

Lecture 19 Safety, Bias, and Fairness

Bias in the Vision and Language of Artificial Intelligence

Margaret Mitchell
Senior Research Scientist
Google AI

Andrew Zaldivar Me Simone Wu Parker Barnes Lucy Vasserman Ben Hutchinson Elena Spitzer Deb Raji Timnit Gebru

Adrian Benton Brian Zhang Dirk Hovy Josh Lovejoy Alex Beutel Blake Lemoine Hee Jung Ryu Hartwig Adam Blaise Agüera y Arcas

What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas
- Bananas with stickers on them
- Bunches of bananas with stickers on them on shelves in a store

...We don't tend to say
Yellow Bananas



What do you see?

Green Bananas

Unripe Bananas



What do you see?

Ripe Bananas

Bananas with spots

Bananas good for banana bread



What do you see?

Yellow Bananas

Yellow is prototypical for bananas



Prototype Theory

分类的目的之一是减少刺激行为和认知上可用的比例的无限差异

物品的一些核心、原型概念可能来自于存储的对象类别的典型属性(Rosch, 1975)

也可以存储范例(Wu & Barsalou, 2009)

A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

How could this be?

Doctor —— Female Doctor

大多数受试者忽视了医生是女性的可能性，包括男性、女性和自称女权主义者的人们。

World learning from text

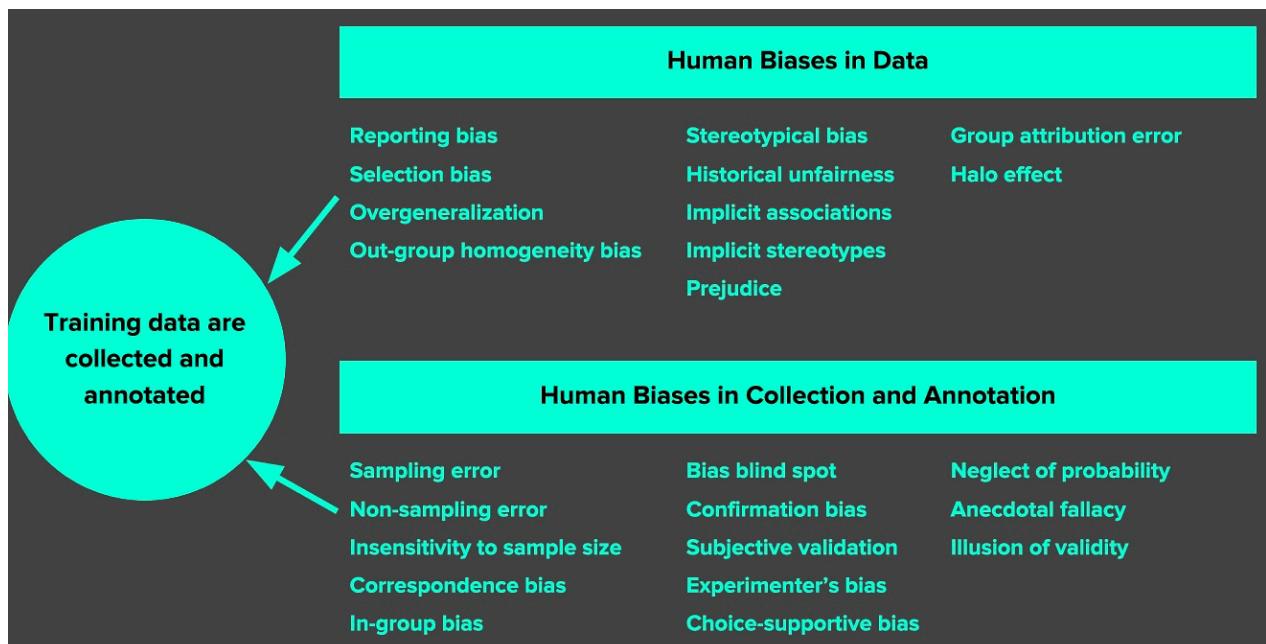
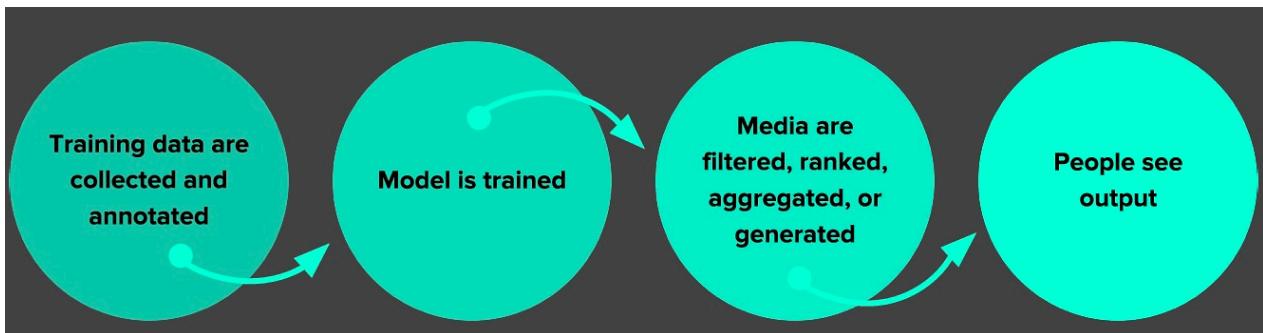
Gordon and Van Durme, 2013

Word	Frequency in corpus
“spoke”	11,577,917
“laughed”	3,904,519
“murdered”	2,834,529
“inhaled”	984,613
“breathed”	725,034
“hugged”	610,040
“blinked”	390,692
“exhale”	168,985

- murdered 是 blinked 出现次数的十倍
- 我们不倾向于提及眨眼和呼吸等事情

Human Reporting Bias: 人们写作中的行为、结果或属性的频率并不反映真实世界的频率，也不反映某一属性在多大程度上是某一类个体的特征。

更多关于我们处理世界和我们认为非凡的东西的实际情况。这影响到我们学习的一切。



- 数据

- Reporting bias 报告偏见:人们分享的并不是真实世界频率的反映
- Selection Bias 选择偏差:选择不反映随机样本
- Out-group homogeneity bias 外群体同质性偏见:People tend to see outgroup members as more alike than ingroup members when comparing attitudes, values, personality traits, and other characteristics

- 解释

- Confirmation bias 确认偏见:倾向于寻找、解释、支持和回忆信息, 以确认一个人先前存在的信念或假设
- Overgeneralization 泛化过度:根据过于笼统和/或不够具体的信息得出结论
- Correlation fallacy 相关性谬误:混淆相关性和因果关系
- Automation bias 自动化偏差:人类倾向于喜欢来自自动化决策系统的建议, 而不是没有自动化的相互矛盾的信息

Biases in Data



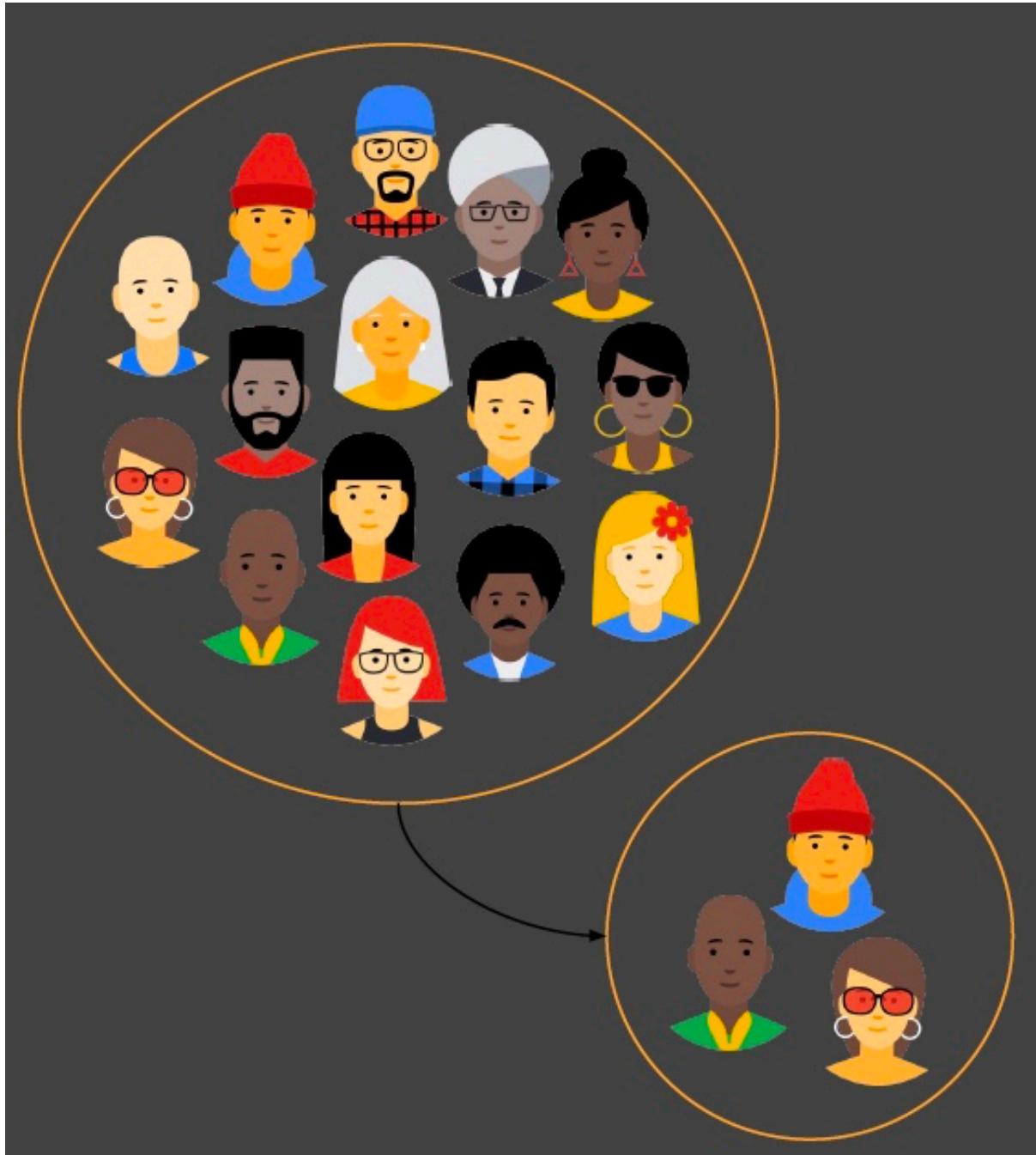
Selection Bias 选择偏差:选择不反映随机样本



Out-group homogeneity bias 外群体同质性偏见：在比较态度、价值观、个性特征和其他特征时，往往群体外的成员认为比群体内的成员更相似

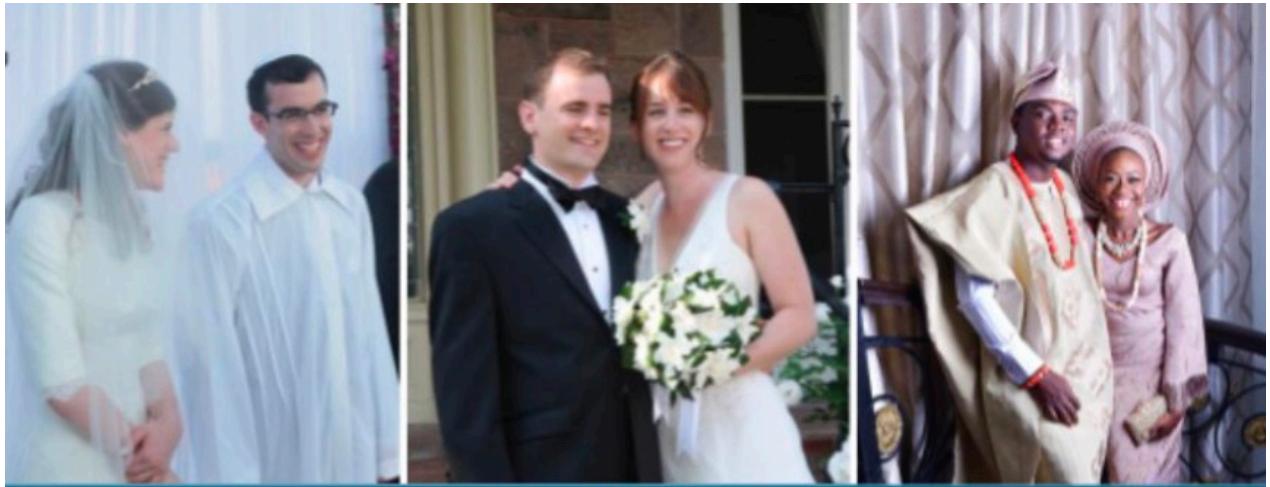
- 这有些难以理解：意思就是左边的四只猫之间是非常不同的，但是在狗的眼里他们是相同的

Biases in Data → Biased Data Representation



你可能对你能想到的每一个群体都有适当数量的数据，但有些群体的表现不如其他群体积极。

Biases in Data → Biased Labels



*ceremony,
wedding, bride,
man, groom,
woman, dress*

*ceremony,
bride, wedding,
man, groom,
woman, dress*

person, people

数据集中的注释将反映注释者的世界观

Biases in Interpretation



Confirmation bias 确认偏见:倾向于寻找、解释、支持和回忆信息,以确认一个人先前存在的信念或假设



Overgeneralization 泛化过度:根据过于笼统和/或不够具体的信息得出结论 (相关: 过拟合)

Post Hoc Ergo Propter Hoc

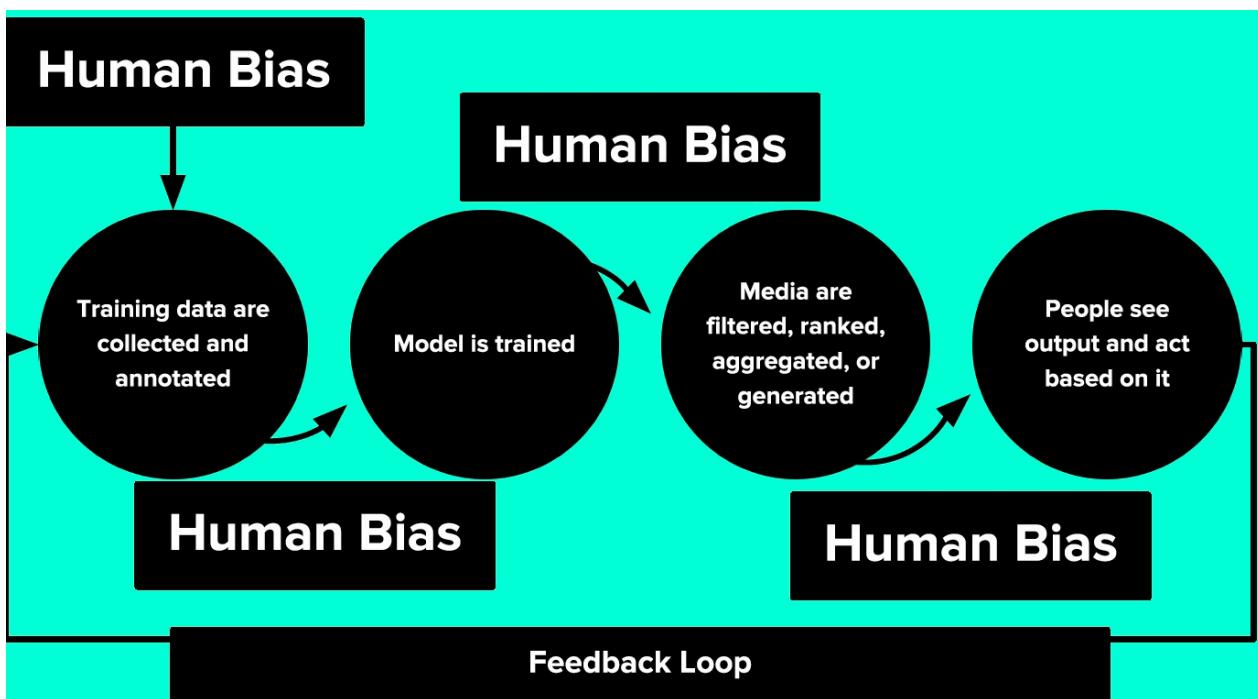
Women were allowed to vote in the early 1900's and then we had two world wars. Clearly giving them the vote was a bad idea.



Correlation fallacy 相关性谬误:混淆相关性和因果关系



Automation bias 自动化偏差:人类倾向于喜欢来自自动化决策系统的建议, 而不是没有自动化的相互矛盾的信息



- 会形成反馈循环
- 这被称为 Bias Network Effect 以及 Bias “Laundering”

人类数据延续了人类的偏见。当ML从人类数据中学习时，结果是一个偏置网络效应。

BIAS = BAD ??

“Bias” can be Good, Bad, Neutral

- 统计以及 ML中的偏差
 - 估计值的偏差：预测值与我们试图预测的正确值之间的差异
 - “偏差”一词b(如 $y = mx + b$)
- 认知偏见
 - 确认性偏差、近因性偏差、乐观性偏差
- 算法偏差
 - 对与种族、收入、性取向、宗教、性别和其他历史上与歧视和边缘化相关的特征相关的人的不公平、不公平或偏见待遇，何时何地在算法系统或算法辅助决策中体现出来”

*“Although neural networks might be said to write their own programs, they do so towards **goals set by humans, using data collected for human purposes**. If the data is skewed, even by accident, the computers will **amplify injustice.**”*

— The Guardian

- 如何避免算法偏差，开发出不会放大差异的算法

Predicting Future Criminal Behavior

- 算法识别潜在的犯罪热点
- 基于之前报道的犯罪的地方，而不是已知发生在哪里
- 从过去预测未来事件
- 预测的是逮捕的地方而不是犯罪的地方

Predicting Sentencing

- Prater (白人)额定低风险入店行窃后,尽管两个武装抢劫;一次持械抢劫未遂。
- Borden (黑色)额定高危后她和一个朋友(但在警察到来之前返回)一辆自行车和摩托车坐在外面。
- 两年后, Borden没有被指控任何新的罪行。Prater因重大盗窃罪被判8年有期徒刑。

系统默认认为黑人的犯罪风险高于拜仁

Automation Bias in face of:

- Overgeneralization
- Feedback Loops
- Correlation Fallacy

Predicting Criminality

以色列启动 Faception

Faception是第一个科技领域的率先面市的，专有的计算机视觉和机器学习技术分析人员和揭示他们的个性只基于他们的面部图像。

提供专业的引擎从脸的形象识别“高智商”、“白领犯罪”、“恋童癖”,和“恐怖分子”。

主要客户为国土安全和公共安全。

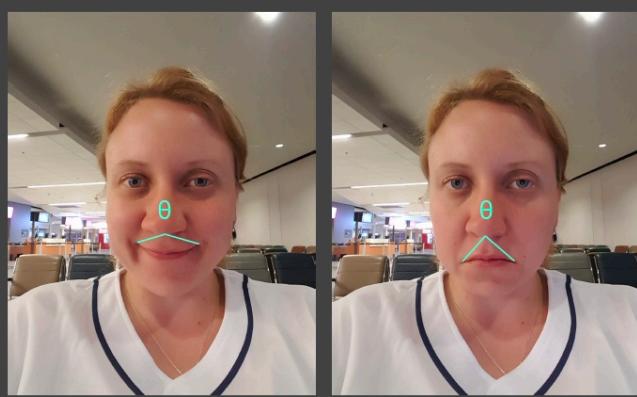
["Automated Inference on Criminality using Face Images"](#) Wu and Zhang, 2016. arXiv

1856个紧密裁剪的面孔的图像,包括“通缉犯”ID特定区域的照片

存在确认偏差和相关性偏差

1,856 closely cropped images of faces;
Includes “wanted suspect” ID pictures
from specific regions.

*[...] angle θ from nose tip to two
mouth corners is on average 19.6%
smaller for criminals than for
non-criminals ...”*



See our longer piece on Medium, [“Physiognomy’s New Clothes”](#)

- Selection Bias + Experimenter’s Bias + Confirmation Bias + Correlation Fallacy + Feedback Loops

Predicting Criminality - The Media Blitz

arXiv Paper Spotlight: Automated Inference on Criminality Using Face ...

www.kdnuggets.com/.../arxiv-spotlight-automated-inference-criminality-face-images.... ▾

A recent paper by Xiaolin Wu (McMaster University, Shanghai Jiao Tong University) and Xi Zhang (Shanghai Jiao Tong University), titled "Automated Inference ...

Automated Inference on Criminality Using Face Images | Hacker News

<https://news.ycombinator.com/item?id=12983827> ▾

Nov 18, 2016 - The automated inference on criminality eliminates the variable of meta-accuracy (the competence of the human judge/examiner) all together.

A New Program Judges If You're a Criminal From Your Facial Features ...

<https://motherboard.vice.com/.../new-program-decides-criminality-from-facial-feature...> ▾

Nov 18, 2016 - In their paper 'Automated Inference on Criminality using Face Images', published on the arXiv pre-print server, Xiaolin Wu and Xi Zhang from ...

Can face classifiers make a reliable inference on criminality?

<https://techxplore.com/.../computer-sciences...> ▾

Nov 23, 2016 - Their paper is titled "Automated Inference on Criminality using Face Images ... face classifiers are able to make reliable inference on criminality.

Troubling Study Says Artificial Intelligence Can Predict Who Will Be ...

<https://theintercept.com/.../troubling-study-says-artificial-intelligence-can-predict-who...> ▾

Nov 18, 2016 - Not so in the modern age of Artificial Intelligence, apparently: In a paper titled "Automated Inference on Criminality using Face Images," two ...

Automated Inference on Criminality using Face Images (via arXiv ...)

<https://computationallegalstudies.com/.../automated-inference-on-criminality-using-fa...> ▾

Dec 6, 2016 - Next Next post: A General Approach for Predicting the Behavior of the Supreme Court of the United States (Paper Version 2.01) (Katz, ...

(Claiming to) Predict Internal Qualities Subject To Discrimination

Predicting Homosexuality

- Wang and Kosinski, [Deep neural networks are more accurate than humans at detecting sexual orientation from facial images](#), 2017.
- “Sexual orientation detector” using 35,326 images from public profiles on a US dating website.
- “与性取向的产前激素理论(PHT)相一致，男同性恋者和女同性恋者往往具有非典型的性别面部形态。”
- 在自拍中，同性恋和异性恋之间的差异与打扮、表现和生活方式有关，也就是说，文化差异，而不是面部结构的差异

See our longer response on Medium, “[Do Algorithms Reveal Sexual Orientation or Just Expose our Stereotypes?](#)”

- Selection Bias + Experimenter’s Bias + Correlation Fallacy

Measuring Algorithmic Bias

评估公平性和包容性

- 分类评估
 - 为每个创建（子组，预测）对。跨子组比较
 - 例如
 - 女性，面部检测
 - 男性，面部检测
- 交叉评估
 - 为每个创建（子组1，子组2，预测）对。跨子组比较
 - 例如
 - 黑人女性，面部检测
 - 白人，面部检测

Evaluate for Fairness & Inclusion: Confusion Matrix

		Model Predictions		
		Positive	Negative	
References	Positive	<ul style="list-style-type: none"> • Exists • Predicted <p>True Positives</p>	<ul style="list-style-type: none"> • Exists • Not predicted <p>False Negatives</p>	Recall, False Negative Rate
	Negative	<ul style="list-style-type: none"> • Doesn’t exist • Predicted <p>False Positives</p>	<ul style="list-style-type: none"> • Doesn’t exist • Not predicted <p>True Negatives</p>	
		<p>Precision, False Discovery Rate</p>		Negative Predictive Value, False Omission Rate
				LR+, LR-

Evaluate for Fairness & Inclusion

Female Patient Results		Male Patient Results	
True Positives (TP) = 10	False Positives (FP) = 1	True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 1	True Negatives (TN) = 488	False Negatives (FN) = 5	True Negatives (TN) = 48
Precision = $\frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$		Precision = $\frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$	
Recall = $\frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$		Recall = $\frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$	
Recall = $\frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$		Recall = $\frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$	
“Equality of Opportunity” fairness criterion: Recall is equal across subgroups			

- “机会平等”公平准则：子组的 recall 是相等的

Precision = $\frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$	Precision = $\frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$
Recall = $\frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$	Recall = $\frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$
“Predictive Parity” fairness criterion: Precision is equal across subgroups	

- “预测平价”公平准则：子组的 precision 是相等

选择评价指标的可接受的假阳性和假阴性之间的权衡

False Positives Might be Better than False Negatives

- Privacy in Images

False Positive: Something that doesn't need to be blurred gets blurred.

Can be a bummer.



False Negative: Something that needs to be blurred is not blurred.

Identity theft.



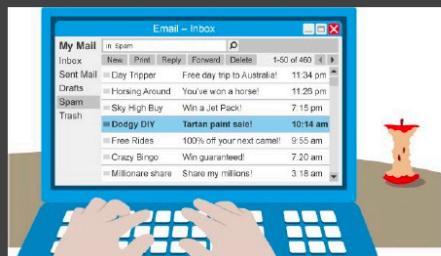
- Spam Filtering

False Negative: Email that is SPAM is not caught, so you see it in your inbox.

Usually just a bit annoying.

False Positive: Email flagged as SPAM is removed from your inbox.

If it's from a friend or loved one, it's a loss!

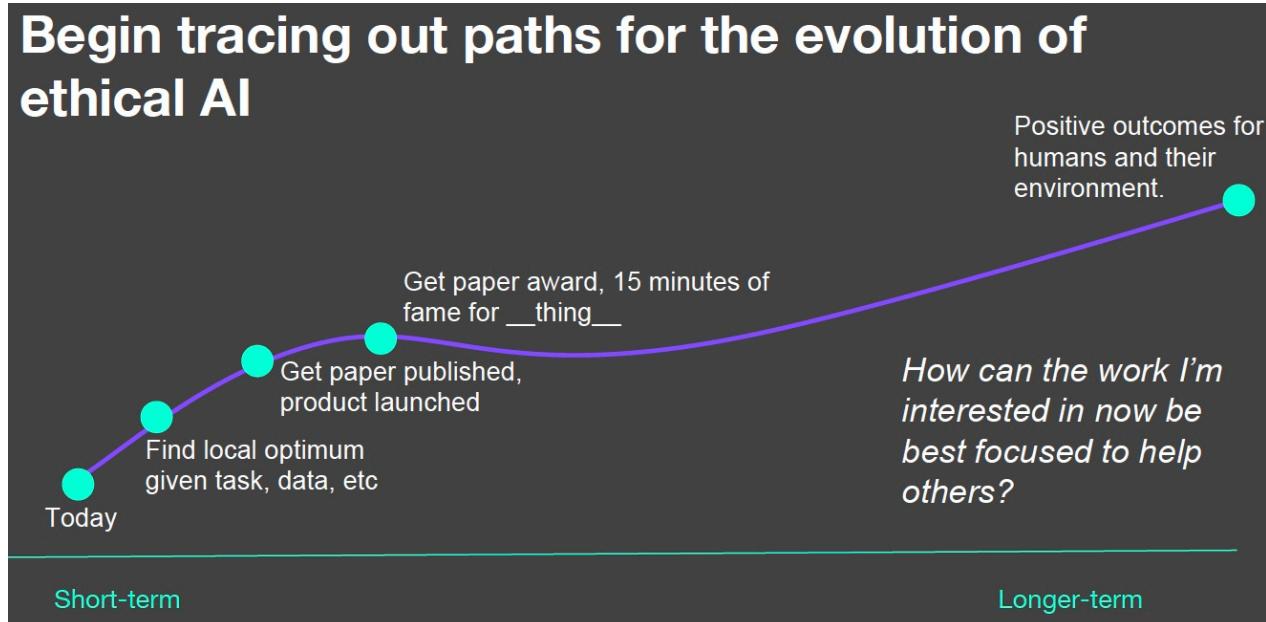


AI Can Unintentionally Lead to Unjust Outcomes

- 缺乏对数据和模型中的偏见来源的洞察力
- 缺乏对反馈循环的洞察力
- 缺乏细心，分类的评价
- 人类偏见在解释和接受结果

It's up to us to influence how AI evolves.

Begin tracing out paths for the evolution of ethical AI



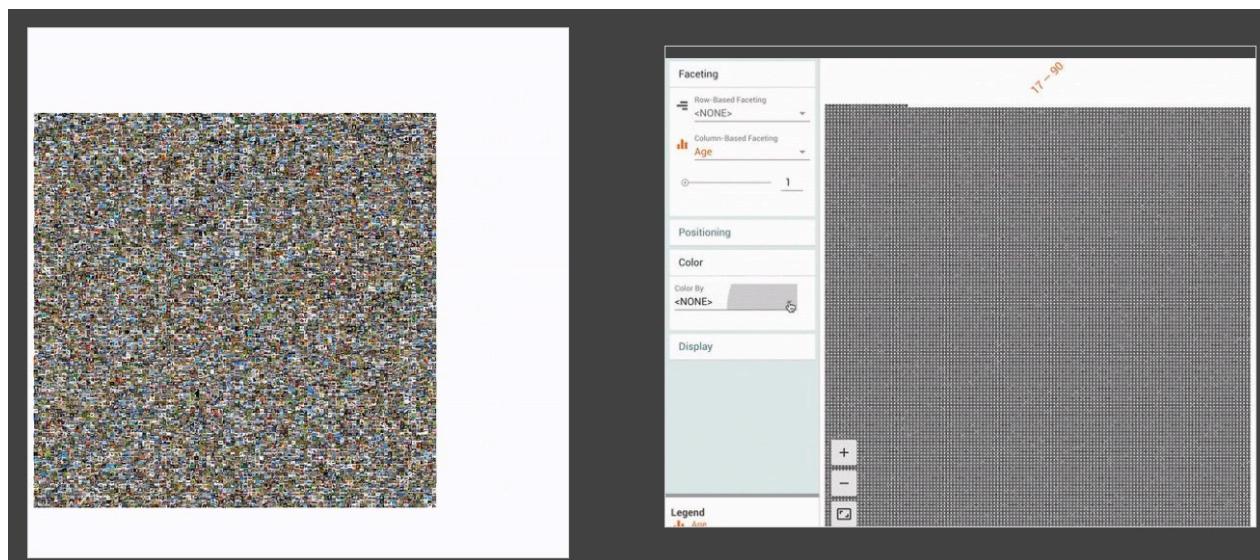
Here are some things we can do.

Data

Data Really, Really Matters

- 了解您的数据：偏差，相关性
- 从类似的分布放弃单一训练集/测试集
- 结合来自多个来源的输入
- 对于困难的用例使用held-out测试集
- 与专家讨论其他信号

Understand Your Data Skews



Datasheets for Datasets

Timnit Gebru¹ Jamie Morgenstern² Briana Vecchione³ Jennifer Wortman Vaughan¹ Hanna Wallach¹
Hal Daumé III^{1,4} Kate Crawford^{1,5}

没有一个数据集是没有偏差的，因为这是一个有偏差的世界。重点是知道是什么偏差。

Machine Learning

Use ML Techniques for Bias Mitigation and Inclusion

Bias Mitigation 偏差缓解

- 删掉有问题的输出的信号
 - 刻板印象
 - 性别歧视,种族歧视,*-ism
 - 又称为“debiasing”

Inclusion

- 添加信号所需的变量
 - 增加模型性能

- 注意性能很差的子组或数据片

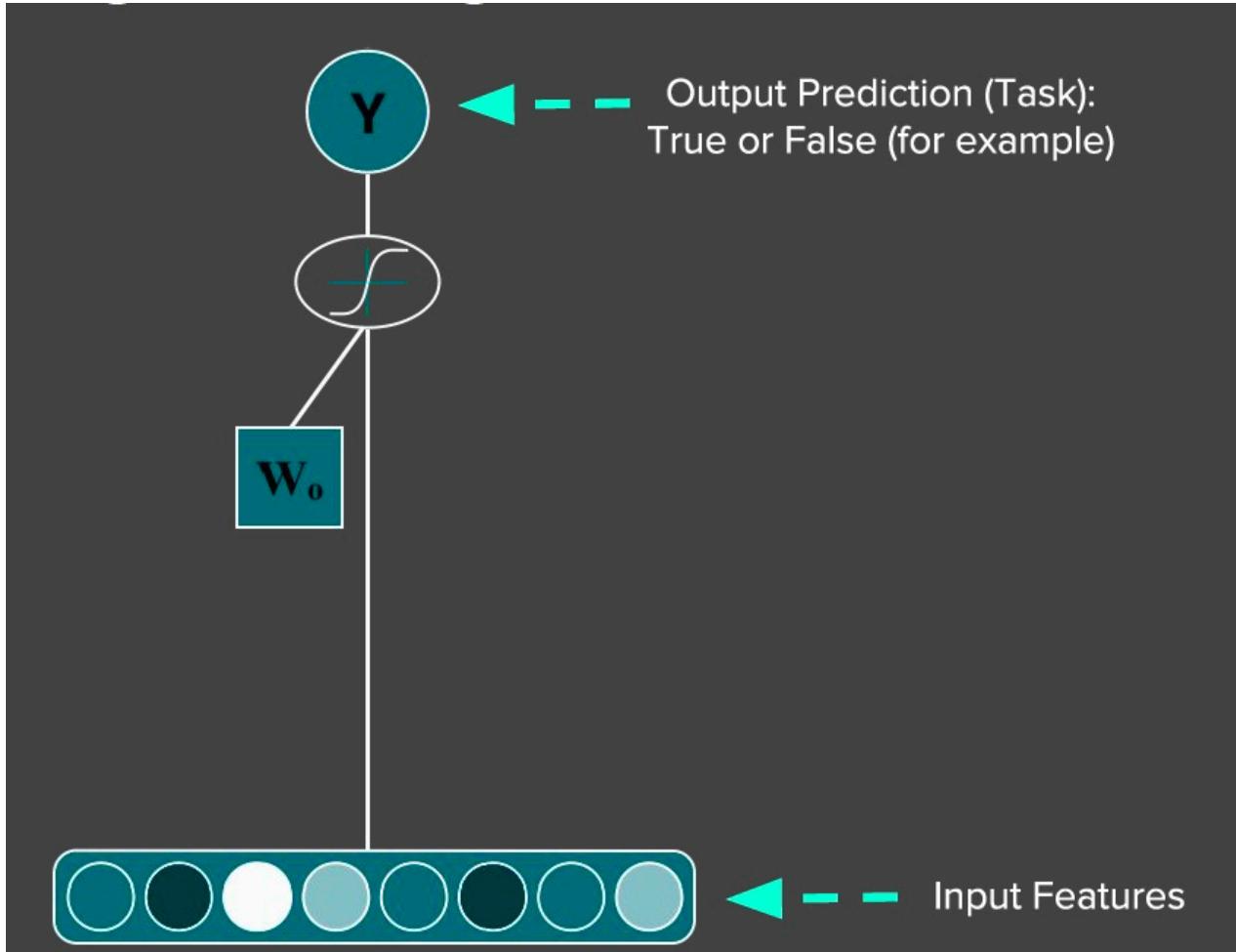
Multi-task Learning to Increase Inclusion

Multiple Tasks + Deep Learning for Inclusion: Multi-task Learning Example

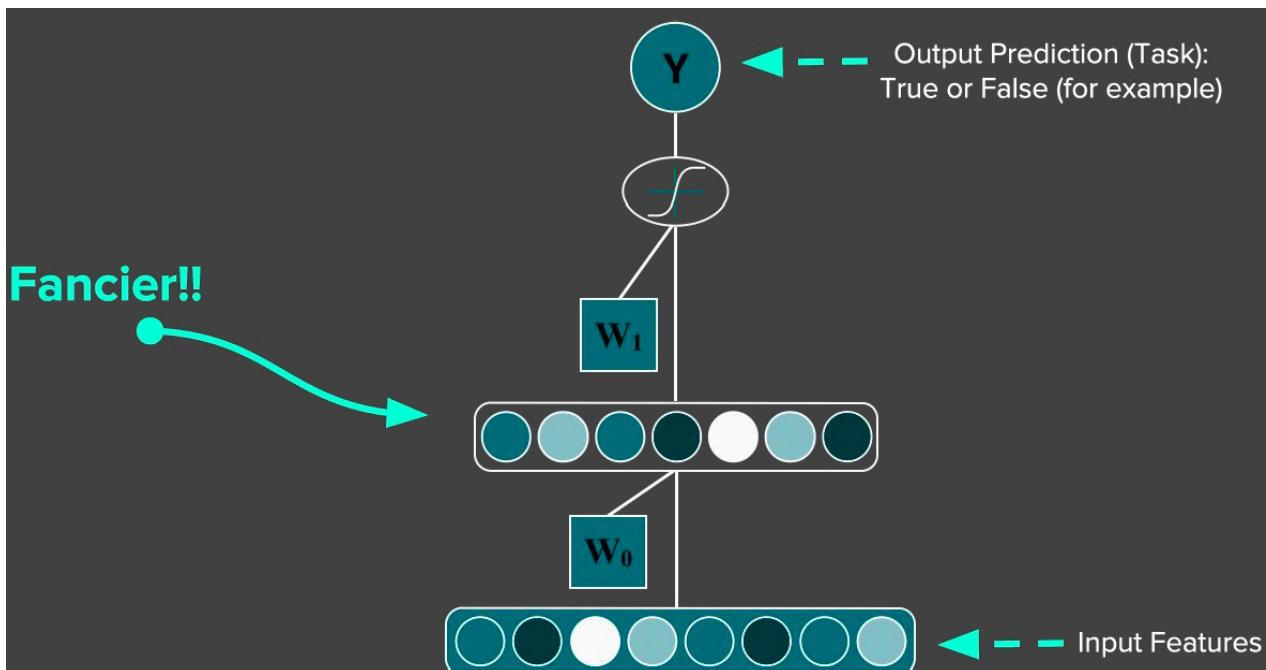
Benton, Mitchell, Hovy. [Multi-task learning for Mental Health Conditions with Limited Social Media Data](#). EACL, 2017.

- 与宾夕法尼亚大学WWP合作
- 直接与临床医生合作
- 目标
 - 系统警告临床医生如果企图自杀迫在眉睫
 - 几个训练实例可用时诊断的可行性
- 内部数据
 - 电子健康记录
 - 病人或病人家属提供
 - 包括心理健康诊断,自杀企图,竞赛
 - 社交媒体数据
- 代理数据
 - Twitter 媒体数据
 - 代理心理健康诊断中使用自称诊断
 - 我被诊断出患有 X
 - 我试图自杀

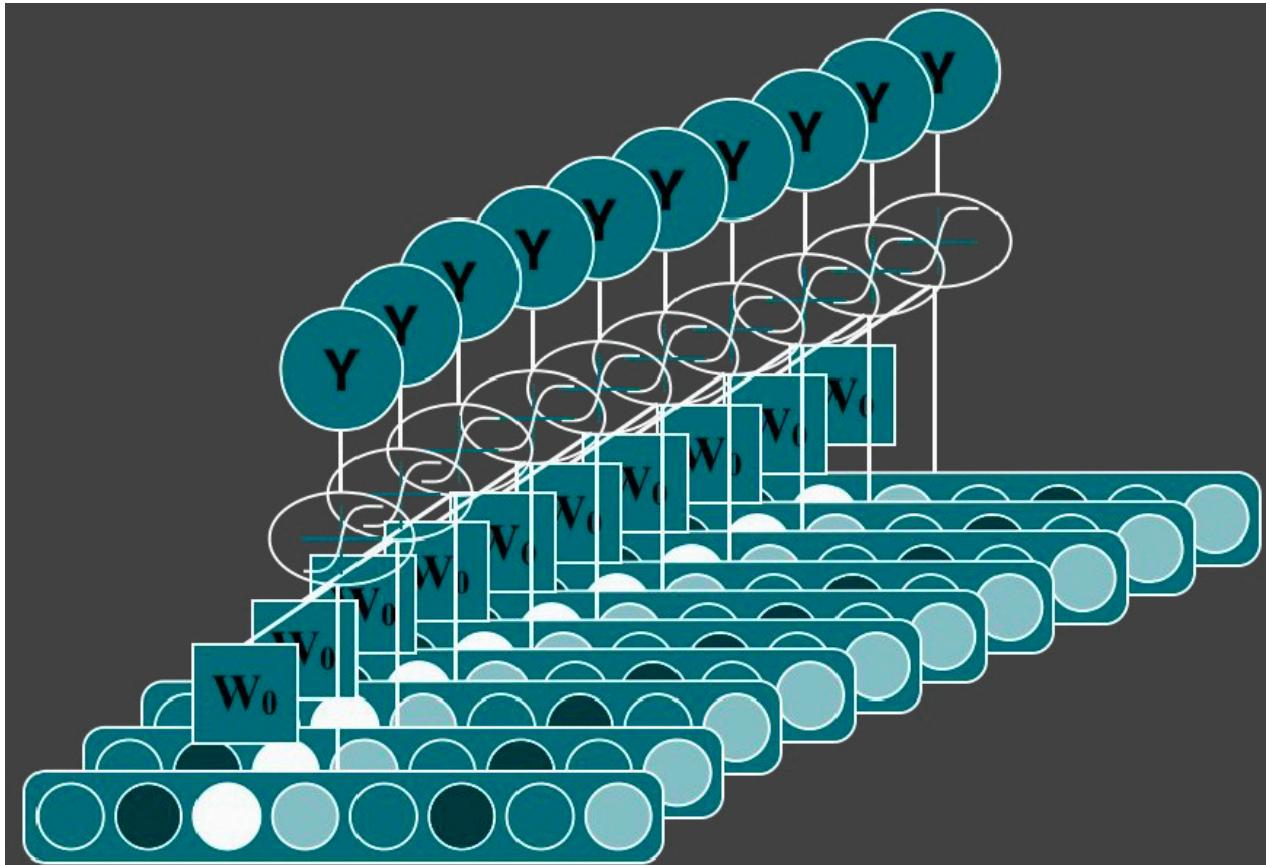
Single-Task: Logistic Regression



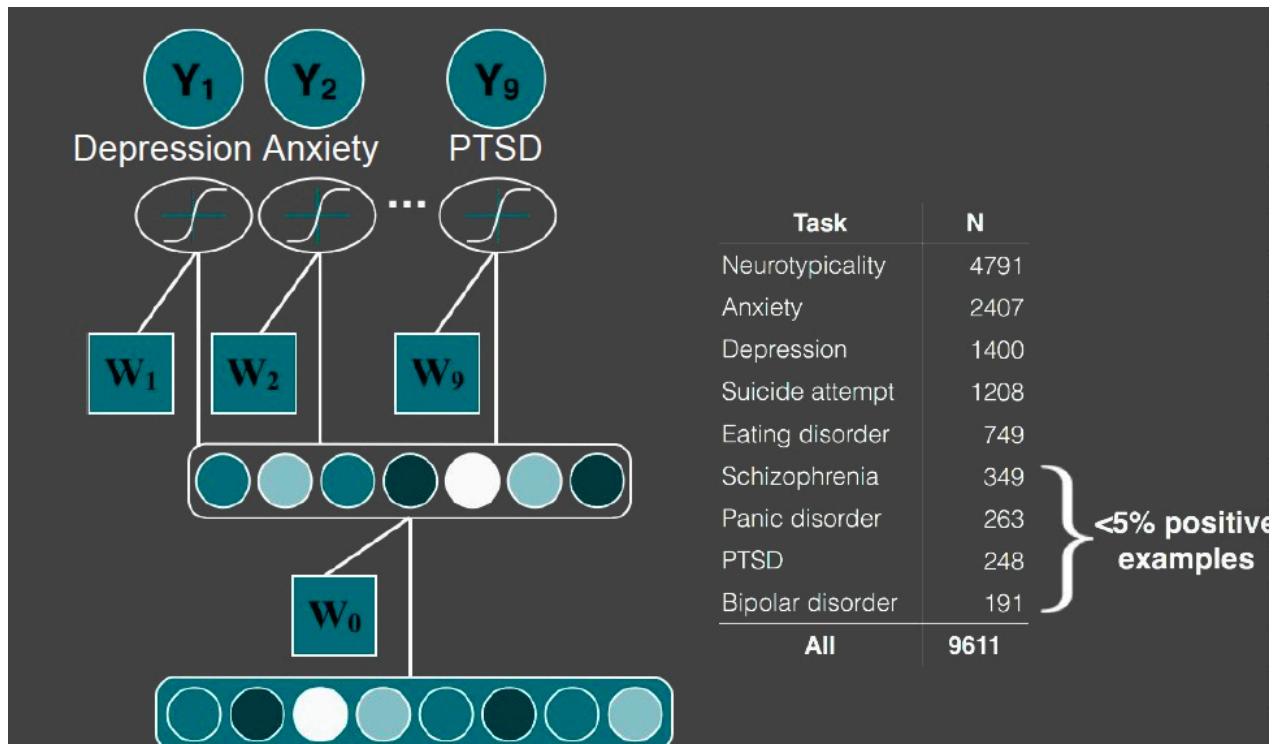
Single-Task: Deep Learning



Multiple Tasks with Basic Logistic Regression



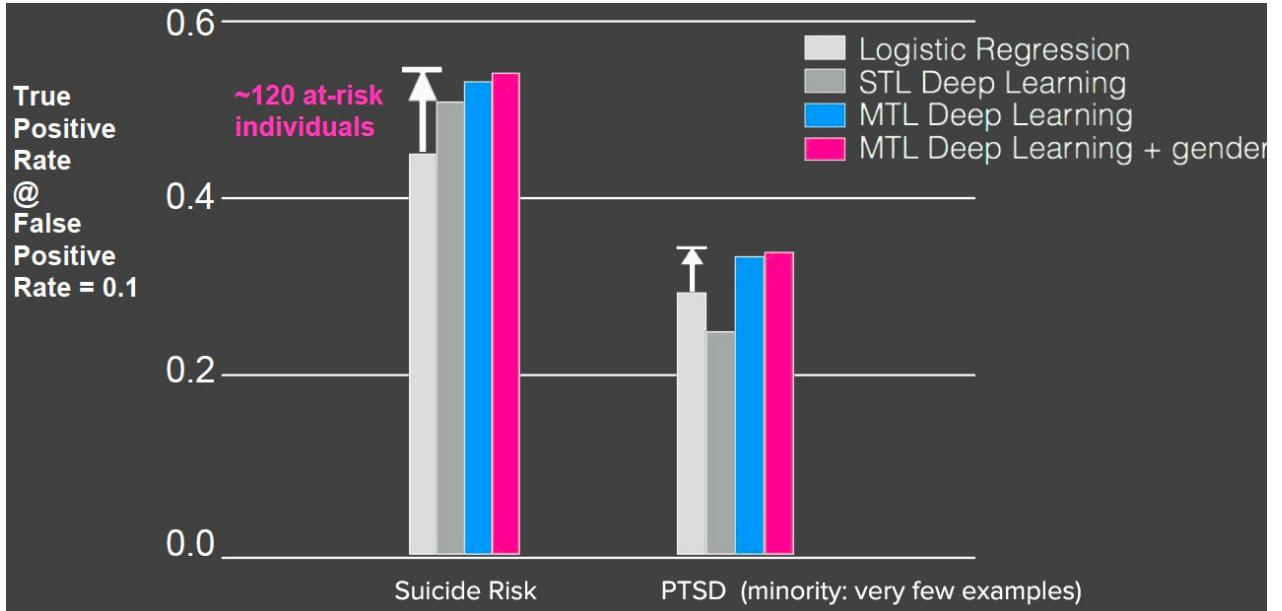
Multi-task Learning



Multitask, given comorbidity

Benton, Mitchell, Hovy. [Multi-task learning for Mental Health Conditions with Limited Social Media Data](#). EACL, 2017.

Improved Performance across Subgroups



Reading for the masses....

Multi-Task Learning for Mental Health using Social Media Text

Adrian Benton
Johns Hopkins University
adrian@cs.jhu.edu

Margaret Mitchell
Microsoft Research*
mmitchellai@google.com

Dirk Hovy
University of Copenhagen
mail@dirkhovy.com

2 Ethical Considerations

As with any author-attribute detection, there is the danger of abusing the model to single out people (*overgeneralization*, see Hovy and Spruit (2016)). We are aware of this danger, and sought to minimize the risk. For this reason, we don't provide a selection of features or representative examples. The experiments in this paper were performed with a clinical application in mind, and use carefully matched (but anonymized) data, so the distribution is not representative of the population as a whole. The results of this paper should therefore *not* be interpreted as a means to assess mental health conditions in social media in general, but as a test for the applicability of MTL in a well-defined clinical setting.

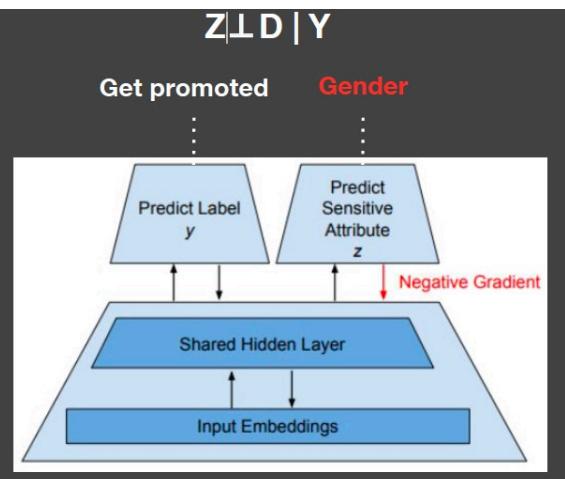
Contextualizing and considering ethical dimensions

Adversarial Multi-task Learning to Mitigate Bias

Multitask Adversarial Learning

- Basic idea: Jointly predict:
 - Output decision D
 - Attribute you'd like to remove from decision Z
 - Negate the effect of the undesired attribute

$$P(\hat{Y} = 1 | Y = 1, Z = 1) = P(\hat{Y} = 1 | Y = 1, Z = 0)$$



Beutel, Chen, Zhao, Chi. [Data Decisions and Theoretical Implications when Adversarially Learning Fair Representations](#). FAT/ML, 2017.

Zhang, Lemoine, Mitchell. [Mitigating Unwanted Biases with Adversarial Learning](#). AIES, 2018.

Equality of Opportunity in Supervised Learning

考虑到真正正确的决策，分类器的输出决策应该在敏感特征之间是相同的。

Case Study: Conversation AI Toxicity

Measuring and Mitigating Unintended Bias in Text Classification

Lucas Dixon
ldixon@google.com

John Li
jetpack@google.com

Jeffrey Sorensen
sorenj@google.com

Nithum Thain
nthain@google.com

Lucy Vasserman
lucyvasserman@google.com



AIES, 2018 and FAT*, 2019

- Conversation-AI
 - ML 提高在线对话的规模
- Research Collaboration
 - Jigsaw, CAT, several Google-internal teams, and external partners (NYTimes, Wikimedia, etc)

Perspective API

“You’re a dork!”
Toxicity: 0.91



Data + ML
Toxicity,
Severe Toxicity,
Threat, Off-topic,
+ dozens other
models

<http://perspectiveapi.com/>

Unintended Bias

模型错误地将经常受到攻击的身份与毒性联系起来：False Positive Bias

<u>Sentence</u>	<u>model score</u>
"i'm a proud tall person"	0.18
"i'm a proud lesbian person"	0.51
"i'm a proud gay person"	0.69

Bias Source and Mitigation

- 偏见造成的数据不平衡
 - 经常袭击了有毒的身份所占比例评论长度问题
- 添加维基百科文章中假定的无毒数据来修复这种不平衡
 - 原始数据集有127820个例子
 - 4620个补充的无毒例子

Term	Comment Length				
	20-59	60-179	180-539	540-1619	1620-4859
ALL	17%	12%	7%	5%	5%
gay	88%	77%	51%	30%	19%
queer	75%	83%	45%	56%	0%
homosexual	78%	72%	43%	16%	15%
black	50%	30%	12%	8%	4%
white	20%	24%	16%	12%	2%
wikipedia	39%	20%	14%	11%	7%
atheist	0%	20%	9%	6%	0%
lesbian	33%	50%	42%	21%	0%
feminist	0%	20%	25%	0%	0%
islam	50%	43%	12%	12%	0%
muslim	0%	25%	21%	12%	17%
race	20%	25%	12%	10%	6%
news	0%	1%	4%	3%	3%
daughter	0%	7%	0%	7%	0%

Measuring Unintended Bias - Synthetic Datasets

挑战与真实数据

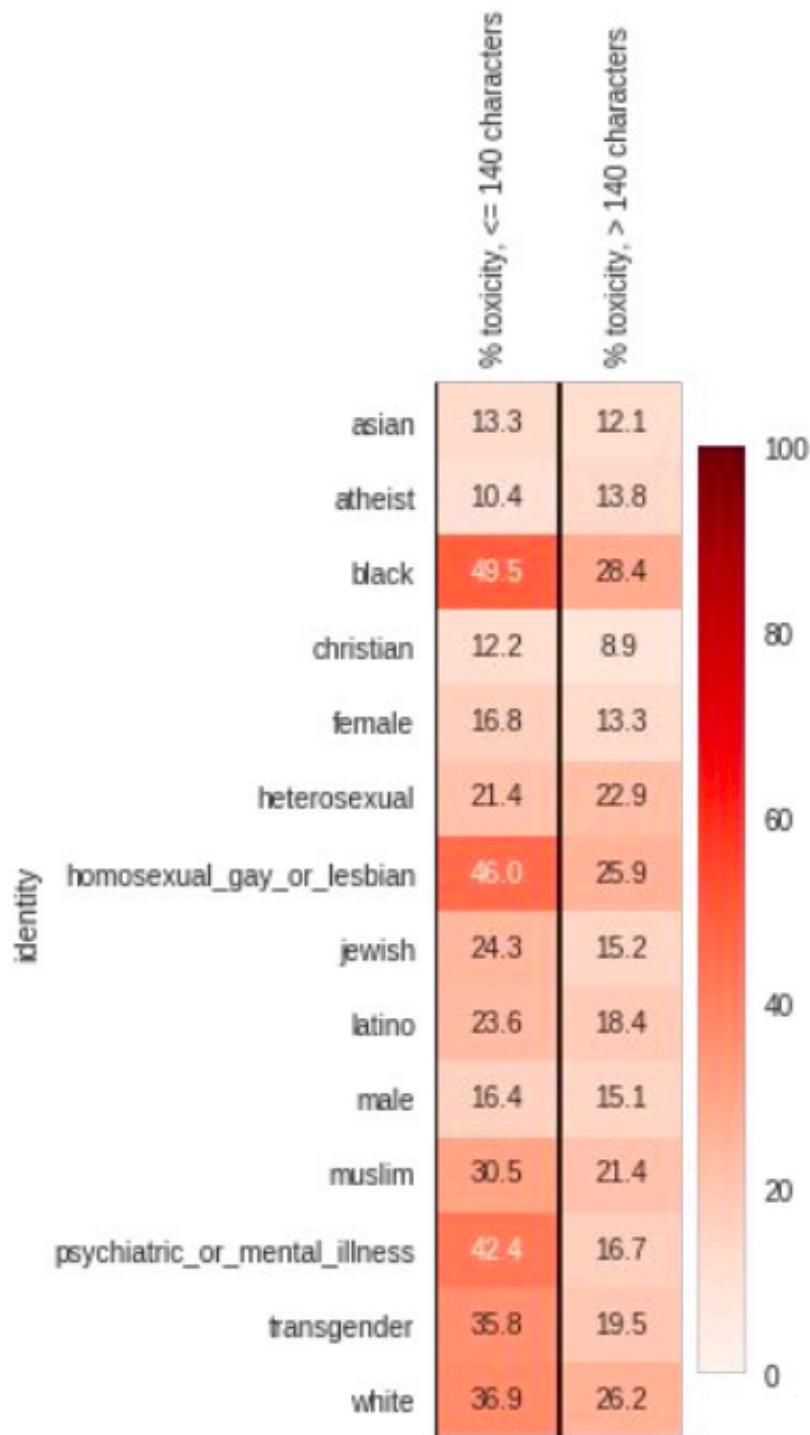
- 现有数据集是小 and/or 有错误的相关性

- 每个例子是完全独特的

Approach: "bias madlibs": 一个综合生成的模板化数据集进行评估

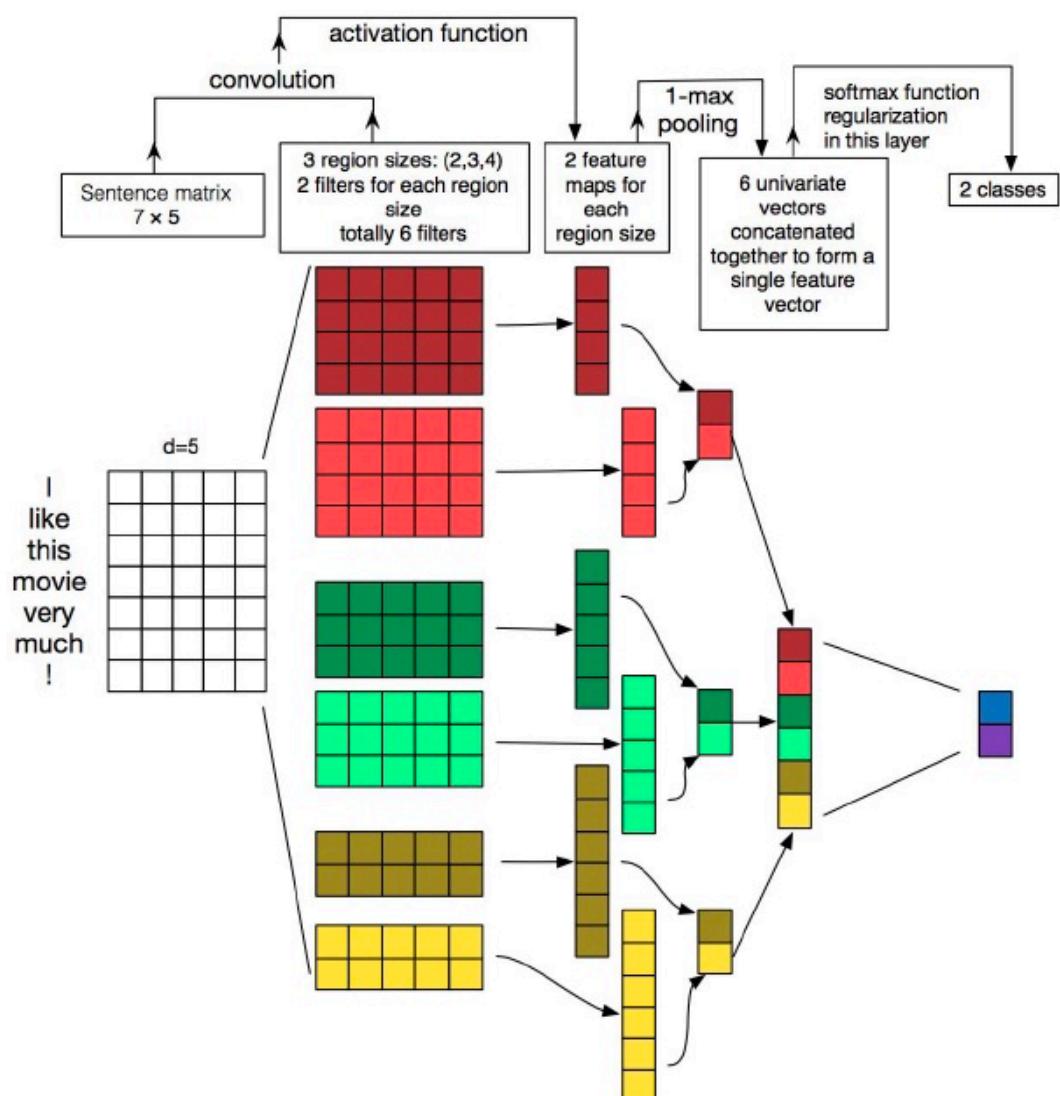
<u>Sentence</u>	<u>model score</u>
"i'm a proud tall person"	0.18
"i'm a proud lesbian person"	0.51
"i'm a proud gay person"	0.69
"audre is a brazilian computer programmer"	0.02
"audre is a muslim computer programmer"	0.08
"audre is a transgender computer programmer"	0.56

Assumptions



- 数据集是可靠的
 - 和产品相似的分布
 - 忽略注释器偏见
 - 没有因果分析

Deep Learning Model



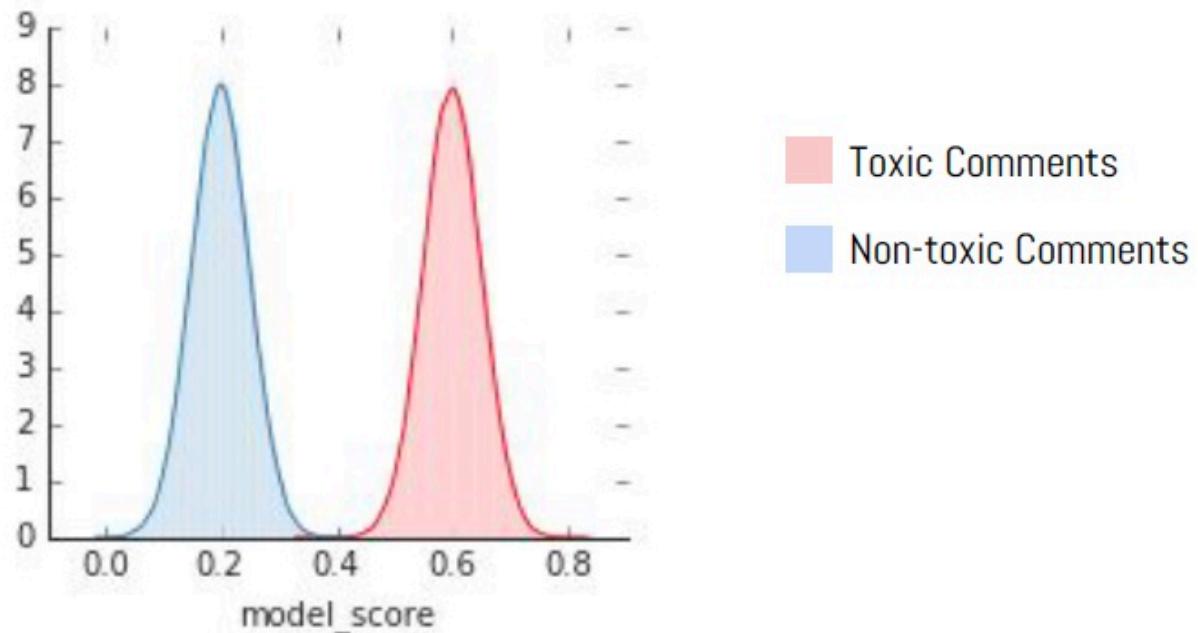
Source: Zhang, Y., & Wallace, B. (2015). A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. arXiv preprint arXiv:1510.03820.

- CNN 架构
- 预训练的 GloVe 嵌入
- Keras 实现

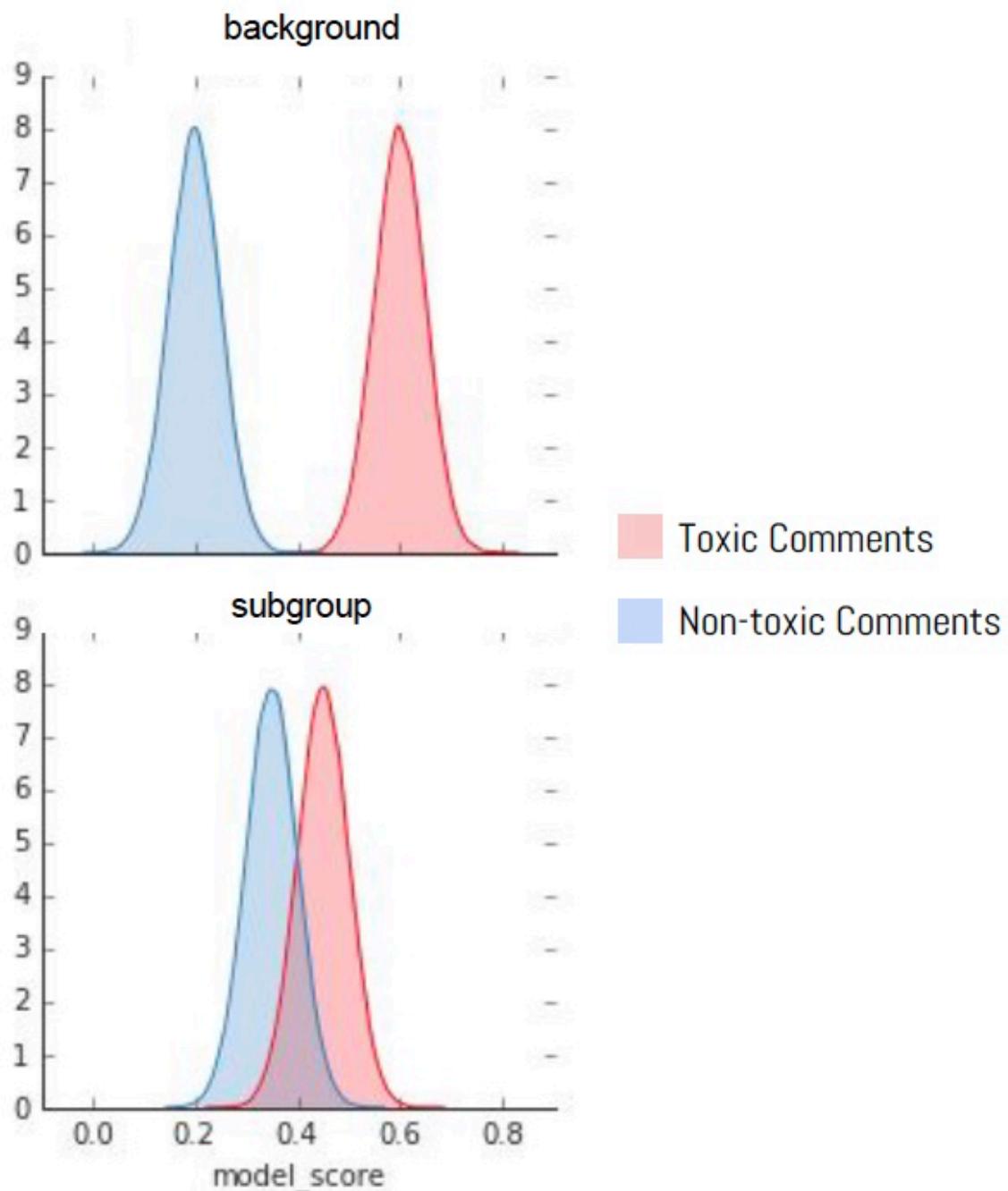
Measuring Model Performance

How good is the model at distinguishing good from bad examples? (ROC-AUC)

AUC (for a given test set) = Given two randomly chosen examples, one in-class (e.g. one is toxic and the other is not), AUC is the probability that the model will give the in-class example the higher score.



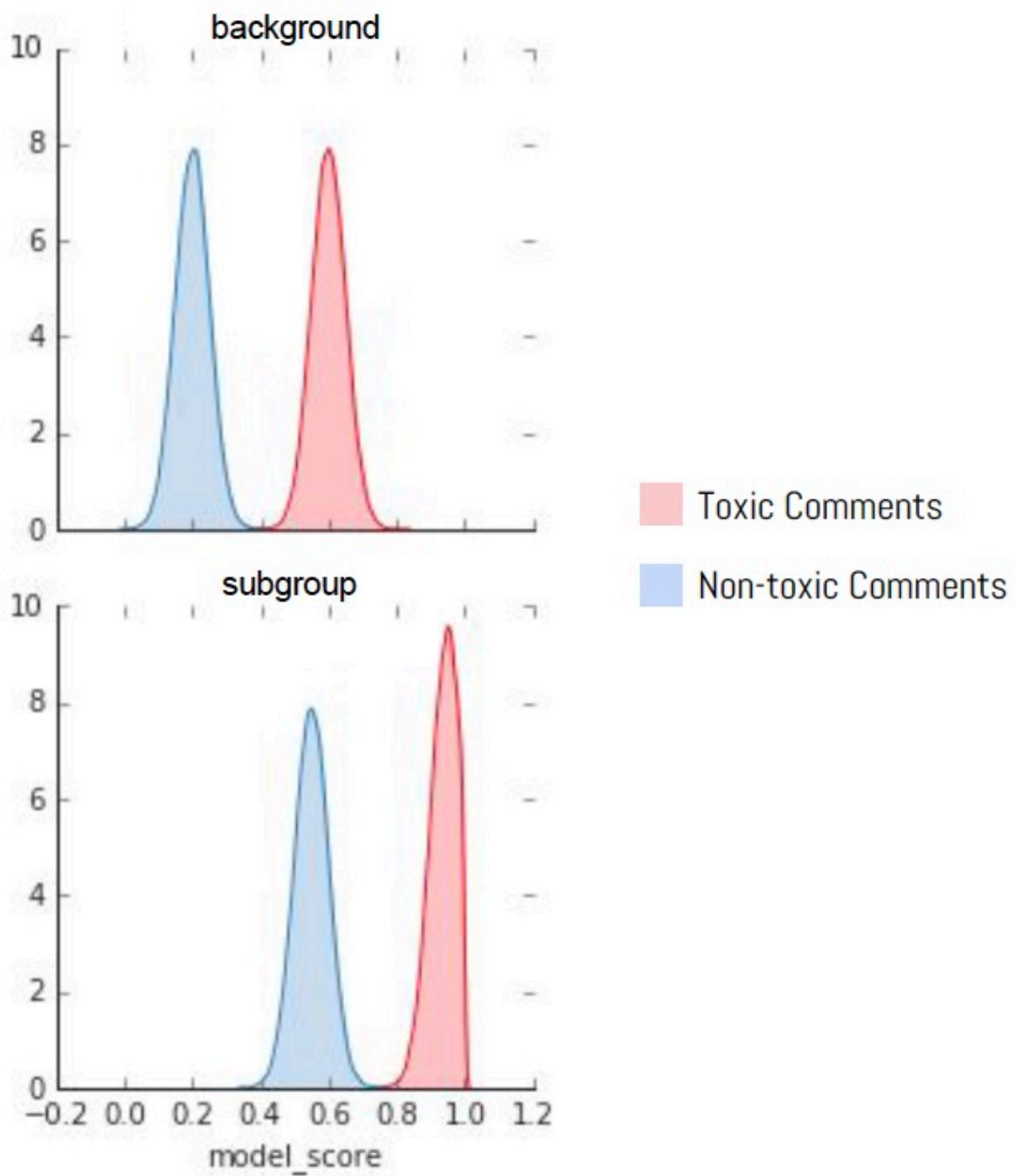
Types of Bias



Low Subgroup Performance

模型在子组注释上的性能比在总体注释上差

Metric : Subgroup AUC

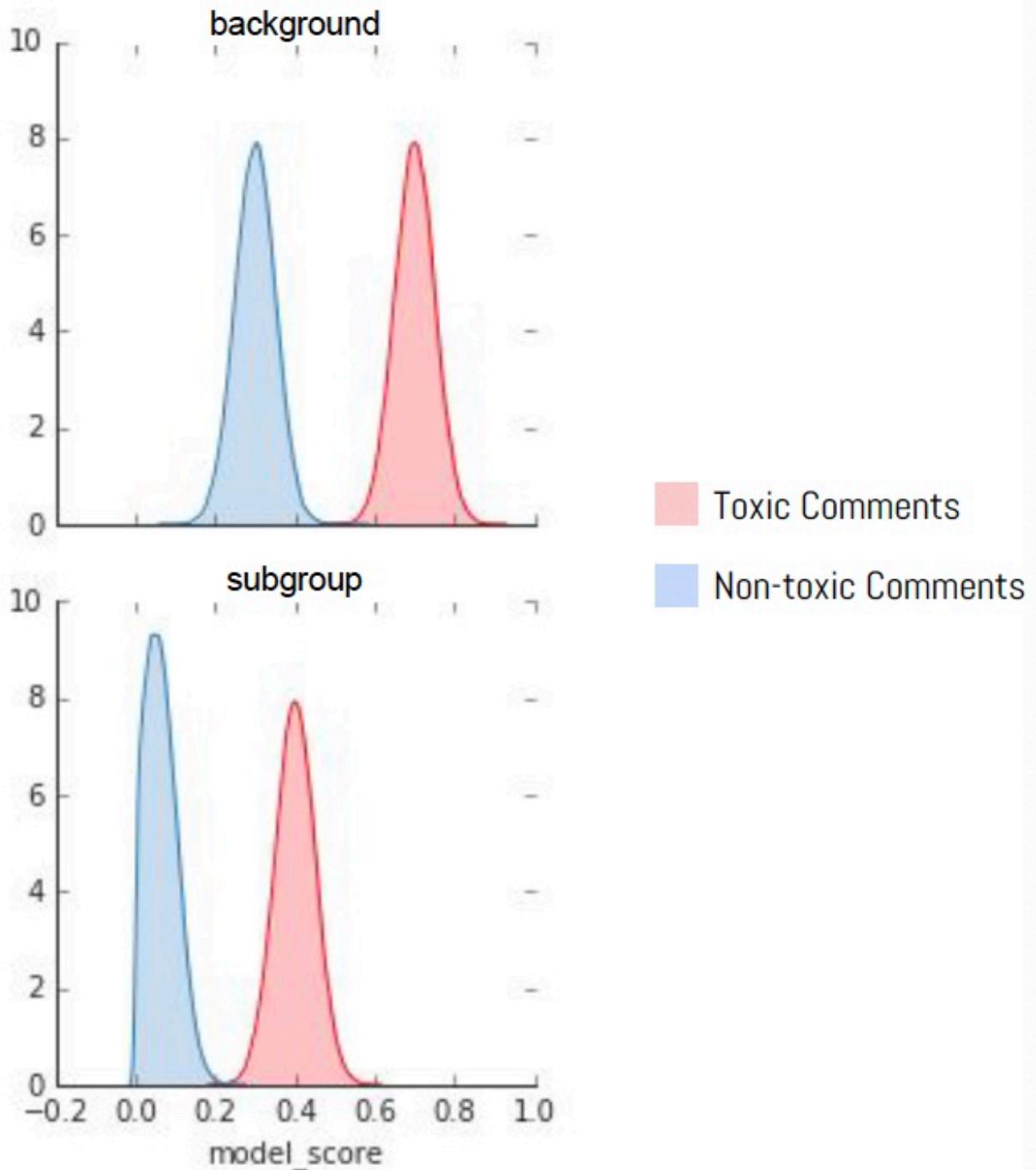


Subgroup Shift (Right)

该模型系统地对来自子组的评价打分更高

Metric: BPSN AUC

(Background Positive Subgroup Negative)



Subgroup Shift (Left)

该模型系统地对来自子组的评价打分较低。

Metric: BNSP AUC

(Background Negative Subgroup Positive)

Results

Toxicity @1

Identity groups	Subgroup AUC	BPSN AUC	BPSP AUC
lesbian	0.93	0.74	0.98
gay	0.94	0.65	0.99
queer	0.98	0.96	0.93
straight	0.99	1.00	0.87
bisexual	0.96	0.95	0.92
homosexual	0.87	0.53	0.99
heterosexual	0.96	0.94	0.92
cis	0.99	1.00	0.87
trans	0.97	0.96	0.91
nonbinary	0.99	0.99	0.99
black	0.91	0.85	0.95
white	0.91	0.88	0.94



Toxicity @6

Identity groups	Subgroup AUC	BPSN AUC	BPSP AUC
lesbian	1.00	0.98	1.00
gay	1.00	0.94	1.00
queer	0.99	0.98	0.99
straight	1.00	1.00	0.97
bisexual	0.98	0.98	0.99
homosexual	1.00	0.96	1.00
heterosexual	1.00	0.99	1.00
cis	1.00	1.00	0.98
trans	1.00	1.00	1.00
nonbinary	1.00	1.00	0.98
black	0.98	0.97	1.00
white	0.99	0.99	0.99



Release Responsibly

Model Cards for Model Reporting

目前还没有模型发布时报告模型效果的common practice

- What It Does
 - 一份关注模型性能透明度的报告，以鼓励负责任的人工智能的采用和应用。
- How It Works
 - 这是一个容易发现的和可用的工件在用户旅程中重要的步骤为一组不同的用户和公共利益相关者。
- Why It Matters
 - 它使模型开发人员有责任发布高质量和公平的模型。

Intended Use, Factors and Subgroups

Example Model Card - Toxicity in Text	
Model Details	Developed by Jigsaw in 2017 as a convolutional neural network trained to predict the likelihood that a comment will be perceived as toxic.
Intended Use	Supporting human moderation, providing feedback to comment authors, and allowing comment viewers to control their experience.
Factors	Identity terms referencing frequently attacked groups focusing on the categories of sexual orientation, gender identity and race.

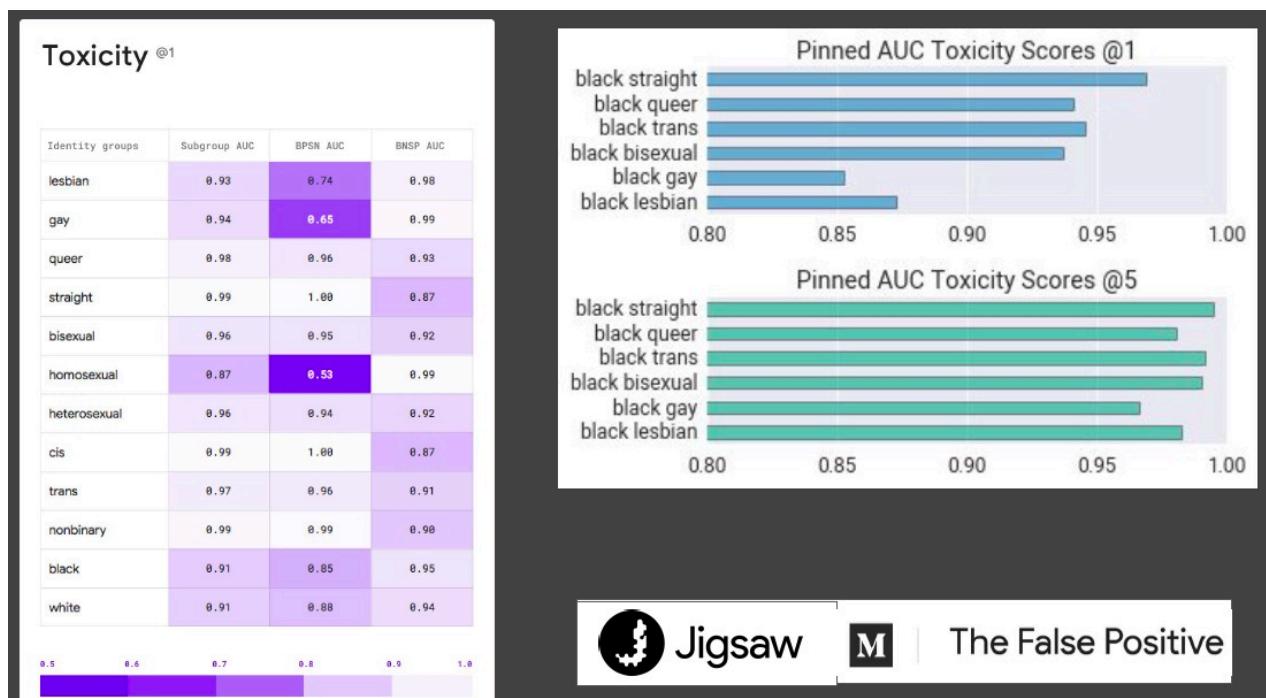
Metrics and Data

Metrics	Pinned AUC, which measures threshold-agnostic separability of toxic and non-toxic comments for each group, within the context of a background distribution of other groups.
Evaluation Data	A synthetic test set generated using a template-based approach, where identity terms are swapped into a variety of template sentences.
Training Data	Includes comments from a variety of online forums with crowdsourced labels of whether the comment is “toxic”. “Toxic” is defined as, “a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion”.

Considerations, Recommendations

Ethical Considerations	A set of values around community, transparency, inclusivity, privacy and topic-neutrality to guide their work.
Caveats & Recommendations	Synthetic test data covers only a small set of very specific comments. While these are designed to be representative of common use cases and concerns, it is not comprehensive.

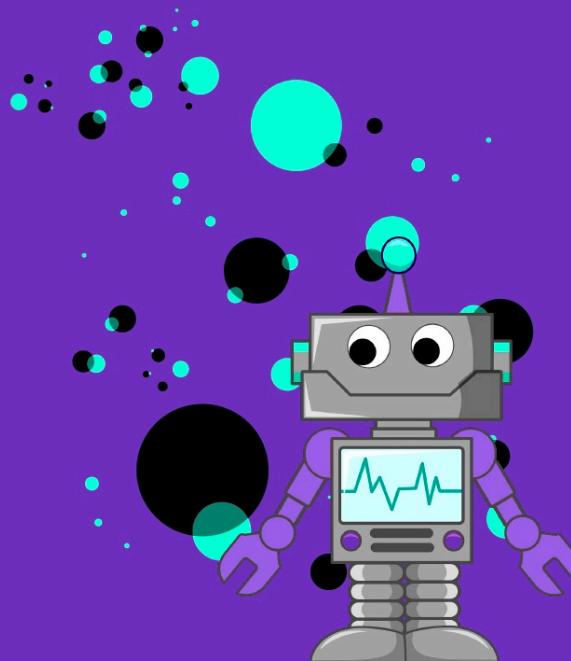
Disaggregated Intersectional Evaluation



Moving from majority representation...

...to diverse representation

...for ethical AI



Reference

以下是学习本课程时的可用参考书籍：

[《基于深度学习的自然语言处理》](#) (车万翔老师等翻译)

[《神经网络与深度学习》](#)

以下是整理笔记的过程中参考的博客：

[斯坦福CS224N深度学习自然语言处理2019冬学习笔记目录](#) (课件核心内容的提炼，并包含作者的见解与建议)

[斯坦福大学 CS224n自然语言处理与深度学习笔记汇总](#) {>>这是针对note部分的翻译<<}