Fairness in visual recognition

Olga Russakovsky

http://visualai.princeton.edu





Computer vision model learns to "increase attractiveness" by manipulating skin color

THE VERGE April 25, 2017

"Machines taught by photos learn a sexist view of women"

WIRED Aug 21, 2017

"Facial recognition is accurate, if you're a white guy"

The New York Times Feb 9, 2018

See: Ruha Benjamin's "Race after technology" — an excellent book



See: Ruha Benjamin's "Race after technology" – an excellent book



Large scale \neq fair representation

Race diversity in face datasets



[Joy Buolamwini & Timnit Gebru. FAT* 2018]

Geographic diversity

(in ImageNet and OpenImages)



[Shreya Shankar et al. NeurIPS 2017]

Object diversity



[Terrance DeVries et al. CVPRW 2019]

Diversity in image search results



[Matthew Kay et al. CHI 2015]

Safiya Umoja Noble's "Algorithms of Oppression" — another excellent book

REVISE: REvealing VIsual biaSEs tool



Fig. 1: Our tool takes in as input a visual dataset and its annotations, and outputs metrics, seeking to produce insights and possible actions.

Key contributions:

- Provides transparency into *large-scale* datasets
- Aids dataset creators&users: fairness ultimately requires manual intervention
- Leverages the (1) available annotations, (2) existing pre-trained models, and
 (3) models trained on the data
- Goes beyond noting *underrepresentation* to analyzing differences in *portrayal* of different objects, genders and geographic regions

Context in images depicting males and females*

*Gender annotations of *binarized sociallyperceived* gender *expression* derived from image captions [J. Zhao et al. EMNLP'17]

Analysis: correlate the presence of different genders in COCO [Lin et al. ECCV'14] with

Objects, using ground truth object annotations **Scenes**, computed with the pre-trained grouped manually into super-categories Places network [B. Zhou et al. TPAMI '17] mountains, desert, sky sports water, ice, snow vehicle male industrial and construction female outdoor outdoor transportation outdoor sports fields, parks animal outdoor man-made elements Object Category accessory cultural or historical place Scene indoor transportation electronic forest, field, jungle furniture cabins, gardens, farms commercial buildings, towns food indoor sports and leisure appliance workplace indoor cultural indoor male home or hotel female kitchen shopping and dining 0.1 0.2 0.3 0.0 0.4 0.5 0.6 0.00 0.05 0.10 0.15 0.20 Fraction of Images that contain this Category Fraction of Images with this Scene

Actionable insight: collect images of the underrepresented gender with the corresponding objects and scenes

["REVISE: A Tool for Measuring and Mitigating Bias in Visual Datasets."

Angelina Wang, Arvind Narayanan, Olga Russakovsky. ECCV 2020 (spotlight). https://github.com/princetonvisualai/revise-tool]

Interaction between objects and people of different genders*

*Gender annotations of *binarized sociallyperceived* gender *expression* derived from image captions [J. Zhao et al. EMNLP'17]

Analysis: use the person-object distance as a proxy for interaction



Actionable insight: consider equalizing the level of interaction with the object (if warranted)

["REVISE: A Tool for Measuring and Mitigating Bias in Visual Datasets." Angelina Wang, Arvind Narayanan, Olga Russakovsky. ECCV 2020 (spotlight). <u>https://github.com/princetonvisualai/revise-tool]</u>

Differences in *portrayal* of different genders*

*Gender annotations of *binarized sociallyperceived* gender *expression* derived from image captions [J. Zhao et al. EMNLP'17]

Analysis:

- for each object class, learn visual classifiers for recognizing this object when it's present with females vs present with males
- identify classes with most stark differences between genders



Actionable insight: collect more images of each gender with the particular object in more diverse situations











The problem: teaching a classifier to ignore a known spurious correlation in the data

Training: skewed distribution (e.g., the target class is correlated with a protected attribute like gender or race)



Testing: classifying images while ignoring this undesirable correlation

hat

no hat





hat is not correlated with glasses

Baseline may learn this undesirable correlation



 $\mathcal{L} = -\sum_{i} log P(y_i | x_i)$ Target classifier

CNN

• Will learn the most discriminative feature(s) in the training dataset

 x_i = image *i* y_i = target class for image *i*

(1) Over/under-sampling may not work well



 $\mathcal{L} = -\sum_{i} log P(y_i | x_i)$



 $x_i = \text{image } i$ $y_i = \text{target class for image } i$

- Undersampling discards data
- Oversampling/reweighting is prone to overfitting and numerical instability

(2) Adversarial de-biasing does not work



(3) Domain-independent training works very well



$$\mathcal{L} = -\sum_{i} log P(y_i | d_i, x_i)$$



 d_i = protected attribute for image *i*

without glasses





Inference:

 $\arg \max_{y} \sum_{d} s(y, d, x)$

 $s = class \ score$

 Works especially well in settings with high skew (high correlation between target and protected attributes)

These findings hold across a variety of settings

Model	CIFAR-grayscale	CIFAR-ImageNet	CelebA	
Baseline Oversampling Adversarial DomIndep	88.5 ± 0.3 89.1 ± 0.4 83.8 ± 1.1 92.0 ± 0.1	79.4 ± 0.4 78.6 ± 0.4 74.1 ± 0.6 $\mathbf{83.5 \pm 0.3}$	74.7 77.6 71.9 76.3	
	Data: Modified CIFAR-10 [Krizhevsky et al. 2009]	Data: Modified CIFAR-10 + ImageNet	Data: CelebA dataset [Liu et al. ICCV 2015]	
	Target: 10-way object classification	[Deng et al. CVPR 2009] Target: 10-way object classification	Target: multi-label face attribute recognition	
	color/grayscale	Protected attribute: CIFAR vs ImageNet	Protected attribute: gender	
	ResNet-18 Metric: Mean per-		Architecture: ResNet-50, pre- trained	
	accuracy	Metric: Mean per- class per-attribute accuracy	Metric: Weighted mean average precision	

Can we generate the data we need?

Real-world data



Paired augmentation















[E. Denton, B. Hutchinson, M. Mitchell, and T. Gebru. "Image counterfactual sensitivity analysis for detecting unintended bias." In CVPRW 2019.]



[E. Denton, B. Hutchinson, M. Mitchell, and T. Gebru. "Image counterfactual sensitivity analysis for detecting unintended bias." In CVPRW 2019.]



[E. Denton, B. Hutchinson, M. Mitchell, and T. Gebru. "Image counterfactual sensitivity analysis for detecting unintended bias." In CVPRW 2019.]



[E. Denton, B. Hutchinson, M. Mitchell, and T. Gebru. "Image counterfactual sensitivity analysis for detecting unintended bias." In CVPRW 2019.]

Generated:



Original:





Original:





Original:









Original:





Updated:

Original:



Pros:

- 1. Single GAN for multiple target/protected attributes
- 2. Preserves intra-class variation
- 3. Ensures independence between attributes in the generated data
- 4. Can be generalized to multiple protected attributes

Cons:

- 1. The classifiers need to be at least reasonable
- 2. Might add other correlations
- 3. GANs are tricky...



Original:



Evaluating classifiers trained with our generated data

Setup: CelebA dataset, protected attribute = gender, target attributes = 26 other face attributes (e.g., has-hat)

vs baseline: much better on all 3 fairness metrics

Attr. type	Model	$DEO\downarrow$	$BA\downarrow$	$KL\downarrow$
Incons	Baseline	21.5 ± 4.4	2.1 ± 0.6	1.7 ± 0.3
meons.	Ours	$\bf 16.5 \pm 4.2$	0.5 ± 0.6	1.3 ± 0.4
G-dep.	Baseline	25.7 ± 3.5	2.3 ± 0.5	1.3 ± 0.2
	Ours	23.4 ± 3.6	1.6 ± 0.5	1.2 ± 0.2
G-indep.	Baseline	16.7 ± 5.0	0.3 ± 0.6	1.1 ± 0.5
	Ours	13.9 ± 5.2	0.0 ± 0.5	0.9 ± 0.6

vs weighted and adversarial:

generally better than adversarial; not as good as weighted

(but weighted strongly amplifies bias in the negative direction)

Method	DEO ↓	BA↓	$KL\downarrow$	
Weighted	5.7 ± 4.2	$\textbf{-2.8}\pm\textbf{0.5}$	0.5 ± 0.4	
Adversarial	23.9 ± 4.4	1.5 ± 0.5	0.6 ± 0.5	
Ours	16.7 ± 4.7	0.5 ± 0.5	1.0 ± 0.5	

vs domain-independent: better when skew is low, but domain-independent still best when skew is high

Skew	Method	$DEO\downarrow$	$BA\downarrow$	$KL\downarrow$
Low/	Dom. Ind. 7.0 ± 3.1		$\textbf{-0.1}\pm\textbf{0.5}$	0.8 ± 0.7
Mod.	Ours	$\textbf{6.0} \pm \textbf{3.0}$	-0.1 \pm 0.5	$\textbf{0.3} \pm \textbf{0.1}$
High	Dom. Ind.	$\textbf{14.9} \pm \textbf{5.6}$	-0.4 \pm 0.5	$\textbf{0.8} \pm \textbf{1.0}$
	Ours	23.9 ± 5.5	0.9 ± 0.4	1.5 ± 0.6

vs others like Fairness GAN: better

	Fairness GAN [35]			Ours		
	Dem. Par.		Eq. Opp.		(Synthetic only)	
Gender exp. g	g = -1	g=1	g = -1	g=1	g = -1	g=1
$FPR\downarrow$	0.52	0.26	0.42	0.17	0.22	0.39
$FNR\downarrow$	0.18	0.41	0.21	0.44	0.06	0.27
Error \downarrow	0.30	0.28	0.29	0.23	0.21	0.18
Error Rate \downarrow	0.22		0.29		0.20	

[P. Sattigeri, et al. IBM Journal of Research and Development, 2019]

[Vikram V. Ramaswamy, Sunnie S. Y. Kim, Olga Russakovsky.

"Fair Attribute Classification through Latent Space De-biasing." https://github.com/princetonvisualai/gan-debiasing]





AI4ALL: a non-profit dedicated to increasing diversity and inclusion in AI

-



"Until this program, I never thought that people who look like me could succeed in computer science and AI."

- AI4ALL 2016 student

- Our goal is to educate and support a diverse next generation of AI *leaders*
- Partnered with 16 universities to run summer programs for high school students from underrepresented groups
 - https://ai-4-all.org/summer-programs/
 - Launched a free online OpenLearning platform
 - https://ai-4-all.org/open-learning
- Providing life-long training, mentorship, opportunities and community for our alumni
 - Amazingly, many alumni have taken on Al leadership roles while still in high school/early college by starting research projects, working groups, panels, clubs, outreach events, ...
 - https://medium.com/ai4allorg/alumni/





http://visualai.princeton.edu

olgarus@cs.princeton.edu



accuracy and fairness

[Vikram V. Ramaswamy, Sunnie S. Y. Kim and Olga Russakovsky. https://github.com/princetonvisualai/gan-debiasing]

and Arvind Narayanan's "Fairness and machine learning"