

KEDL2019

Algorithmic Bias: What is it, Why Does it Matter and What's Being Done About it?

Paul Clough, University of Sheffield and Peak Indicators



Talk overview

- Introduction to algorithmic bias
- Study of gender biases in search engines
- What can be done?
- Summary

PEAK indicators



Introduction

Roducion

Data-driven decision making

- Algorithms can support decision-making by
 - Prioritising and ranking
 - Making predictions (regression and classification)
 - Finding patterns and associations
 - Filtering

- Predictive models used in
 - Personalised pricing and recommendations
 - Credit scoring
 - Automated CV screening of job applicants
 - Profiling of potential suspects by the police



S. Hajian, F. Bonchi and C. Castillo (2016). <u>Algorithmic Bias: From Discrimination Discovery to</u> <u>Fairness-aware Data Mining</u>. In KDD, pp. 2125-2126.

"Bias: inclination or prejudice for or against one person or group, especially in a way considered to be unfair"

"Predictive models can discriminate people, even if the computing process is fair and well-intentioned"

THE DAILY NEWSLETTER Sign up to our daily email newsletter



Gender-recognition AI tools correctly identify white men more accurately than BAME women.

Google's online advertising system showed high-income jobs to men much more often than to women.

Facebook's automatic translation software chose the wrong translation for Hebrew "good morning" vs. "attack them".

Discriminating algorithms: 5 times Al showed prejudice

Artificial intelligence is supposed to make life easier for us all – but it is also prone to amplify sexist and racist biases from the real world

TECHNOLOGY 12 April 2018, updated 27 April 2018



John Lamb/Getty

By Daniel Cossins

Modern life runs on intelligent algorithms. The data-devouring, self-improving computer programmes that underlie the artificial intelligence revolution already determine Google search results, Facebook news feeds and online shopping recommendations. Increasingly, they also decide how easily we get a mortgage or a job interview, the chances we will get stopped and searched by the police on our way home, and what penalties we face if we commit a crime, too. Researchers found that COMPAS predicts that black defendants pose a higher risk of recidivism than they do, and the reverse for white defendants.

> In 2016, the Human Rights Data Analysis Group found that PredPol could lead police to unfairly target certain neighbourhoods.



https://www.oxfordinsights.com/racial-bias-in-natural-language-processing

MATTHEW REIDSMA

ARTICLES TALKS WORK NOTES

ALGORITHMIC BIAS IN LIBRARY DISCOVERY SYSTEMS

March 11, 2016

<u>« Prev</u>

More and more academic libraries have invested in discovery layers, the centralized "Google-like" search tool that returns results from different services and providers by searching a centralized index. The move to discovery has been driven by the ascendence of Google as well as libraries' increasing focus on user experience. Unlike the vendor-specific search tools or federated searches of the previous decade, discovery presents a simplified picture of the library research process. It has the familiar single search box, and the results are not broken out by provider or format but are all shown together in a list, aping the Google model for search results.

Discovery's promise of a simple search experience works for users, more often than not. But discovery's external simplicity hides a complex system running in the background, making decisions for our users. And it is the rare user that questions these decisions. As Sherry Turkle (1997) observed, users approach complex systems

https://matthew.reidsrow.com/articles/173

It's not all bad news

"It is a myth to think that algorithms are objective, but also a myth to think that human processes are not subject to biases on par with algorithms."





A study of gender bias in search engines

Popula

Competent Men and Warm Women: Gender Stereotypes and Backlash in Image Search Results

Jo Bates

Information School

University of Sheffield, UK

jo.bates@sheffield.ac.uk

Jahna Otterbacher Social Information Systems Open University of Cyprus jahna.otterbacher@ouc.ac.cy

ABSTRACT

PEAK

There is much concern about algorithms that underlie information services and the view of the world they present. We develop a novel method for examining the content and strength of gender stereotypes in image search, inspired by the trait adjective checklist method. We compare the gender distribution in photos retrieved by Bing for the query "person" and for queries based on 68 character traits (e.g., "intelligent person") in four regional markets. Photos of men are more often retrieved for "person," as compared to women. As predicted, photos of women are more often retrieved for warm traits (e.g., "emotional") whereas agentic traits (e.g., "rational") are represented by photos of men. A backlash effect, where stereotype-incongruent individuals are penalized, is observed. However, backlash is more prevalent for "competent women" than "warm men." Results underline the need to understand how and why biases enter search algorithms and at which stages of the engineering process.

Author Keywords

Algorithmic bias; "Big Two" dimensions of social perception; gender stereotypes; image search.

Otterbacher, J., Bates, J., and Clough P. (2017), Competent Men and Warm Women: Gender Stereotypes and Backlash in Image Search Results, In Proceedings of CHI'2017, pp. 6620-6631. participation in public life [20]. Even when users are intimately familiar with a system, they are often unaware that algorithms filter their access to information [14] and users hold beliefs about algorithms, which, true or not, influence how they use systems [39].

Paul Clough

Information School

University of Sheffield, UK

p.d.clough@sheffield.ac.uk

Machines running algorithmic processes have become the new gatekeepers, largely determining what and whom we see, and do not see [8]. Given the power that algorithms exert, researchers must scrutinize these processes, their potential biases and social impact; in other words, we must work toward "algorithmic accountability" [11] and "algorithmic transparency" [10]





Saarah	
Search	

11

Did you know that if you do an image search for 'person' the results contain twice as many men as women? RECENT POSTS

In this episode we investigate why this is the case by looking at bias and how it has permeated our online world, and in particular, search engines. Search engines now have

EPISODE 1: SEARCH ENGINES AND BIAS

MACHINE MINDS

http://www.machinemindspodcast.com/



http://www.slate.com/articles/technology/future_ tense/2015/12/why_google_search_results_favor_ democrats.html



"And there's the illusion of neutrality. About two-thirds of Americans who use search engines believe they are completely unbiased, according to a 2012 Pew study. The study showed that "73 percent of search engine users say that most or all the information they find as they use search engines is accurate and trustworthy." Search engines are more trusted than the news media itself."

Who is a nurse?



Matthew Kay, Cynthia Matuszek, and Sean A. Munson. 2015. Unequal Representation and Gender Stereotypes in Image Search Results for Occupations. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems* (CHI '15). ACM, New York, NY, USA, 3819-3828.



Who is a nurse?



Male nurse

Ь	Ma	le Nurse			م										Sign in	R 100	r =
	We	b Images	• Videos	Maps New	s										SafeSearch:	Moderate -	Filter \bigtriangledown
Hot Male Nu	irse	Male Nurse with Patient	Male Nurse Cartoon	Male Nurse Clip Art	Cute Male Nurse	Male Nurse Stereotypes	Funny Male Nurse	Male Nurse Jokes	Male Male Teacher Libraria	Male Nurse Meme	Black Male Nurse	Male and Female Nur	Male Nurses se Working	Male Nurse	Male Nurse Costume	Male Doctor P	Male atlent
								1			5			5			
a la		k	and a second		4					Received and a second				1			
	1								MALE					i al			Bac.
				2												-	
f								3			2		fr "	26			5



Intelligent person



Stereotypes beyond occupation – personality traits Who does Bing say represents a 'person'?

Shy person



Shy person



Gender distribution in images of top-ranked 50 images

Women/girls:		25
(50%)		
Men/boys:		5
(10%)		
Mixed gender:	0	
Unknown/none:	20 (40%)	

Can we automatically identify gender distribution in results?

Stereotypes: "Big Two" of person perception

- Personality traits captured by the 'Big five'
- Our perceptions of others are based on two dimensions [Fiske et al., 2002]
 - 1) <u>agency (or competence)</u>: whether or not we perceive someone as being capable of achieving his/her goals
 - 2) <u>warmth (or communality)</u>: whether or not we think someone has pro-social intentions or is a threat to us
- Stereotypes are captured by combinations of the two dimensions [Cuddy et al., 2008]
 - <u>Women</u>: [low agency, high warmth]
 - <u>Men</u>: [high agency, low warmth]

Trait adjective checklist method

- How do we measure content and strength of given social stereotype?
 - Trait adjective checklist method
- Used in the Princeton Trilogy studies of ethnic and racial stereotypes [Katz & Braly, 1933]
- Participants describe target social groups using list of trait adjectives
- 68 traits developed in cross-lingual study across five countries [Abele et al., 2008]

able	egoistic	persistent					
active	emotional	polite					
affectionate	energetic	rational			son	0	
altruistic	expressive	reliable			501	~	
ambitious	fair	reserved		Web Images •	Videos Maps News		
assertive	friendly	self-confident		How Consistieus Consistieus	Consistions	Conselectious	
boastful	gullible	self-critical		Are You Personality Type	Clip Art Conscientious	Student Jobs for Personality	
capable	harmonious	self-reliant		Conscientiousness		WOLL ACHIEVER	
caring	hardhearted	self-sacrificing		High Low Low Few Goals Many Goals Many Goals			
chaotic	helpful	sensitive		Set-disciplined Cavely Cavely Responsible Disorganized	The	CONSCIENTIOUSNESS	
communicative	honest	shy				SELF-DISCIPLINED CAREFUL DISCHARGE	
competent	independent	sociable					
competitive	industrious	striving	(
conceited	insecure	strong-minded			Search mai	rkets	
conscientious	intelligent	supportive				1	
considerate	lazy	sympathetic					
consistent	loyal	tolerant			US-EN	4	
creative	moral	trustworthy			IN-EN		
decisive	obstinate	understanding			ZA-EN	l	
detached	open	vigorous					
determined	open-minded	vulnerable					
dogmatic	outgoing	warm					
dominant	perfectionistic						

Research Questions

- RQ1: Baseline Representation bias
 - In a search for "person" which genders are depicted?
- RQ2: Stereotype content and strength
 - Which character traits are most often associated with which genders?
 - Are these associations consistent across bing search markets? (UK, US, IN, ZA)
- RQ3: Backlash effects
 - How are stereotype-incongruent individuals depicted?

shy person		٦			
Web Images -	Videos Maps News	12			
Shy Person Shy Person Clip Art Cartoon	Shy Person Shy Person In Class Drawing	How Many People Are Shy	Another Word for Shy Person	Quiet Person	Nic Pers
9	l'm Shy!		has		L
2 PP	5- F	19		14-3	
WOMAN/GIRL	WOMAN/GIRL	WOMA	N/GIRL	MAN/B	OY
I'm actually a really shy person.		PA		10	
Like when you first meet me, our conversation is going to be awkward no matter what because I resulden't have any sites what to take about. It's also worse	2 😂 🗾			-	R
WOMAN/GIR		WOMAN/GIRL	WON	AN/GIRL	2
	Shy People Problems	" FRIENELY AM WORE SHY TH	anan Shy	people	
1 1	When people ask why you' so quiet.	AND FAMILIES AND FAMILIES AND FAMILIES	but	ce everyth they do no	ning ot
WOMAN/GIRL	NONE	NONE	get	NONE	

Challenges in automating the process – how hard is recognising gender in images for people?

Pilot study on Crowdflower

- 1,000 "person" images from UK market
- 3 annotators per image
- Is the image:

1) a photograph, 2) a sketch/illustration, 3) some other type?

• Does the image depict:

only women/girls, 2) only men/boys, 3) mixed gender group,
 gender ambiguous person(s), 5) no person(s)?



Classifying image type

	# Images	Inter-judge agreement	
Photos	576	0.97	
Sketches	346	0.96	
Other	22	0.74	
No longer accessible	56	1.00	

High

degree of

agreement



Classifying gender

	Women /girls	Men/boys	Mixed gender	Unknown	No persons	Inter-judge agreement
Photos	0.27	0.55	0.10	0.07	0.01	0.94
Sketches	0.08	0.28	0.05	0.55	0.04	0.91

Automating gender recognition

- Clarify API
 - General image recognition tool
 - Coverage: 95% of images from bing
 - Provides 20 textual concept tags
- Linguistic Inquiry and Wordcount (LIWC) [Pennebaker et al., 2015]
 - Female references: mom, girl
 - Male references: dad, boy





Performance on gender classification

	Ν	Precision	Recall	F ₁
Recognizing photographs	473	0.91	0.75	0.822
Women/girls	130	0.89	0.60	0.717
Men/boys	282	0.95	0.67	0.786
Other	61	0.68	0.82	0.743

RQ1: who represents a "person"?





Consistent gendering of traits across regions

Men/boys:

ambitious, boastful, competent, conceited, conscientious, consistent, decisive, determined, gullible, independent, industrious, intelligent, lazy, persistent, rational, self-critical, vigorous

Women/girls:

detached, emotional, expressive, fair, insecure, open-minded, outgoing, perfectionistic, self-confident, sensitive, shy, warm

Gender-neutral:

able, active, affectionate, caring, communicative, competitive, friendly, helpful, self-sacrificing, sociable, supportive, understanding, vulnerable

- The bing algorithm is not itself gender-biased
- However, bing image results do perpetuate gendered perceptions of personhood



http://www.websci16.org/sites/websci16/files/keynotes/keynote_baeza-yates.pdf

What can be done?

Roducion

Engineering for equity during all phases of ML design



Technical solutions

- Tools to identify data and algorithmic bias
- Tools to reduce discrimination
- Explainable AI
- Tools for algorithmic auditing



NEWS

Google Cloud AI/ML customers

A survey on measuring indirect discrimination in machine learning

INDRĖ ŽLIOBAITĖ, Aalto University and Helsinki Institute for Information Technology HIIT

Nowadays, many decisions are made using predictive models built on historical data. Predictive models may systematically discriminate groups of people even if the computing process is fair and well-intentioned. Discrimination-aware data mining studies how to make predictive models free from discrimination, when historical data, on which they are built, may be biased, incomplete, or even contain past discriminatory decisions. Discrimination refers to disadvantageous treatment of a person based on belonging to a category rather than on individual merit. In this survey we review and organize various discrimination measures that have been used for measuring discrimination in data, as well as in evaluating performance of discrimination-aware predictive models. We also discuss related measures from other disciplines, which have not been used for measuring discrimination, but potentially could be suitable for this purpose. We computationally analyze properties of selected measures. We also review and discuss measuring procedures, and present recommendations for practitioners. The primary target audience is data mining, machine learning, pattern recognition, statistical modeling researchers developing new methods for non-discriminatory predictive modeling. In addition, practitioners and policy makers would use the survey for diagnosing potential discrimination by predictive models.

General Terms: fairness in machine learning, predictive modeling, non-discrimination, discriminationaware data mining

Discrimination discovery

Discrimination prevention

Educational solutions

- Algorithmic literacy to help people gain a broad understanding of the algorithmic 'value chain' (people perceive algorithms as unbiased)
- Educational programmes for raised awareness of discrimination and algorithmic bias

ACM Conference on Fairness, Accountability, and Transparency (ACM FAT*)

A computer science conference with a cross-disciplinary focus that brings together researchers and practitioners interested in fairness, accountability, and transparency in socio-technical systems.





Overview Career Opportunities

As we harness the power of AI, machine learning, and data science throughout many aspects of society and Microsoft systems and products, we need to consider the larger issues with AI.

Governance and transparency

- Clear governance procedures and algorithm accountability (role of CDO?)
- Developing an algorithmic audit trail –

"knowing very well the data collected [and the sources], identifying which pieces of data are used by algorithms, and made known how this data is weighted or used in the algorithm"

• Ethics and governance frameworks





Guidance Data Ethics Framework

Department for Digital, Culture, Media & Sport

Public sector organisations should use the Data Ethics Framework to guide the appropriate use of data to inform policy and service design.

Summary

na na

Roducion

Summary

- Decision-making increasingly data-driven
- Algorithms commonly employed to support or automate decisions and processes
- However, like humans, machines introduce biases and one area of current concern is algorithmic bias
- Algorithms and digital systems can perpetuate bias— one example being gender bias in image search
- Solutions include technical, educational and governance

Note: we can only expose, measure and try to reduce bias; we cannot completely eradicate it



Data. Insight. Action.

Questions?

WEAPONS OF Math destruction

CATHY O'NEIL

 Understand, Manage, and Prevent Algorithmic Bias

> A Guide for Business Users and Data Scientists

Tobias Baer

Apress*