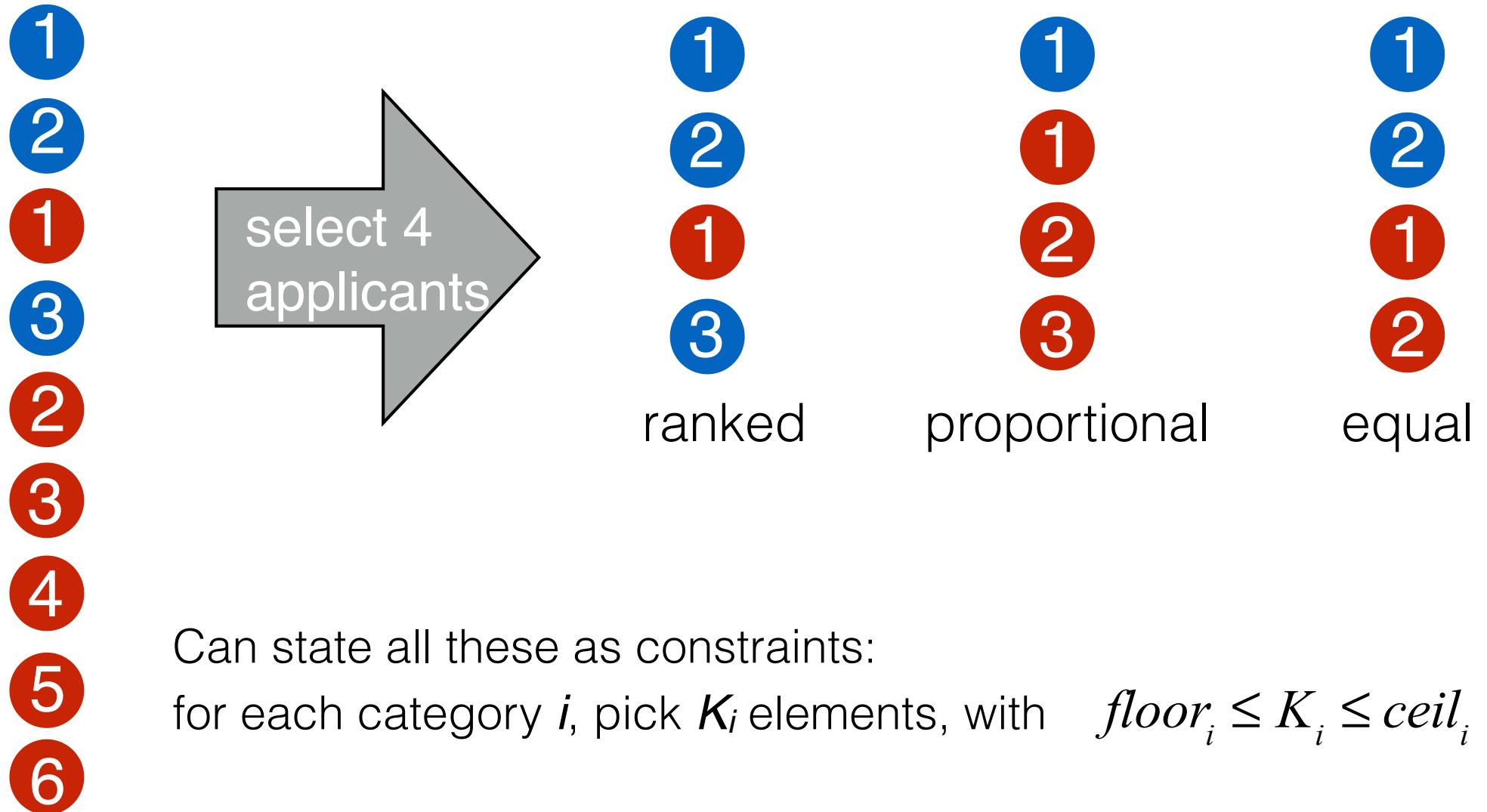
- Multi-dimensional indexing methods “arrangement construction”
- Sampling of items (does work), sampling of functions (doesn’t work) to speed up index construction
- Experiments on COMPAS and on US Department of Transportation (DOT) - flights / airlines

Follow-up work on designing fair ranking functions

Diversity

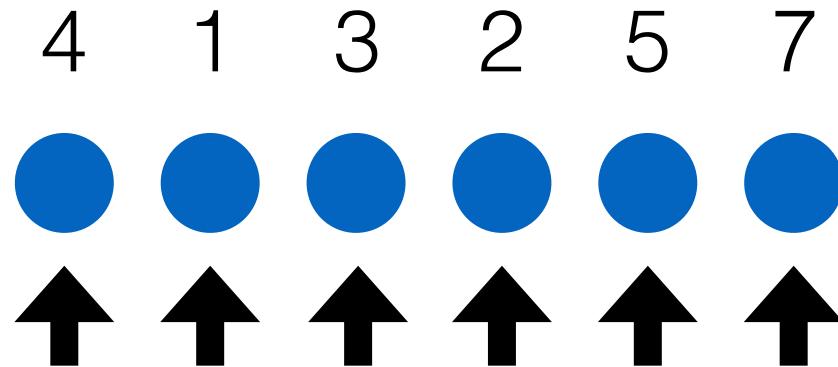


Job applicant selection



Hiring a job candidate

Goal: Hire a candidate with a high score



Candidates arrive one-by-one

A candidate's score is revealed when the candidate arrives

Decision to accept or reject a candidate made on the spot

The Secretary Problem

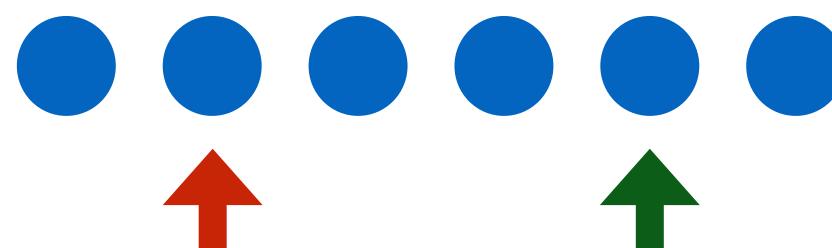
Goal: Design an algorithm for picking **one** element of a **randomly ordered sequence**, to maximize the probability of picking the **maximum element** of the entire sequence.

$$N = 6$$

$$S = \left\lfloor \frac{N}{e} \right\rfloor = 2$$

$$T = 4$$

4 1 3 2 5 7



Competitive ratio

$$\frac{1}{e}$$

the best possible!

Consider, and reject, the first S candidates

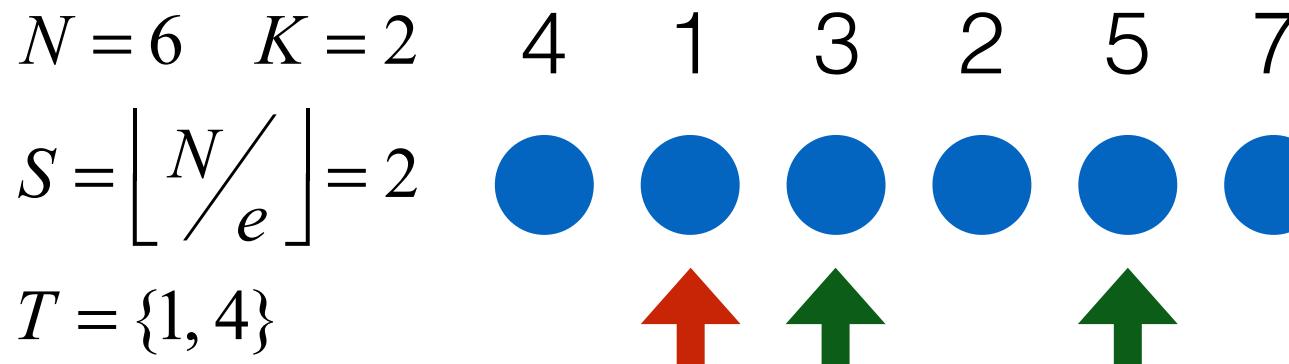
Record T , the best seen score among the first S candidates

Accept the next candidate with score better than T

K-choice Secretary

[Babaioff et al., 2007]

Goal: Design an algorithm for picking K elements of a randomly ordered sequence, to maximize their **expected sum**.



Competitive ratio

$$\frac{1}{e}$$

far from optimal

Consider, and reject, the first S candidates

Record K best scores among the first S candidates, call this T

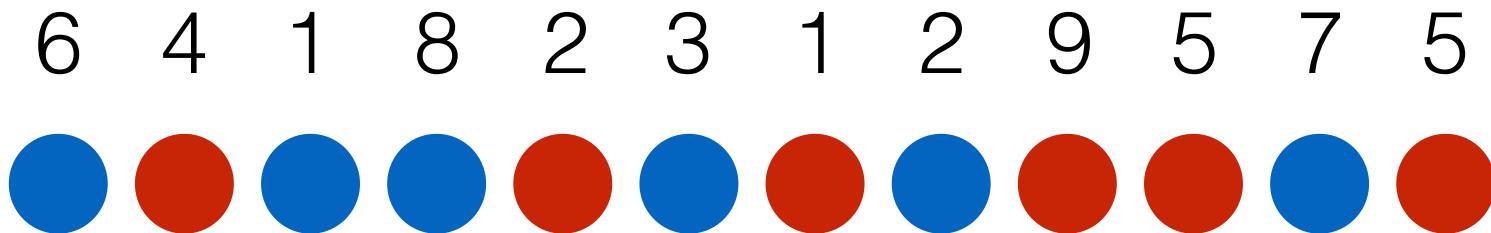
Whenever a candidate arrives whose score is higher than the minimum in T , accept the candidate and delete the minimum from T

K-choice Secretary

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

Goal: Design an algorithm for picking K elements of a randomly ordered sequence, to maximize their expected sum.

For each category i , pick K_i elements, with $\text{floor}_i \leq K_i \leq \text{ceil}_i$



$$N_{red} = N_{blue} = 6$$

$$K = 3$$

$$1 \leq K_{red}, K_{blue} \leq 2$$

Accept floor items for each category from per-category streams

$$\textit{slack} = K - (\text{floor}_{red} + \text{floor}_{blue})$$

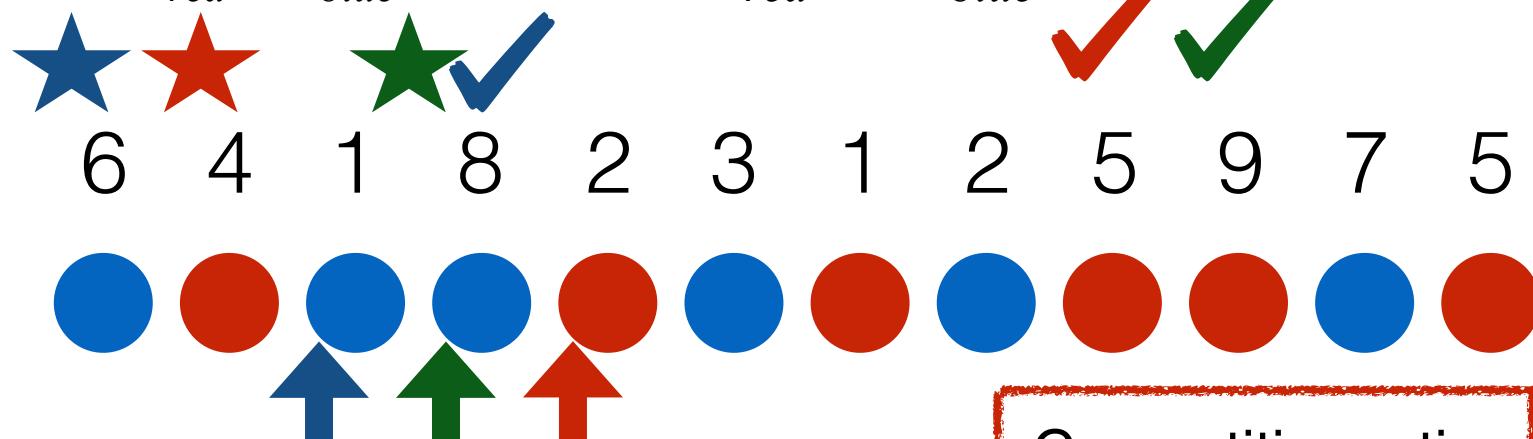
Accept the remaining \textit{slack} items irrespective of category membership, but subject to \textit{ceil}

Diverse K-choice Secretary

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

$$N_{red} = N_{blue} = 6$$

$$K = 3 \quad 1 \leq K_{red}, K_{blue} \leq 2$$



$$slack = 1$$

$$S_{red} = S_{blue} = 2 \quad S = 4$$



Competitive ratio

$$\frac{1}{e}$$

far from optimal

Adding a deferred list

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

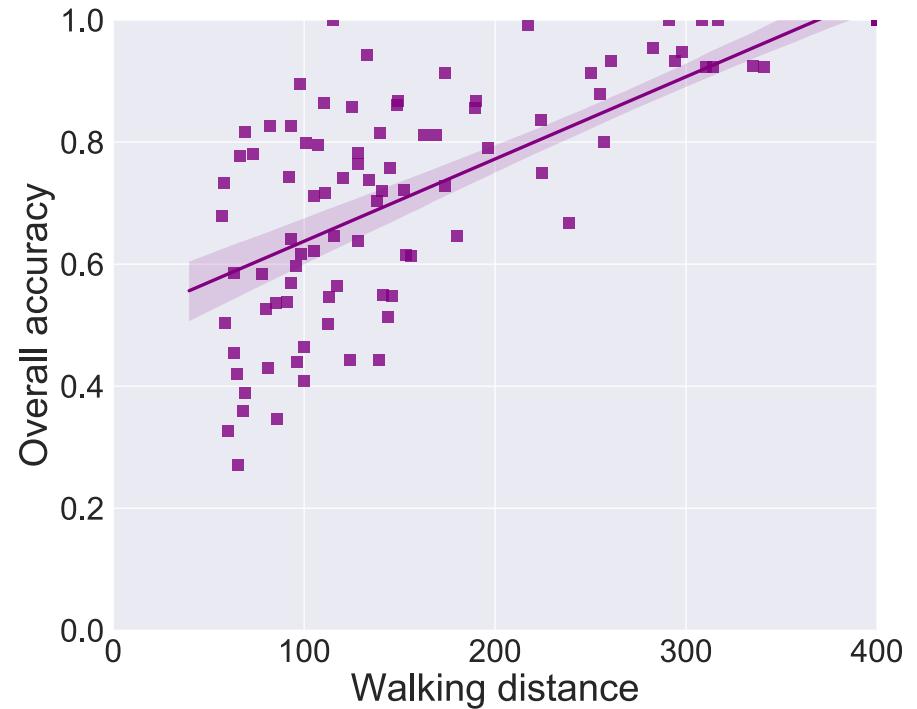
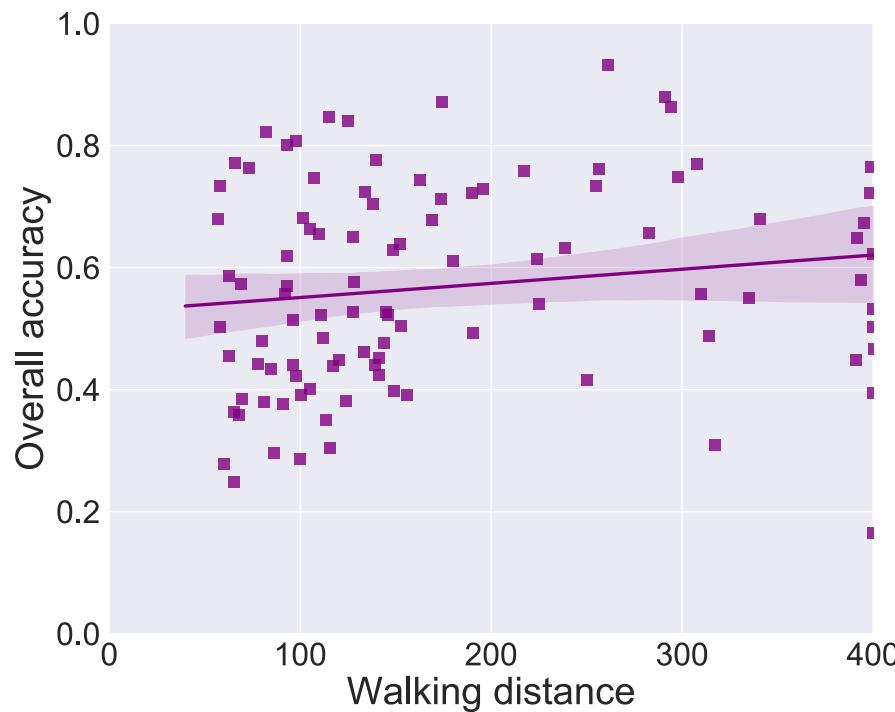
- An improvement on Diverse K-choice Secretary
- Do not immediately reject or accept items: keep a deferred list D_i per category i of size up to ceil_i
- Stop reading the input, post-warm up, once all floor_i constraints are met, and once there are K items in the union of the deferred lists
- Main advantage: often avoids reading items from the end of the stream

Diversity is achievable

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

deferred list

with deferred list

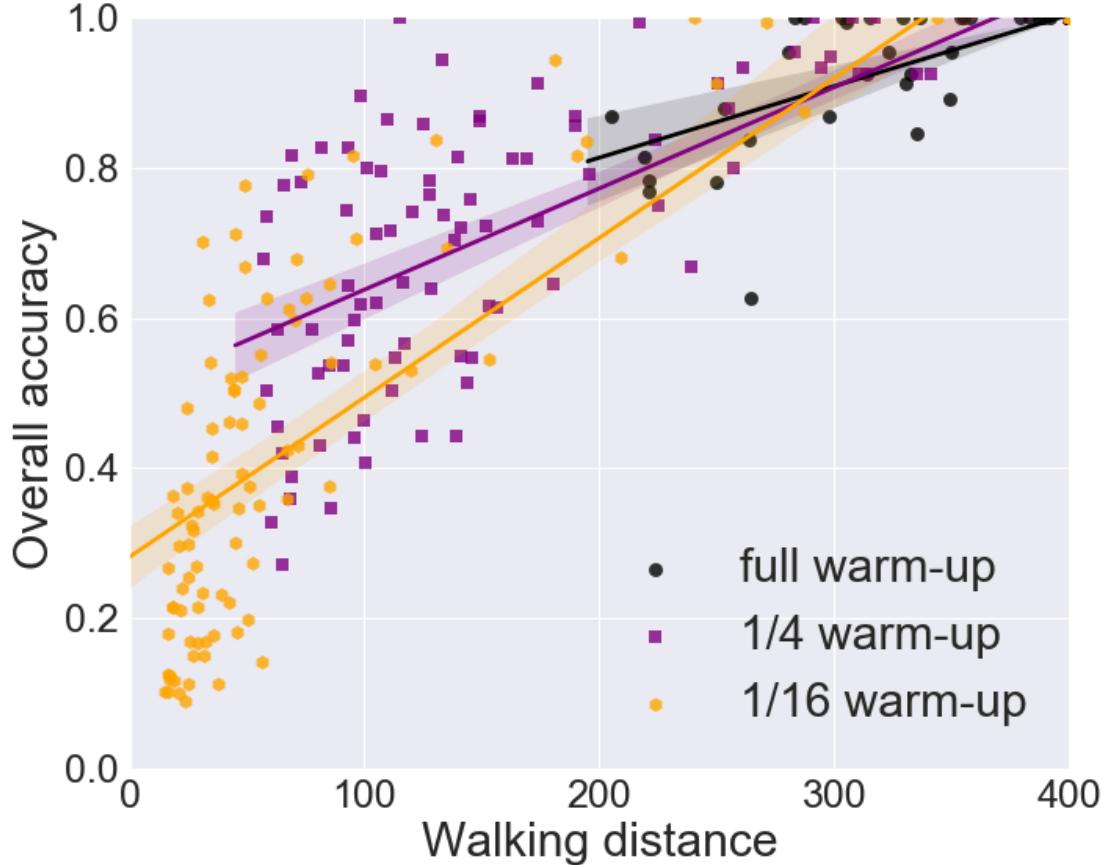


Forbes US Richest: N=400, K=4 (27 female, 373 male)

diversity on gender: select 2 per gender

Warm-up can be shorter

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

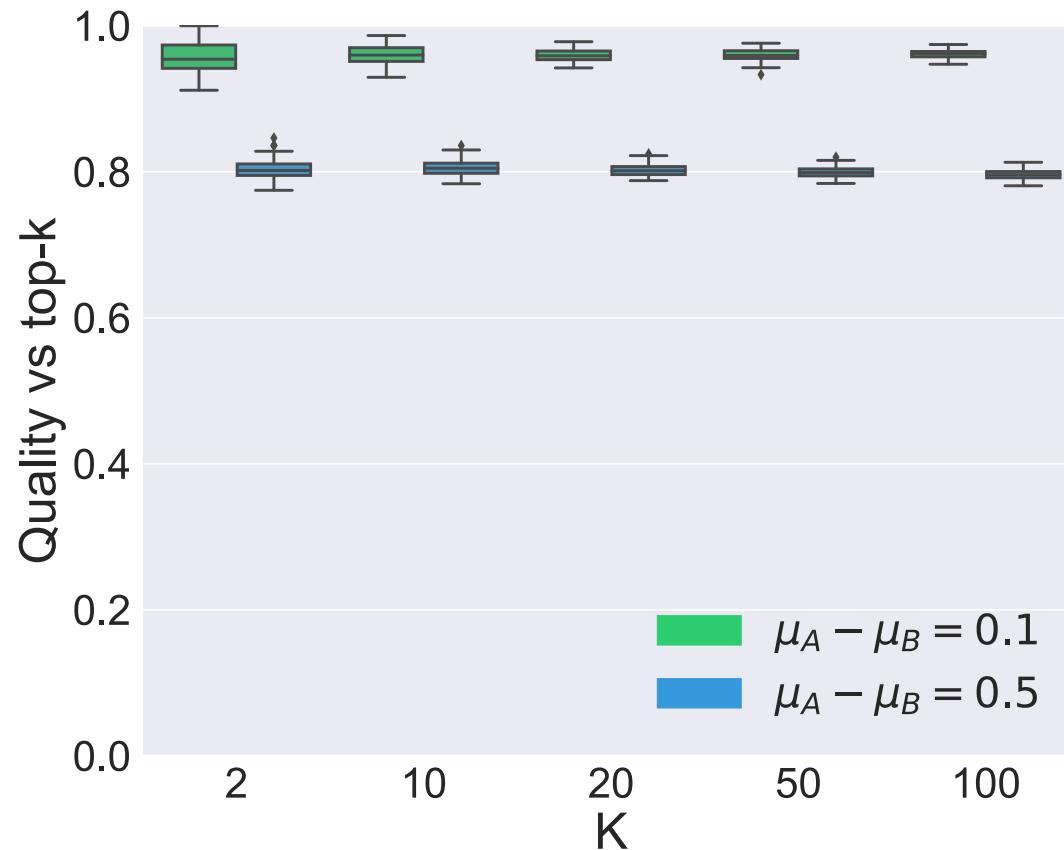


Forbes US Richest: N=400, K=4 (27 female, 373 male)

deferred list variant, diversity on gender: select 2 per gender

The cost of diversity

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

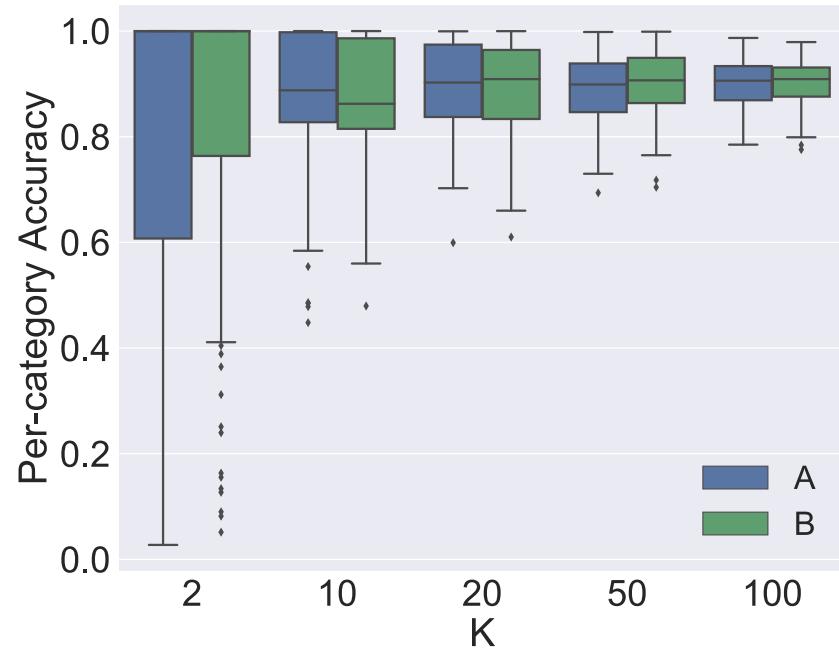


static variant (see paper), synthetic data in categories A and B, score lower for B

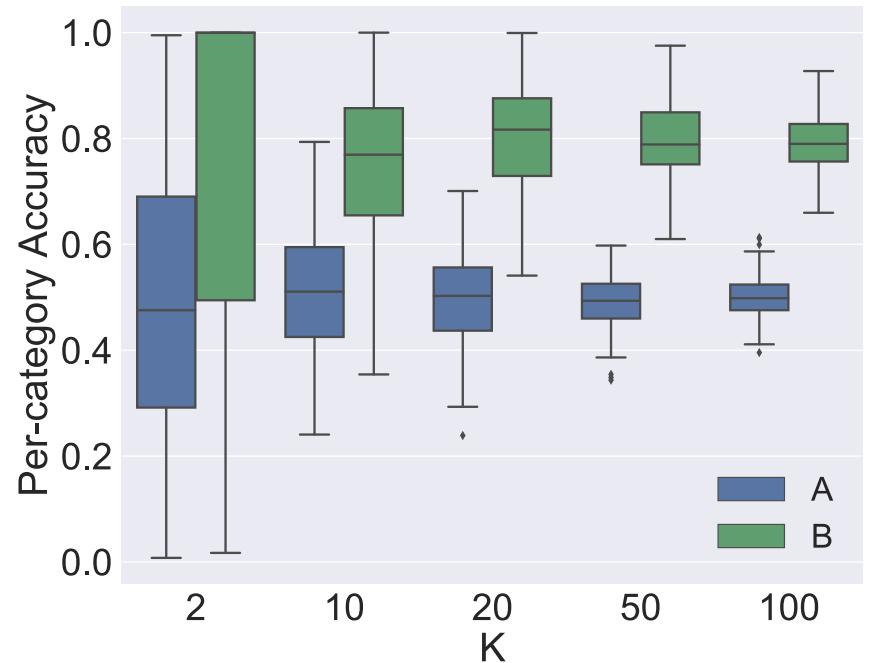
Per-category warm-up is crucial

[J. Stoyanovich, K. Yang, HV Jagadish, EDBT (2018)]

Per-category warm-up period



Common warm-up period



synthetic data with categories A and B, score depends on category, lower for A

diversity by design

AI's White Guy Problem

The New York Times

Artificial Intelligence's White Guy Problem

By KATE CRAWFORD JUNE 25, 2016



Like all technologies before it, artificial intelligence will reflect the values of its creators. So **inclusivity matters** — from who designs it to who sits on the company boards and which ethical perspectives are included.

Otherwise, **we risk constructing machine intelligence that mirrors a narrow and privileged vision of society**, with its old, familiar biases and stereotypes.

A technical review paper

REVIEW

Diversity in Big Data: A Review

Marina Drosou¹, H.V. Jagadish², Evangelia Pitoura¹, and Julia Stoyanovich^{3,*}

Big Data

Volume 5 Number 2, 2017

© Mary Ann Liebert, Inc.

DOI: 10.1089/big.2016.0054

Abstract

Big data technology offers unprecedented opportunities to society as a whole and also to its individual members. At the same time, this technology poses significant risks to those it overlooks. In this article, we give an overview of recent technical work on diversity, particularly in selection tasks, discuss connections between diversity and fairness, and identify promising directions for future work that will position diversity as an important component of a data-responsible society. We argue that diversity should come to the forefront of our discourse, for reasons that are both ethical—to mitigate the risks of exclusion—and utilitarian, to enable more powerful, accurate, and engaging data analysis and use.

Keywords: data; diversity; empirical studies; models and algorithms; responsibly

Diversity measures

[M. Drosou, HV Jagadish, E. Pitoura, J. Stoyanovich; BigData 2017]

- **Distance-based**: the most common
 - **MaxSum**: maximize total (or average) pair-wise distance in S
 - **MaxMin**: maximize lowest pair-wise distance in S , a variant of the p-dispersion problem
- **Coverage-based**: “Noah’s Arc”, based on a set of pre-defined discrete categories (topics, demographics...)
- **Novelty-based**: relative to elements seen in the past (e.g., Maximal Marginal Relevance - MMR)

Diversity models

[M. Drosou, HV Jagadish, E. Pitoura, J. Stoyanovich; BigData 2017]

Diversity is an aspect of **quality of a collection** of items S . It is often traded off with **per-item quality** (utility).

$$S = \operatorname{argmax}_{S' \subseteq I, |S'|=k} \operatorname{div}(S')$$

Given a set of items I , select a diverse set of items S of size k , as quantified by diversity measure div .

Other variants: **aggregate diversity** (e.g., diversify recommendations to each user and across users) and **bundle diversity** (e.g., dinner and a movie).

Diversity & friends

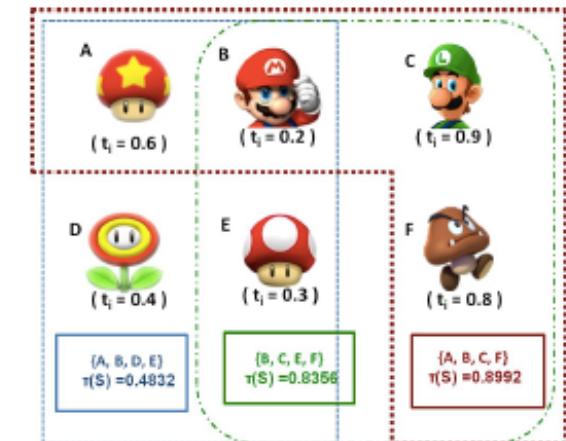
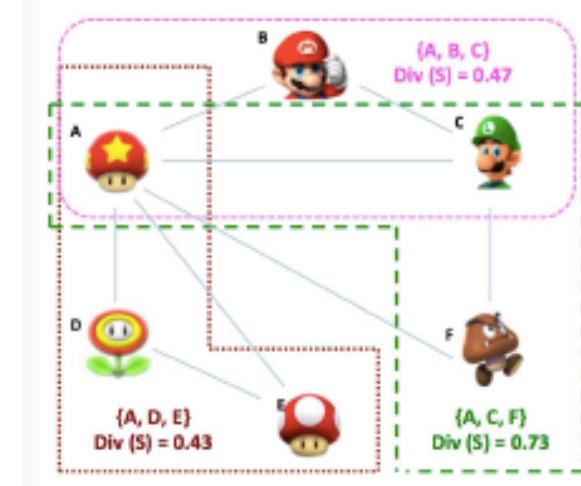
- For a given user consuming information in search and recommendation, relevance is important, but so are:
 - **diversity** - avoid returning similar items
 - **novelty** - avoid returning known items
 - **serendipity** - surprise the user with unexpected items
- For a set of users
 - uncommon information needs must be met: **less popular** “in the tail” queries constitute the overwhelming majority
 - lack of diversity can lead to **exclusion**

Jonas Lerman: “... the nonrandom, systematic omission of people who live on big data’s margins, whether due to poverty, geography, or lifestyle...”

Diversity can improve accuracy

[T. Wu, L. Chen, P. Hui, C.J. Zhang, W. Li; PVLDB 2015]

- Importance of diversity of opinion for **accuracy** is well-understood in the social sciences
 - Diversity is crucial in crowdsourcing, see Surowiecki “*The Wisdom of the Crowds*” 2005
 - The “Diversity trumps ability theorem”
- Crowd diversity: an aggregate of pair-wise diversity
- S-Model:** similarity-driven / task-independent
- T-Model:** task-driven, opinions are probabilistic



Online dating

[J. Stoyanovich, S. Amer-Yahia, T. Milo; EDBT 2011]

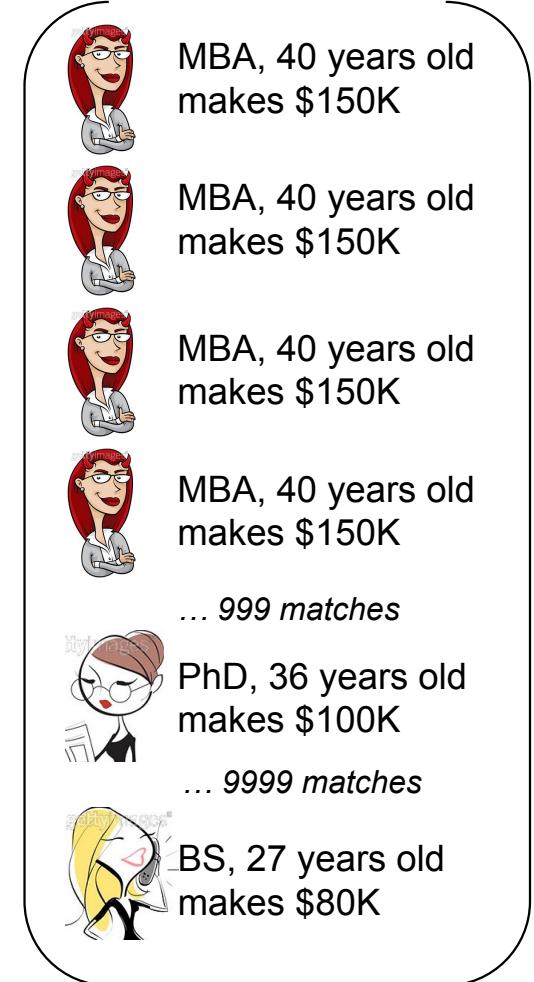
Dating query: female, 40 or younger, at least some college, in order of decreasing income

Results are homogeneous at top ranks

Both the seeker (asking the query) and the matches (results) are dissatisfied

Crowdsourcing, crowdfunding, ranking of Web search results, ... - all subject to this problem

the rich get richer, the poor get poorer

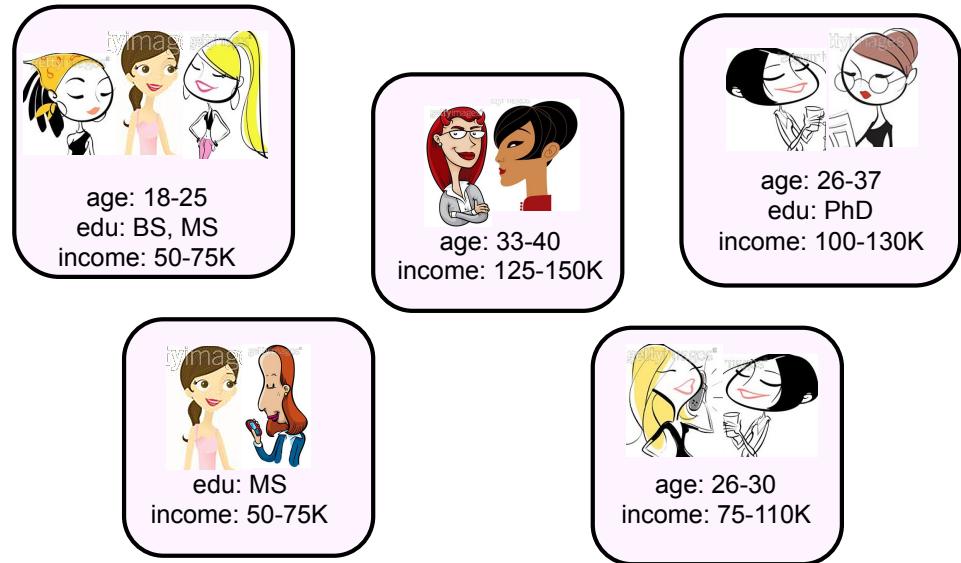


Diversity can improve user engagement

[J. Stoyanovich, S. Amer-Yahia, T. Milo; EDBT 2011]



Return clusters that expose **best from among comparable** items (profiles) w.r.t. user preferences



More diverse items seen, and liked, by users

Users are more engaged with the system

Why is diversity important?

- Unlike fairness, there is **no legal reason** to enforce diversity
- However, there are strong **utilitarian reasons**: diversity leads to better user satisfaction (IR, recommendation), higher quality of results (crowdsourcing, team formation), more efficient resource allocation (matchmaking)
- Further, diversity levels the playing field and **improves fairness in the long run**