Fairness in visual recognition

Olga Russakovsky



Vikram Ramaswamy



Angelina Wang



Zeyu Wang













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

1HE VERGE

April 25, 2017







Large scale \neq fair representation

Race diversity in face datasets

Dataset

IJB-A Adience



[Joy Buolamwini & Timnit Gebru. FAT* 2018]

Geographic diversity (in ImageNet and OpenImages)

[Shreya Shankar et al. NeurIPS 2017 Workshop]

Diversity in image search results



[Matthew Kay et al. CHI 2015]

Stereotyped representation in datasets

person+flower



[Angelina Wang et al. ECCV 2020]

Counteracting the disparities by annotating demographics

Annotated demographics on 139 people synsets (categories) in ImageNet 13,900 images; 109,545 worker judgments.





["Towards Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in the ImageNet Hierarchy." Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, Olga Russakovsky. FAT* 2020. <u>http://image-net.org/filtering-and-balancing</u>]

Counteracting the disparities by annotating demographics

Annotated demographics on 139 people synsets (categories) in ImageNet 13,900 images; 109,545 worker judgments.



Subtleties:

- Rebalancing \implies removing data, changing the original distribution
- Accuracy/validity of these labels
- The implication of including people categories in a dataset (cf. the FAT* paper)
- Privacy of subjects, esp. minors; consent of content creators (working on this)
- The representation of folks of different genders (skin colors, ages) within a synset

["Towards Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in the ImageNet Hierarchy." Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, Olga Russakovsky. FAT* 2020. <u>http://image-net.org/filtering-and-balancing</u>]

Revealing and mitigating dataset biases with 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:

- Goes beyond *underrepresentation* to analyzing differences in *portrayal*
- Allows for semi-automatic analysis of *large-scale* datasets
- Aids dataset creators&users: fairness ultimately requires manual intervention
- Integrates bias mitigation throughout the dataset *construction* process

Inner workings of the REVISE tool

Implementation:

- Freely available Python notebooks
- Analyzes portrayal of objects, people and geographic regions
- Uses provided annotations, pre-trained models, and models trained on the data

In this talk:

- Focus specifically on portrayal of different *genders*
- Caveat: use of binarized socially-perceived gender expression
- Analysis on COCO [T. Y. Lin et al. ECCV '14] and OpenImages [I. Krasin et al. '17]
- Gender annotations derived from image captions [J. Zhao et al. EMNLP'17]

Co-occurrence of males and females with different objects and in different scenes

Analysis: correlate the presence of different genders in COCO with



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

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

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

Annotated gender in datasets defaults to "male"

Analysis: investigate occurrences where gender is annotated but the person is too small or no face is detected in the image



Man and boats on the sand in low tide.



The group of buses are parked along the city street as a man crosses the street in the background.



A man is kiteboarding in the open ocean.



A man riding a kiteboard on top of a wave in the ocean.

Actionable insight: prune these gender labels

["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>]









Re-visiting many existing problems in this context

Long tail distributions



Learning with constraints



Domain adaptation



Interpretability



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

Toy illustration on CIFAR, to temporarily simplify the exploration

Training: skewed distributions (correlates class with color/grayscale)



Testing: classifying images into one of 10 object classes (no correlation)



Testing on color images

Training on skewed data: 89% accuracy

Training on all-grayscale: 93% accuracy

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

Toy illustration on CIFAR, to temporarily simplify the exploration

Training: skewed distributions (correlates class with color/grayscale)



Testing: classifying images into one of 10 object classes (no correlation)



Classes primarily in color during training airpane 0.8 bird 0.4 deet frog 0.2 **Frue** label ship automobile cat dog horse truck Predicted label

["Towards Fairness in Visual Recognition: Effective Strategies for Bias Mitigation." Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, Olga Russakovsky. CVPR 2020. <u>https://github.com/princetonvisualai/DomainBiasMitigation]</u>

Testing on color images

Domain-independent training works very well





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



 $\begin{aligned} x_i &= \text{image } i \\ y_i &= \text{object class for image } i \\ d_i &= \text{domain (c or g) for image } i \end{aligned}$

Inference:

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

s = pre-softmax score

Domain-independent training works very well, especially with high class/domain correlation

Training data: CIFAR-10, skewed color/grayscale distribution

Architecture:

ResNet-18

Testing metric: Mean per-class per-domain accuracy (i.e., equal color/grayscale distribution within classes)



Adversarial de-biasing does not work



Adversarial de-biasing does not work



- It is very difficult to train a visual representation that can't classify the domain (visual representations are powerful!)
- In cases of high correlation between classes and domains, one can be inferred from the other during training so adversarial de-biasing may not be appropriate
- Beware of fairness literature that reports equal error rates but not overall accuracy

These findings hold across a variety of settings

Model	CIFAR-grayscale	CIFAR-ImageNet	CelebA
BASELINE Adversarial DomIndep	$\begin{array}{c} 88.5 \pm 0.3 \\ 83.8 \pm 1.1 \\ \textbf{92.0} \pm \textbf{0.1} \end{array}$	79.4 ± 0.4 74.1 ± 0.6 $\mathbf{83.5 \pm 0.3}$	74.7 71.9 76.3
	Training data: CIFAR-10, skewed color/grayscale distribution	Same, except more subtle domain shift (substituting in images of similar classes from ImageNet instead of converting to grayscale)	Task: multi-label face attribute recognition, wher presence&appearance of an attribute may be
	Architecture: ResNet-18		correlated with gender Architecture: ResNet-50,
	per-class per- domain accuracy (i.e., equal color/ grayscale distribution within classes)		Testing metric : Weighted mean average precision (i.e., equal gender distribution within classes)

Coming soon: can we generate the data we need?





Fig. 2: Consider a GAN trained on a biased dataset of faces where the presence of hats is correlated with the presence of sunglasses. Moving in the latent space towards adding glasses additionally adds a hat (top). In contrast, our method allows for adding glasses without including the correlated hat attribute (bottom). Note that other attributes (apart from the target attribute) can change







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

- Celebrated our 3rd birthday on March 8, 2020
- 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
- Summer program alumni have started Al research projects, internships, working groups, panels, clubs, ... (while still in high school/early college)
 - https://medium.com/ai4allorg/alumni/
- Long-term vision is to foster a community of diverse *leaders* in Al



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

- AI4ALL 2016 student



