



Evaluating Group Fairness Metrics for Rankings

Marlene Lutz, Tobias Schumacher, Sandipan Sikdar and Markus Strohmaier

8th International Conference on Computational Social Science

Ranking applications



Hiring



Admission



Lending



E-commerce

...

Algorithmic Ranking

- Ranking algorithms assist in decisions that impact peoples wellbeing and success
- Rankings should not only be accurate but also **fair**
- Growing body of research for designing and deploying **fairness metrics** for rankings

■ **How can fairness metrics for rankings be compared and evaluated?**

Contribution

1

Proposal of 13 properties

We propose properties that are informative for the construction and evaluation of group fairness metrics for rankings.

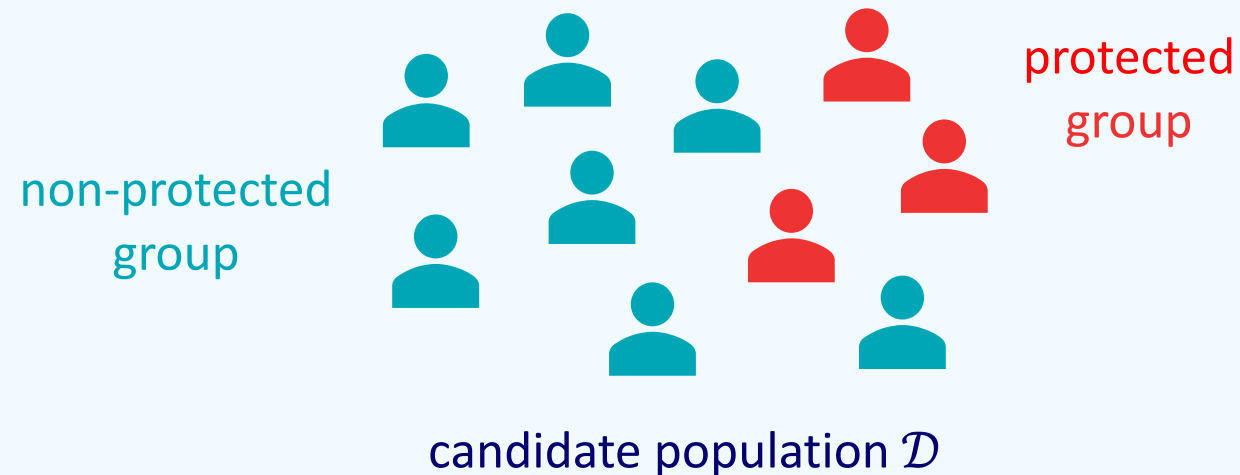
2

Application to fairness metrics

We deploy the properties to 10 existing group fairness metrics for rankings and study the extent to which they satisfy them.

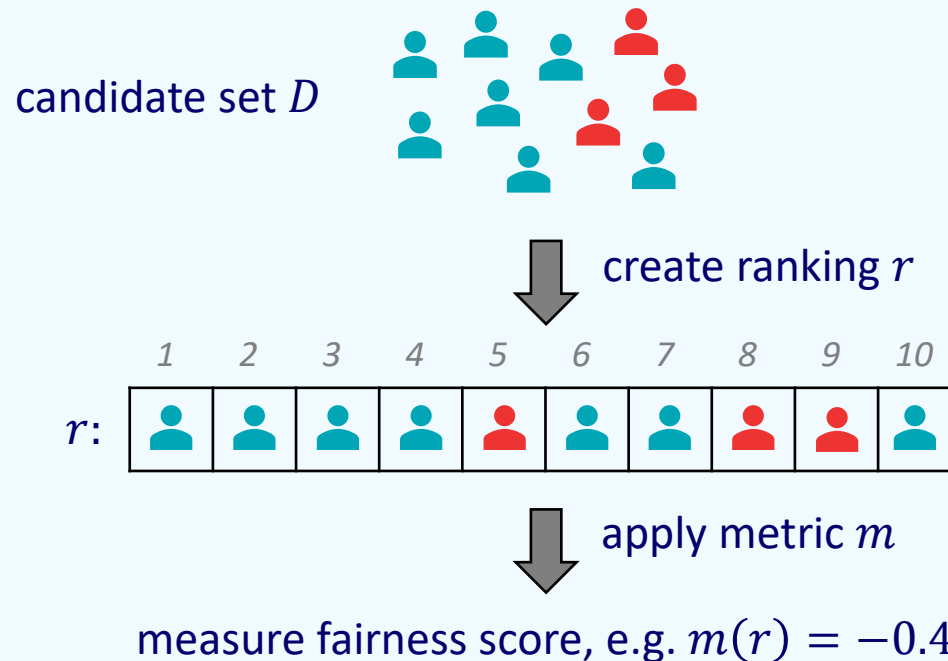
Fair Ranking Setup

Goal: Ranking a set of candidates $D \subseteq \mathcal{D}$ s.t. the ranking r is fair with respect to a protected group (*group fairness*).



Fair Ranking Setup

Goal: Ranking a set of candidates $D \subseteq \mathcal{D}$ s.t. the ranking r is fair with respect to a protected group (*group fairness*).



Higher fairness score is better

Properties for Fair Ranking Metrics

P1: Group Distinctiveness

P2: Boundedness

P3: Monotonicity

P4: Deepness

P5: Intra-group Fairness

P6: Invariance to Linear Transformation of Relevance Scores

Universal Properties
(both ranking settings)

P7: Optimality of Random Rankings

P8: Invariance to Ranking Length

P9: Invariance to Group Proportions

P10: Symmetric Penalties for all Groups

Ranking the full popul.

P11: Deepness Threshold

P12: Closeness Threshold

P13: Confidence

Ranking a subset

Properties for Fair Ranking Metrics

P1: Group Distinctiveness

P2: Boundedness

P3: Monotonicity

P4: Deepness

P5: Intra-group Fairness

P6: Invariance to Linear Transformation of Relevance Scores

Universal Properties
(both ranking settings)

P7: Optimality of Random Rankings

P8: Invariance to Ranking Length

P9: Invariance to Group Proportions

P10: Symmetric Penalties for all Groups

Ranking the full popul.

P11: Deepness Threshold

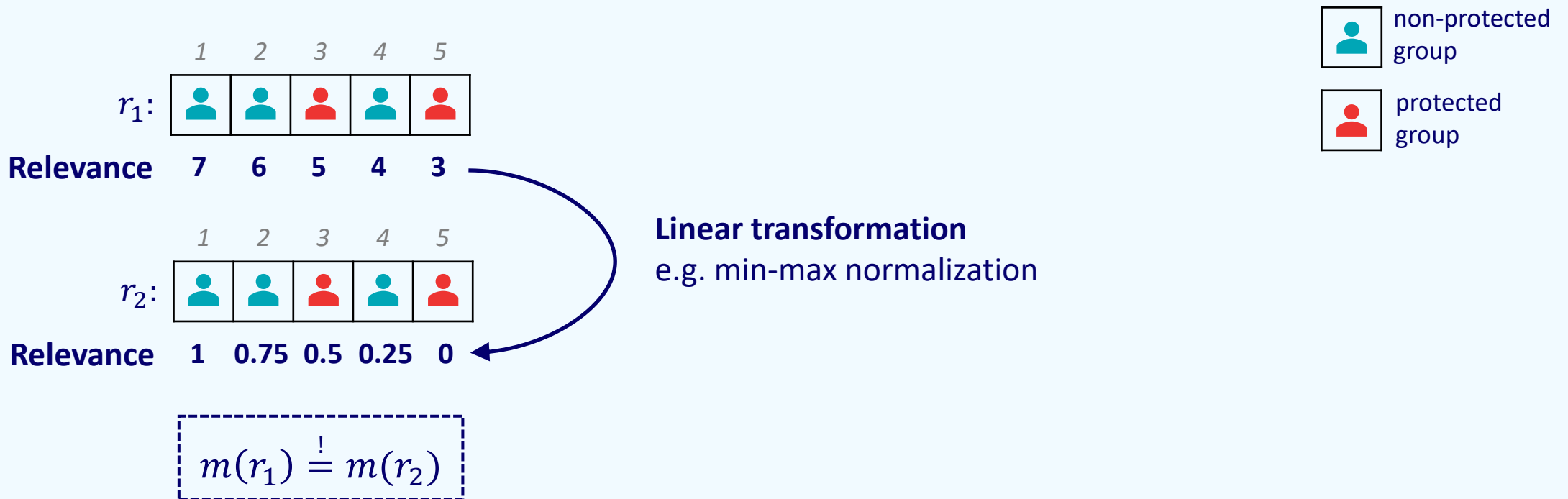
P12: Closeness Threshold

P13: Confidence

Ranking a subset








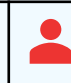


Invariance to Linear Transformation of Relevance Scores

A metric m is **invariant to linear transformation of relevance scores** if its values do not change after transforming the rel. scores of a candidate set.



Invariance to Linear Transformation of Relevance Scores

Example: Only transformed relevance scores are accessible (e.g. for privacy reasons).

	1	2	3	4	5
r_1 :					
Relevance (true)	7	6	5	4	3
	1	2	3	4	5
r_2 :					
Relevance (transformed)	1	0.75	0.5	0.25	0

Fairness:

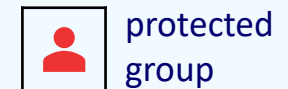
$$DTD(r_1) = -0.01$$

in favour of non-protected group

Fairness:

$$DTD(r_2) = 0.74$$

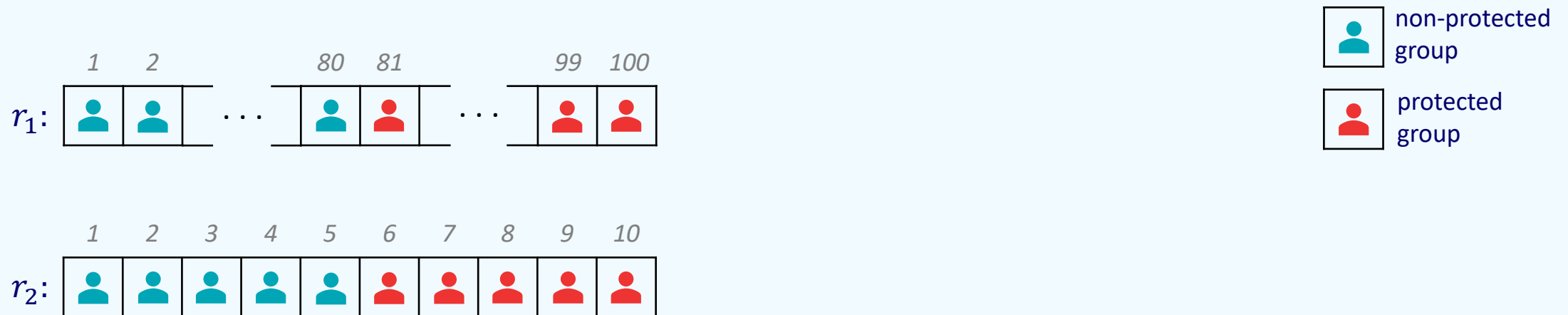
in favour of protected group



DTD Disparate treatment difference
0 - perfect fairness

Invariance to Ranking Length

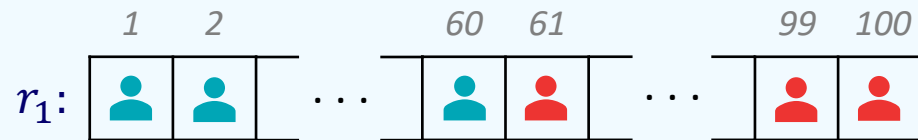
A metric m is **invariant to ranking length** if its *worst-case values* (total disadvantage of one group) do not change for different ranking lengths.



$$m(r_1) \stackrel{!}{=} m(r_2)$$

Invariance to Ranking Length


Example: Comparing the fairness of hiring processes (= rankings of job applicants) at different companies.




$ED(r_1) = -0.08$
e.g. a job at PopularCompany



$ED(r_2) = -0.24$
e.g. a job at UnpopularCompany

 non-protected group (60%)

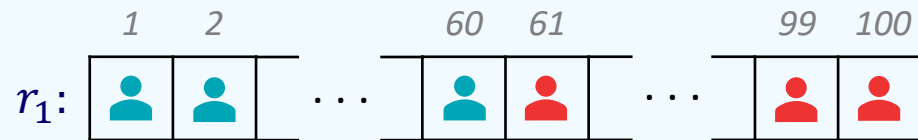
 protected group (40%)

ED Exposure difference
0 - perfect fairness


Higher fairness score is better


Invariance to Ranking Length

Example: Comparing the fairness of hiring processes (= rankings of job applicants) at different companies.

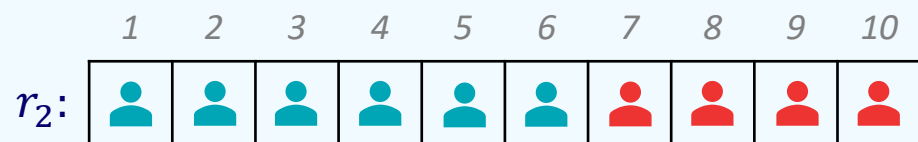


$ED(r_1) = -0.08$
e.g. a job at PopularCompany

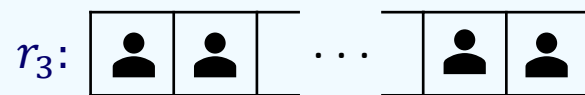
 non-protected group (60%)

 protected group (40%)

ED Exposure difference, 0 - perfect fairness



$ED(r_2) = -0.24$
e.g. a job at UnpopularCompany



$ED(r_3) = -0.07$
How can this value be interpreted?

Higher fairness score is better

Application to Fair Ranking Metrics

P6: Invariance to Linear Transformation of Relevance Scores

P8: Invariance to Ranking Length

Metric	P 1	P 2	P 3	P 4	P 5	P 6	P 7	P 8	P 9	P 10	P 11	P 12	P 13
Normalized discounted difference (<i>rND</i>)	×	✓	×	×	N/A	N/A	×	×	×	×	×	×	×
Normalized discounted KL-divergence (<i>rKL</i>)	×	✓	×	×	N/A	N/A	×	×	×	×	×	×	×
Normalized discounted ratio (<i>rDR</i>)	×	✓	×	×	N/A	N/A	×	×	×	×	×	×	×
Exposure difference (<i>ED</i>)	✓	✓	✓	✓	N/A	N/A	✓	×	×	×	✓	✓	✓
Exposure ratio (<i>ER</i>)	✓	×	✓	✓	N/A	N/A	×	×	×	×	✓	✓	✓
Disparate treatment difference (<i>DTD</i>)	✓	×	✓	✓	×	×	✓	×	×	×	✓	✓	✓
Disparate treatment ratio (<i>DTR</i>)	✓	×	✓	✓	×	×	×	×	×	×	✓	✓	✓
Disparate impact difference (<i>DID</i>)	✓	×	✓	✓	×	×	✓	×	×	×	✓	✓	✓
Disparate impact ratio (<i>DIR</i>)	✓	×	✓	✓	×	×	×	×	×	×	✓	✓	✓
Pairwise statistical parity (<i>PSP</i>)	✓	✓	✓	×	N/A	N/A	✓	✓	✓	✓	N/A	N/A	N/A

Universal Properties

Ranking full popul.

Ranking a subset

- ✓ property satisfied
- ✗ property not satisfied
- N/A property not applicable

Conclusion



Not every application requires satisfaction of all properties!

- Highlight limitations of existing metrics
 - Lack of interpretability and comparability
 - Unexpected side effects
- Support informed evaluation and design of group fairness metrics for rankings
- Guide practitioners in choosing appropriate metrics

Thank You!

Contact:

Marlene Lutz – marlene.lutz@uni-mannheim.de

Tobias Schumacher – tobias.schumacher@uni-mannheim.de