

Algorithm Fairness in Online Information Systems

Ansgar Koene¹, Sofia Ceppi², Helena Webb³, Menisha Patel³, Elvira Perez¹, Giles Lane⁴, Marina Jirotko³, Michael Rovatsos²

(1) University of Nottingham, (2) University of Edinburgh, (3) University of Oxford, (4) Proboscis

Introduction to UnBias

UnBias investigates the user experience of algorithm driven services and the processes of algorithm design to provide policy recommendations, ethical guidelines and a stakeholder co-produced 'fairness toolkit'. We focus on a wide range of stakeholders and carry out activities that

- 1) support user understanding of algorithm mediated services
- 2) raise awareness among 'smart' systems providers about the concerns and rights of users
- 3) generate debate about 'fair' operation of algorithms in modern life.

Multi-Stakeholder Engagement

Engaging with academics, civil-society, industry, regulators:

- What constitutes a fair algorithm?
- What kinds of (legal and ethical) responsibilities do internet companies have to ensure fair and unbiased algorithmic decisions?
- What might enhance user awareness/trust in online algorithms?
- How might concepts of fairness be built into algorithmic design?

The discussion was structured around four case studies:

- 1) gaming the system - anti-Semitic autocomplete and search results;
- 2) news recommendation and fake news;
- 3) personalisation algorithms;
- 4) algorithmic transparency.

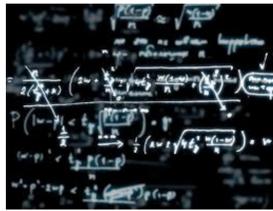


How to bump Holocaust deniers off Google's top spot? Pay Google

Facebook fires trending team, and algorithm without humans goes crazy

Google is unhappy with Holocaust denial beating the truth in its search results - but it probably makes more money that way

Module pushes out false story about Fox's Megyn Kelly, offensive Ann Coulter headline and a story link about a man masturbating with a McDonald's sandwich



Defining Algorithm Fairness

Candidate definition: "A context-dependent evaluation of the algorithm and/or outcomes against socio-cultural values".

E.g.: equal treatment without prejudice due to task-irrelevant factors.

Criteria relating to social norms and values:

- Sometimes disparate outcome are acceptable if based on individual lifestyle choices over which people have control.
- Ethical precautions are more important than higher accuracy.
- There needs to be a balancing of individual values and socio-cultural values. Problem: How to weigh relevant social-cultural value?

Criteria relating to system reliability:

- Results must be balanced with due regard for trustworthiness.
- Needs for independent system evaluation and monitoring over time.

Criteria relating to (non-)interference with user control:

- Subjective fairness experience depends on user objectives at time of use, therefore requires an ability to tune the data and algorithm.
- Users should be able to limit data collection about them and its use. Inferred personal data is still personal data. Meaning assigned to the data must be justified towards the user.
- Reasoning/behaviour of algorithm demonstrated/explained in a way that can be understood by the data subject.
- If not vital to the task, there should be opt-out of the algorithm
- Users must have freedom to explore algorithm effects, even if this would increase the ability to "game the system"
- Need for clear means of appeal/redress for impact of the algorithmic system that the user cannot control.

Dissenting opinion

"I struggle with the underlying anthropomorphisation of algorithms when one speaks of "algorithm fairness". In my view, the concept of fairness is as you rightly noted deeply enshrined in the specific socio-cultural codes of the respective group of actors. Since algorithms as such constitute only the tools that actors in the social context use to achieve (some of) their objectives, one should also judge fairness probably more along the behaviour of these actors, their objectives, and methods. This implies that both the process of the algorithm and its outcomes need to be taken into account."

Collective Fairness

Many multi-user/limited 'resource' scenarios are characterised by a combinatorial nature, i.e., decisions applied to one user also affect other users such that decision optimization must consider the requirements and preferences of all users at the same time. E.g. problem of allocating course-work topics to students when each student must have exactly one coursework topic and each topic can be assigned to at most one student.

- Each student expressed preferences over each topic by assigning them a score from 1 to 7 representing how desirable that topic is.
- 5 Algorithms used to determine 'optimal' course-work assignment using different fairness criteria:
 - ❖ A1: minimise total distance while guaranteeing at least 70% of maximum possible utility
 - ❖ A2: maximise the minimum individual student utility while guaranteeing at least 70% of maximum possible total utility
 - ❖ A3: maximise total utility
 - ❖ A4: maximise the minimum individual student utility
 - ❖ A5: minimise total distance

Where: utility = score the student gave to the assigned topic
distance = sum of absolute difference between student's utility and utility for each other student

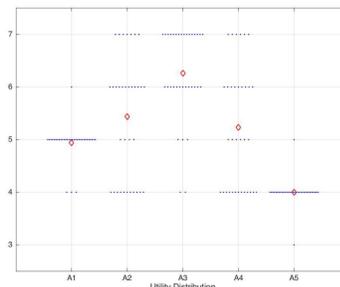


Figure 1: Student's utility distribution (blue dots) and mean (red diamond)

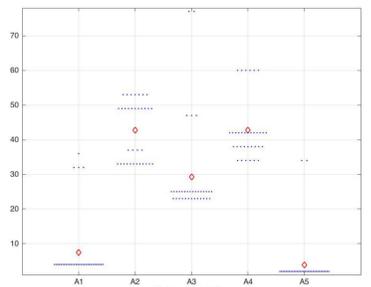


Figure 2: Student's distance distribution (blue dots) and mean (red diamond)

We asked the student to rank their preference for each algorithm:

1. on the basis of knowing only the topic assignment outcomes
2. knowing the fairness criteria used by each algorithm

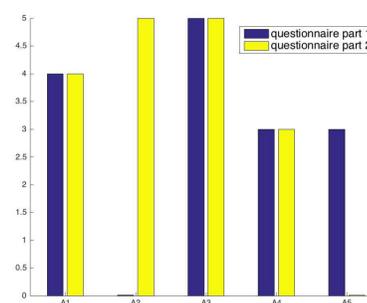


Figure 3: Algorithms most preferred by the students when they are not aware of the values embedded in the algorithm (blue) and when a description of the algorithms is provided (yellow) to them. Note that students were allowed to select more than one algorithm.

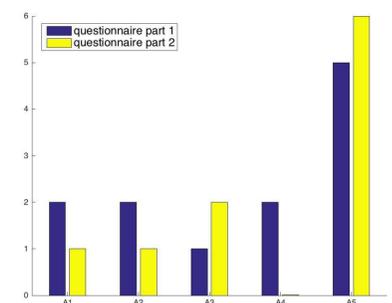


Figure 4: Algorithms least preferred by the students when they are not aware of the values embedded in the algorithm (blue) and when a description of the algorithms is provided (yellow) to them. Note that students were allowed to select more than one algorithm.

Individual students preferred different algorithms. Preferences differed with/without knowledge of optimisation criteria.