## Machine Learning 2

DS 4420 - Spring 2020

### Bias and fairness

Byron C. Wallace

Material in this lecture modified from materials created by Jay Alammar (<a href="http://jalammar.github.io/illustrated-transformer/">http://jalammar.github.io/illustrated-transformer/</a>) and Sasha Rush (<a href="https://nlp.seas.harvard.edu/2018/04/03/attention.html">https://nlp.seas.harvard.edu/2018/04/03/attention.html</a>).



## Intro

### Today

 We will talk about bias and fairness, which are critically important to understand if you go out and apply models in real-world settings

## Examples [from CIML, Daume III]

• Early speech recognition systems failed on female voices.

## Examples [from CIML, Daume III]

- Early speech recognition systems failed on female voices.
- Models to predict criminal recidivism biased against minorities.



IAN WALDIE/GETTY IMAGES

#### **Tech Policy / AI Ethics**

## Al is sending people to jail —and getting it wrong

Using historical data to train risk assessment tools could mean that machines are copying the mistakes of the past.

by Karen Hao

Jan 21, 2019





Gender was misidentified in **up to 1 percent of lighter-skinned males** in a set of 385 photos.



Gender was misidentified in **up to 7 percent of lighter-skinned females** in a set of 296 photos.



Gender was misidentified in **up to 12 percent of darker-skinned males** in a set of 318 photos.



Gender was misidentified in **35 percent of darker-skinned females** in a set of 271 photos.

### Can word vectors be sexist?

## Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

 $\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}$ 

#### Gender stereotype she-he analogies.

sewing-carpentry	register-nurse-physician	housewife-shopkeeper
nurse-surgeon	interior designer-architect	softball-baseball
blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
giggle-chuckle	vocalist-guitarist	petite-lanky
sassy-snappy	diva-superstar	charming-affable
volleyball-football	cupcakes-pizzas	hairdresser-barber

#### Gender appropriate she-he analogies.

queen-king	sister-brother	mother-father
waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

#### Extreme *she* occupations

1. homemaker	2. nurse	3. receptionist
4. librarian	5. socialite	6. hairdresser

7. nanny 8. bookkeeper 9. stylist

10. housekeeper 11. interior designer 12. guidance counselor

#### Extreme *he* occupations

maestro
 skipper
 philosopher
 captain
 architect
 financier
 warrior
 broadcaster

10. magician 11. figher pilot 12. boss

word2vec resume

Q

#### Scholar

About 93 results (0.02 sec)

#### Articles

#### Machine Learned **Resume**-Job Matching Solution

Case law

My library

Y Lin, H Lei, PC Addo, X Li - arXiv preprint arXiv:1607.07657, 2016 - arxiv.org

... We use LDA to classify **resumes** into 32 and 64 topics respectively. ... each Chinese phrase as a word and each list of phrases as a sentence, after **word2vec** training, each ... In this paper, we have considered the **resume**-job matching problem and pro- posed a solution by using ... Cite Save

#### Any time

Since 2016

Since 2015

Since 2012

Custom range...

#### [PDF] SKILL: A System for Skill Identification and Normalization.

M Zhao, F Javed, F Jacob, M McNair - AAAI, 2015 - pdfs.semanticscholar.org

... ThiS dictionary capacitateS 90% of noiSe exhibited in reSume SkillS SectionS. ... iS initiated firSt for the input queY ry (aka, Seed Skill phraSeS from reSumeS) for proper ... implement and produce highly precise and relevant skills recognition system, we utilize word2vec (Mikolov et ... Cited by 4 Related articles All 3 versions Cite Save More

#### Sort by relevance

Sort by date

✓ include patents✓ include citations

Create alert

Word2Vec vs DBnary ou comment (ré) concilier représentations distribuées et réseaux lexico-sémantiques? Le cas de l'évaluation en traduction automatique C Servan, Z Elloumi, H Blanchon, L Besacier - TALN 2016, 2016 - hal.archives-ouvertes.fr

... Page 2. Word2Vec vs DBnary ou comment (ré)concilier représentations ... RÉSUMÉ Cet article présente une approche associant réseaux lexico-sémantiques et représentations distribuées de mots appliquée à l'évaluation de la traduction automatique. ...

Cite Save

#### Macau: Large-scale skill sense disambiguation in the online recruitment domain

Q Luo, M Zhao, F Javed, F Jacob - Big Data (Big Data), 2015 ..., 2015 - ieeexplore.ieee.org

... Contexts are extracted from either skill section(s) of resumes or requirement section(s) of job postings. We used a popular tool word2vec [12] with parameter

Bolukbasi et al. '16 Slides: Adam Kalai



Figure from: https://towardsdatascience.com/named-entity-recognition-with-nltk-and-spacy-8c4a7d88e7da

## Recognizing names in text

Country	,	(Huang et al., 2015) GloVe words		(Lample et al., 2016) GloVe words+chars		(Devlin et al., 2019) BERT subwords			
	P	R	F1	P	R	F1	P	R	F1
Original	96.9	96.5	96.7	97.1	98.1	97.6	98.3	98.1	98.2
US	96.9	99.6	98.2	96.9	99.6	98.3	98.4	99.7	99.1
Russia	96.8	99.5	98.1	97.1	99.8	98.4	98.4	99.3	98.9
India	96.5	99.5	98.0	97.1	99.3	98.2	98.4	98.8	98.6
Mexico	96.7	98.9	97.8	97.1	98.9	98.0	98.4	99.2	98.8
China-Taiwan	95.4	93.2	93.9	97.0	94.9	95.6	98.3	92.0	94.8
US (Difficult)	95.9	87.4	90.2	96.6	87.9	90.7	98.1	88.5	92.3
Indonesia	95.3	84.6	88.7	96.5	91.0	93.3	97.8	85.8	92.0
Vietnam	94.6	78.2	84.2	96.0	78.5	84.5	98.0	84.2	89.8

### Intermezzo 1

Before moving on to the next part of lecture, let's walk through this notebook tutorial

https://nbviewer.jupyter.org/github/Azure-Samples/learnAnalytics-DeepLearning-Azure/blob/master/Students/12-biased-embeddings/how-to-make-a-racist-ai-without-really-trying.ipynb

## Domain adaptation

# One potential cause: Train/test mismatch

• If the train set is drawn from a different distribution than the test set, this introduces a *bias* such that the model will do better on examples that look like train set instances

# One potential cause: Train/test mismatch

- If the train set is drawn from a different distribution than the test set, this introduces a bias such that the model will do better on examples that look like train set instances
- If the speech recognition model has been trained on mostly male voices and optimized well, it will tend to do better on male voices.

### Unsupervised adaptation

• Given training data from distribution D<sup>old</sup>, learn a classifier that performs well on a related, but distinct, distribution D<sup>new</sup>

### Unsupervised adaptation

- Given training data from distribution Dold, learn a classifier that performs well on a related, but distinct, distribution Dnew
- Assumption is that we have train data from D<sup>old</sup> but what we actually care about is loss on D<sup>new</sup>

## Unsupervised adaptation

- Given training data from distribution D<sup>old</sup>, learn a classifier that performs well on a related, but distinct, distribution D<sup>new</sup>
- Assumption is that we have train data from D<sup>old</sup> but what we actually care about is loss on D<sup>new</sup>
- What can we do here?

Test loss = 
$$\mathbb{E}_{(x,y)\sim\mathcal{D}^{\text{new}}}\left[\ell(y,f(x))\right]$$

definition (8.2)

Test loss = 
$$\mathbb{E}_{(x,y)\sim\mathcal{D}^{\mathrm{new}}}\left[\ell(y,f(x))\right]$$
 definition (8.2)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{new}}(x,y)\ell(y,f(x))$  expand expectation (8.3)

Test loss = 
$$\mathbb{E}_{(x,y)\sim\mathcal{D}^{\mathrm{new}}}\left[\ell(y,f(x))\right]$$
 definition (8.2)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{new}}(x,y)\ell(y,f(x))$  expand expectation (8.3)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{new}}(x,y) \frac{\mathcal{D}^{\mathrm{old}}(x,y)}{\mathcal{D}^{\mathrm{old}}(x,y)}\ell(y,f(x))$  times one (8.4)

Note: Does this look familiar?!

Test loss = 
$$\mathbb{E}_{(x,y) \sim \mathcal{D}^{\text{new}}} [\ell(y, f(x))]$$
 definition (8.2)  
=  $\sum_{(x,y)} \mathcal{D}^{\text{new}}(x,y) \ell(y,f(x))$  expand expectation (8.3)  
=  $\sum_{(x,y)} \mathcal{D}^{\text{new}}(x,y) \frac{\mathcal{D}^{\text{old}}(x,y)}{\mathcal{D}^{\text{old}}(x,y)} \ell(y,f(x))$  times one (8.4)  
=  $\sum_{(x,y)} \mathcal{D}^{\text{old}}(x,y) \frac{\mathcal{D}^{\text{new}}(x,y)}{\mathcal{D}^{\text{old}}(x,y)} \ell(y,f(x))$  rearrange (8.5)

Test loss = 
$$\mathbb{E}_{(x,y)\sim\mathcal{D}^{\mathrm{new}}}[\ell(y,f(x))]$$
 definition (8.2)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{new}}(x,y)\ell(y,f(x))$  expand expectation (8.3)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{new}}(x,y) \frac{\mathcal{D}^{\mathrm{old}}(x,y)}{\mathcal{D}^{\mathrm{old}}(x,y)}\ell(y,f(x))$  times one (8.4)  
=  $\sum_{(x,y)} \mathcal{D}^{\mathrm{old}}(x,y) \frac{\mathcal{D}^{\mathrm{new}}(x,y)}{\mathcal{D}^{\mathrm{old}}(x,y)}\ell(y,f(x))$  rearrange (8.5)  
=  $\mathbb{E}_{(x,y)\sim\mathcal{D}^{\mathrm{old}}}\left[\frac{\mathcal{D}^{\mathrm{new}}(x,y)}{\mathcal{D}^{\mathrm{old}}(x,y)}\ell(y,f(x))\right]$  definition (8.6)

## Importance weighting

- So we have re-expressed the test loss as an expectation over Dold, which is good because that's what we have for training data
- But we do not have access to D<sup>old</sup> or D<sup>new</sup> directly

### Ratio estimation

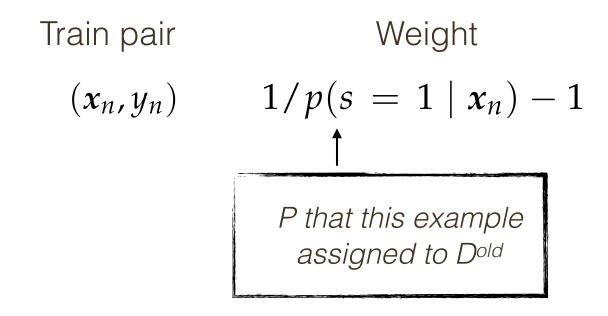
Assume all examples drawn from an underlying shared distribution (base), and then sorted into  $D^{old}/D^{new}$  with some probability depending on x

$$\mathcal{D}^{\text{old}}(x, y) \propto \mathcal{D}^{\text{base}}(x, y) p(s = 1 \mid x)$$

$$\mathcal{D}^{\text{new}}(x, y) \propto \mathcal{D}^{\text{base}}(x, y) p(s = 0 \mid x)$$

### Ratio estimation

Supposing we can estimate p... we can reweight examples:



Intuitively: Upweights instances likely to be from *Dnew* 

## How should we estimate *p*?

Want to estimate:  $p(s = 1 \mid x_n)$ 

This is just a binary classification task!

### Algorithm 23 SELECTION ADAPTATION $(\langle (x_n, y_n) \rangle_{n=1}^N, \langle z_m \rangle_{m=1}^M, \mathcal{A})$

1: 
$$D^{dist} \leftarrow \langle (\boldsymbol{x}_n, +1) \rangle_{n=1}^N \cup \langle (\boldsymbol{z}_m, -1) \rangle_{m=1}^M$$

// assemble data for distinguishing

// between old and new distributions

 $\hat{p} \leftarrow \text{train logistic regression on } D^{\text{dist}}$ 

3: 
$$D^{weighted} \leftarrow \left\langle \left( x_n, y_n, \frac{1}{\hat{p}(x_n)} - 1 \right) \right\rangle_{n=1}^N$$

4: **return**  $A(D^{weighted})$ 

// assemble weight classification // data using selector // train classifier

#### Supervised adaptation

 We were supposing that we had access to labels only in Dold, but wanted to learn a model for Dnew

#### Supervised adaptation

- We were supposing that we had access to *labels* only in D<sup>old</sup>, but wanted to learn a model for D<sup>new</sup>
- In some cases we might have at least some labels from D<sup>new</sup> as well

shared old-only new-only

shared old-only new-only 
$$x_n^{(\text{old})} \mapsto \left\langle \begin{array}{c} x_n^{(\text{old})} & , \ x_n^{(\text{old})} \\ \end{array} \right\rangle \xrightarrow[D-\text{many}]{}$$

shared old-only new-only 
$$x_n^{(\text{old})} \mapsto \left\langle \begin{array}{c} x_n^{(\text{old})} & , \ x_n^{(\text{old})} & , \ x_n^{(\text{old})} & , \ \underbrace{0,0,\ldots,0}_{D\text{-many}} \end{array} \right\rangle$$
  $x_m^{(\text{new})} \mapsto \left\langle \begin{array}{c} x_m^{(\text{new})} & , \ \underbrace{0,0,\ldots,0}_{D\text{-many}} & , \ x_m^{(\text{new})} \end{array} \right\rangle$ 

shared old-only new-only 
$$x_n^{(\text{old})} \mapsto \left\langle \begin{array}{c} x_n^{(\text{old})} & , \ x_n^{(\text{old})} & , \ x_n^{(\text{old})} & , \ \underbrace{0,0,\ldots,0}_{D\text{-many}} \end{array} \right\rangle$$
  $x_m^{(\text{new})} \mapsto \left\langle \begin{array}{c} x_m^{(\text{new})} & , \ \underbrace{0,0,\ldots,0}_{D\text{-many}} & , \ x_m^{(\text{new})} \end{array} \right\rangle$ 

We have seen this trick before!!

Algorithm 24 EasyAdapt( $\langle (x_n^{(old)}, y_n^{(old)}) \rangle_{n=1}^N$ ,  $\langle (x_m^{(new)}, y_m^{(new)}) \rangle_{m=1}^M$ ,  $\mathcal{A}$ )

1:  $D \leftarrow \left\langle (\langle x_n^{(old)}, x_n^{(old)}, \mathbf{0} \rangle, y_n^{(old)}) \right\rangle_{n=1}^N \cup \left\langle (\langle x_m^{(new)}, \mathbf{0}, x_m^{(new)} \rangle, y_m^{(new)}) \right\rangle_{m=1}^M$  // union // of transformed data

2: **return**  $\mathcal{A}(D)$  // train classifier

#### Subtler bias

- What if the distribution is not different, but rather the train data simply reflects biases in society?
- Training models on this and then using them can magnify biases that already exist!

#### Sensitive attributes

In many settings, there may be certain fields/attributes that we know a
priori we don't want to exploit. Great: Let's just remove these!

#### Sensitive attributes

- In many settings, there may be certain fields/attributes that we know a
  priori we don't want to exploit. Great: Let's just remove these!
- Why is this not enough?

#### Sensitive attributes

- In many settings, there may be certain fields/attributes that we know a
  priori we don't want to exploit. Great: Let's just remove these!
- Why is this not enough?

Because other features may correlate strongly with the protected feature!

Let's consider a concrete example: Hiring

Following slides derived from: <a href="https://mrtz.org/nips17">https://mrtz.org/nips17</a>

Solon Barocas and Moritz Hardt









**ROBO RECRUITING** 

#### Can an Algorithm Hire Better Than a Human?



Claire Cain Miller @clairecm JUNE 25, 2015











Hiring and recruiting might seem like some of the least likely jobs to be automated. The whole process seems to need human skills that computers lack, like making conversation and reading social cues.

But people have biases and predilections. They make hiring decisions, often unconsciously, based on similarities that have nothing to do with the job requirements — like whether an applicant has a friend in common, went to the same school or likes the same sports.

That is one reason researchers say traditional job searches are broken. The question is how to make them better.

A new wave of start-ups — including <u>Gild</u>, <u>Entelo</u>, <u>Textio</u>, <u>Doxa</u> and <u>GapJumpers</u> — is trying various ways to automate hiring. They say that software can do the job more effectively and efficiently than people can. Many people are beginning to buy into the idea. Established headhunting firms like Korn Ferry are incorporating algorithms into their work, too.

If they succeed, they say, hiring could become faster and less expensive, and their data could lead recruiters to more highly skilled people who are better matches for their companies. Another potential result: a more diverse workplace. The software relies on data to surface candidates from a wide variety of places and match their skills to the job requirements, free of human biases.



Mayurakshi Ghosh January 7, 2016

Hi Claire, excellent article and really insightful facts on algorithm recruitment. I completely agree how you mentioned the role played by...

#### Deborah Bishop July 2, 2015

This is a very interesting article. Perhaps distinguishing between different aspects in the process of bringing talent into your...

#### Yeti July 2, 2015

Including talents that would be rejected by the subjective biased boss or colleagues does not fuarantee their integration. They will be...

SEE ALL COMMENTS WRITE A COMMENT

"[H]iring could become faster and less expensive, and [...] lead recruiters to more highly skilled people who are better matches for their companies. Another potential result: a more diverse workplace. The software relies on data to surface candidates from a wide variety of places and match their skills to the job requirements, free of human biases."

Miller (2015)





HIDDEN BIAS

#### When Algorithms Discriminate



Claire Cain Miller @clairecm JULY 9, 2015













The online world is shaped by forces beyond our control, determining the stories we read on Facebook, the people we meet on OkCupid and the search results we see on Google. Big data is used to make decisions about health care, employment, housing, education and policing.

But can computer programs be discriminatory?

There is a widespread belief that software and algorithms that rely on data are objective. But software is not free of human influence. Algorithms are written and maintained by people, and machine learning algorithms adjust what they do based on people's behavior. As a result, say researchers in computer science, ethics and law, algorithms can reinforce human prejudices.

Google's online advertising system, for instance, showed an ad for high-income jobs to men much more often than it showed the ad to women,  $\underline{a}$  new study by Carnegie Mellon University researchers found.

Research from Harvard University found that ads for arrest records were significantly more likely to show up on searches for distinctively black names or a historically black fraternity. The Federal Trade Commission said advertisers are able to target people who live in low-income neighborhoods with high-interest loans.

#### RECENT COMMENTS

tom July 10, 2015

Discrimination against women persists in other ways. Take the obituary column of the NYT - on a good week, you will find obits for perhaps...

SierramanCA July 10, 2015

"There is a widespread belief that software and algorithms that rely on data are objective." says Ms. Miller. Well, Ms. Miller, two things:1...

Dalgliesh July 10, 2015

Algorithms are written by people. People are biased, not objective. Daniel Kahneman et al. have proven this.

SEE ALL COMMENTS

# Isn't discrimination the very point of machine learning?

### Discrimination is not a general concept

It is **domain** specific

Concerned with important opportunities that affect people's life chances

It is **feature** specific

Concerned with socially salient qualities that have served as the basis for unjustified and systematically adverse treatment in the past

## Regulated domains

- Credit (Equal Credit Opportunity Act)
- Education (Civil Rights Act of 1964; Education Amendments of 1972)
- Employment (Civil Rights Act of 1964)
- Housing (Fair Housing Act)
- 'Public Accommodation' (Civil Rights Act of 1964)

#### Legally recognized 'protected classes'

Race (Civil Rights Act of 1964); Color (Civil Rights Act of 1964); Sex (Equal Pay Act of 1963; Civil Rights Act of 1964); Religion (Civil Rights Act of 1964);
 National origin (Civil Rights Act of 1964); Citizenship (Immigration Reform and Control Act); Age (Age Discrimination in Employment Act of 1967);
 Pregnancy (Pregnancy Discrimination Act); Familial status (Civil Rights Act of 1968); Disability status (Rehabilitation Act of 1973; Americans with Disabilities Act of 1990); Veteran status (Vietnam Era Veterans'
 Readjustment Assistance Act of 1974; Uniformed Services Employment and Reemployment Rights Act); Genetic Information (Genetic Information Nondiscrimination Act)

#### Discrimination Law: Two Doctrines

**Disparate Treatment** 

**Disparate Impact** 

**Formal** 

Unjustified

or

or

**Intentional** 

**Avoidable** 

## Disparate Treatment

Formal: explicitly considering class membership

Even if it is relevant

Intentional: purposefully attempting to discriminate without direct reference to class membership

Pretext or 'motivating factor'

#### Disparate Impact

1. Plaintiff must first establish that decision procedure has a disparate impact

'Four-fifths rule'

2. Defendant must provide a justification for making decisions in this way

'Business necessity' and 'job-related'

3. Finally, plaintiff has opportunity to show that defendant could achieve same goal using a different procedure that would result in a smaller disparity

'Alternative practice'

# What does discrimination law aim to achieve?

**Disparate Treatment** 

**Disparate Impact** 

**Procedural fairness** 

Distributive justice

**Equality of opportunity** 

Minimized inequality of outcome

#### How machines learn to discriminate

Skewed sample
Tainted examples
Limited features
Sample size disparity
Proxies

B, Selbst (2016)

## Skewed sample

Police records measure "some complex interaction between criminality, policing strategy, and community-policing relations"

Lum, Isaac (2016)

## Skewed sample: feedback loop

Future observations of crime confirm predictions

Fewer opportunities to observe crime that contradicts predictions

Initial bias may compound over time

#### Limited features

Features may be less informative or less reliably collected for certain parts of the population

A feature set that supports accurate predictions for the majority group may not for a minority group

Different models with the same reported accuracy can have a very different distribution of error across population



Running example: Hiring ad for (fictitious?) Al startup

## Formal setup

- $\circ$  X features of an individual (browsing history etc.)
- A sensitive attribute (here, gender)
- $\circ$  C = c(X, A) predictor (here, show ad or not)
- Y target variable (here, SWE)

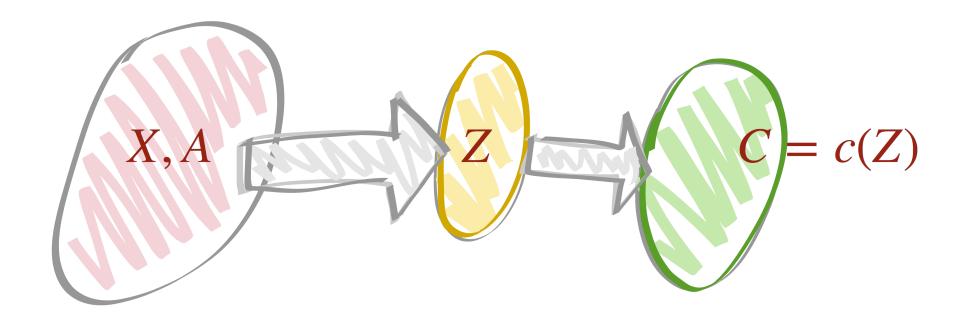
## Formal setup

Score function is any random variable  $R = r(X, A) \in [0, 1]$ .

Can be turned into (binary) predictor by thresholding

Example: Bayes optimal score given by  $r(x, a) = \mathbb{E}[Y \mid X = x, A = a]$ 

#### Representation learning approach



#### Adversarial Learning for Fairness

Beutel et al., 2017; Edwards and Storkey, 2015

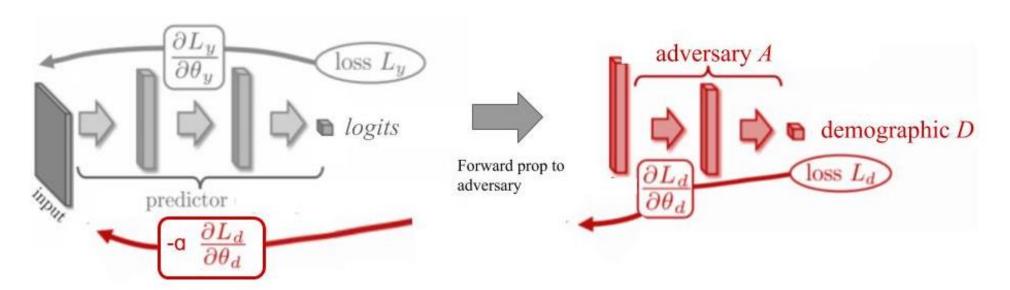


Figure from Wadsworth et al., 2018

#### Three fundamental criteria

**Independence**: *C* independent of *A* 

**Separation**: *C* independent of *A* conditional on *Y* 

**Sufficiency**: *Y* independent of *A* conditional on *C* 

Lots of other criteria are related to these

## First criterion: Independence

Require C and A to be independent, denoted  $C \perp A$ 

That is, for all groups a, b and all values c:

$$\mathbb{P}_a\{C=c\} = \mathbb{P}_b\{C=c\}$$

## Variants of independence

Sometimes called demographic parity, statistical parity

When 
$$C$$
 is binary  $0/1$ -variables, this means  $\mathbb{P}_a\{C=1\} = \mathbb{P}_b\{C=1\}$  for all groups  $a, b$ .

Approximate versions:

$$\frac{\mathbb{P}_a\{C=1\}}{\mathbb{P}_b\{C=1\}} \ge 1 - \epsilon \qquad \qquad |\mathbb{P}_a\{C=1\} - \mathbb{P}_b\{C=1\}| \le \epsilon$$

□ Drawbacks to independence as our criterion? Is this the right objective?

- □ Drawbacks to independence as our criterion? Is this the right objective?
  - May rule out best possible model due to actual correlations in the real-world. Could rule out C = Y (perfect predictor).

- Drawbacks to independence as our criterion? Is this the right objective?
  - May rule out best possible model due to actual correlations in the real-world. Could rule out C = Y (perfect predictor).
  - Can satisfy by just selecting random people form the minority group as "positive" – will not dramatically lower error rate (since they are a minority) but will satisfy constraint.

- Drawbacks to independence as our criterion? Is this the right objective?
  - May rule out best possible model due to actual correlations in the real-world. Could rule out C = Y (perfect predictor).
  - Can satisfy by just selecting random people form the minority group as "positive" – will not dramatically lower error rate (since they are a minority) but will satisfy constraint.
- Other criteria exist let's look at one more: separation

#### Separation

Require R and A to be independent conditional on target variable Y, denoted  $R \perp A \mid Y$ 

**Definition.** Random variable R separated from A if  $R \perp A \mid Y$ .



#### Post-Processing

Method from H, Price, Srebro (2016): Post-processing correct of score function

Post-processing: Any thresholding of R (possibly depending on A)

No retraining/changes to R

