Adversarial Removal of Demographic Attributes from Text Data

Yanai Elazar and Yoav Goldberg Bar-Ilan University / NLP Group November 2, 2018









Text is used for predictions



- For example, consider a text classification setup, where we predict:
 - Hiring decisions
 - Mortgages approval
 - Loans rates





Department of Linguistics & Department of Computer Science, Stanford University Stanford CA 94305-2150

Education

B.A Linguistics, with honors, University of California at Berkeley, 1983
Ph.D. Computer Science, University of California at Berkeley, 1992
Postdoc, International Computer Science Institute, Berkeley, 1992-1995

Academic Employment

Stanford University: Professor and Chair of Linguistics and Professor of Computer Science, 2014-Stanford University: Professor of Linguistics and (by courtesy) of Computer Science, 2010-Stanford University: Associate Professor of Linguistics and (by courtesy) of Computer Science, 2004-2010 University of Colorado, Associate Professor of Linguistics, Computer Science, Cognitive Science, 2001-2003 University of Colorado, Assistant Professor of Linguistics, Computer Science, and Cognitive Science, 1996-2001

This applicant would easily get any NLP job



The common implementation:



Stanford University: Professor and Chair of Linguistics and Professor of Computer Science, 2014-Stanford University: Professor of Linguistics and (by courtesy) of Computer Science, 2010-Stanford University: Associate Professor of Linguistics and (by courtesy) of Computer Science, 2004-2010 University of Colorado, Associate Professor of Linguistics, Computer Science, Cognitive Science, 2001-2003 University of Colorado, Assistant Professor of Linguistics, Computer Science, and Cognitive Science, 1996-2001

Input CV



Hire



ML Model

Don't Hire



The common implementation:



Input CV

6



BUSINESS NEWS

RETAIL

KEUICK3

INVESTING

against women

TECH

APPAREL DISCOUNTERS DEPARTMENT STORES E-COMMERCE FOOD AND BEVERA

recruiting tool that showed bias

Amazon scraps a secret A.I.

POLITICS

CNBC TV

d a big problem: their new

e 2014 to review job search for top talent, five

intelligence to give job uch like shoppers rate

ct 2018

prime





- When deciding on recruiting an applicant based on their writings/CV...
- ...we would like that attributes like the author's:
 - Gender
 - Race
 - Age
- won't be part of the decision.
- In some places, this is even illegal



- We seek to build models which are:
 - Predictive for some main task (e.g. Hiring decision)



• Agnostic to irrelevant/protected attributes (e.g. race, gender, ...)





How do we know we do not condition on some sensitive attribute by mistake?



If we **can** predict protected attributes from the representation...

A talented candidate might suffer from demographic discrimination





Hire



If we **can not** predict protected attributes from the representation...

We don't condition on these protected attributes and... A talented candidate won't suffer from demographic discrimination





In this work:

we do not have access to sensitive tasks like Hiring decisions.

we focus on other tasks, less sensitive

B I U N L P

Let's predict... EMOJIS

We use DeepMoji.

DeepMoji is a model for predicting Emojis from tweets

Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm

Bjarke Felbo¹, Alan Mislove², Anders Søgaard³, Iyad Rahwan¹, Sune Lehmann⁴

¹Media Lab, Massachusetts Institute of Technology
 ²College of Computer and Information Science, Northeastern University
 ³Department of Computer Science, University of Copenhagen
 ⁴DTU Compute, Technical University of Denmark



Let's predict... EMOJIS

I love mom's cooking

I love how you never reply back..

I love cruising with my homies

I love messing with yo mind!!

I love you and now you're just gone..

This is shit

This is the shit





Let's predict... EMOJIS

- DeepMoji is a strong and expressive model
- It also create powerful representations







Let's predict... EMOJIS

- DeepMoji is a strong and expressive model
- It also create powerful representations



results on text classification



Let's predict... EMOJIS

















And use them to predict demographics.

We define: <u>leakage</u> = score above a random guess an "Attacker" achieves We use DeepMoji encoder, to encode tweets, from 3 datasets,
 all binary and balanced

• Each dataset is tied to a different demographic label

• We then train Attackers to predict these attributes



0 1

0 1

··· 1

The dev-set scores above chance level are quite high

Big Surprise?

Not really. This is the core idea in **Transfer-Learning**. We've seen its benefits in pretrained embeddings, language models etc.

Random Guess









- Why do we get this major "help" in predicting other attributes than those we trained for?
- One option is the correlation between attributes in the data.

Fair enough. Let's control for it.



Controlled Setup

25

Rangel et al., 2016

New setup

- We use Twitter data
- We focus on sentiment prediction, emoji based

• With *Race, Gender* and *Age* as protected attributes









New setup





26



Training our own encoder on the balanced datasets



I love messing with yo mind

Balanced Training





Balanced Training - Leakage



We wanted to see something like this:

But instead...



The Attacker manages to extract a substantial amount of sensitive information

Even in a balanced setup, leakage exists





Our objective

B I U N L P

- Create a representation which:
 - Is predictive of the main task (e.g. sentiment)





Our objective

B I U N L P

- Create a representation which:
 - Is predictive of the main task (e.g. sentiment)



Our objective



- Interesting technical problem How to **unlearn** something?
- Interesting technical problem Can we **unlearn** something?









Actively Reducing Leakage

Adversarial Setup (Ganin and Lempitsky, 2015) Classifier 2 - Adv adv(h(x))(Protected Attribute) f(h(x))Classifier 1 gradient reversal layer (Main Task) **Remove stuff** from h/x Representation representation $-\lambda \overline{-\partial L_{adv}}$ Encoder Embeddings

I love messing with yo mind χ

Does it work?



Successfully predicting sentiment

"I love mom's cooking"

Does it work?





Successfully removed demographics?

"I love mom's cooking"

Does it work?





Does it work? Not so quickly...





Does it work? Not so quickly...



Consistent across tasks and protected attributes

Above Chance Scores of Attacker



Does it work? more or less



Well, the adversarial method does help. But not enough





While effective during training, in test time, the adversarial do not remove all the protected information



Can we make stronger adversaries?









Error Analysis

- What are the hard cases, which slip the adversary?
 - We trained the adversarial model 10 times (with random seeds)
 - then, trained the Attacker on each model
 - We collected all examples, which were consistently labeled correctly

• What are the hard cases, which slip the adversary?

AAE("non-hispanic blacks")

Enoy yall day

_ Naw im cool

My Brew Eatting

My momma Bestfrand died

Tonoght was cool

SAE ("non-hispanic whites")

I want to be tan again

Why is it so hot in the house?!

I want to move to california

I wish I was still in Spain

Ahhhh so much homework.

More about the leakage origin can be found in the paper

- When training a text encoder for some task
 - Encoded vectors are also useful for predicting other things ("transfer learning")
 - Including things we did not want to encode ("leakage")
- It is hard to completely prevent such leakage
 - Do not blindly trust adversarial training
 - Verify your model using an "Attacker"

B I U N L P

- We still have a problem
 - During training it seems that the information was removed
 - But the Attacker tells us another story
- Everything we reported was on the dev-set
- Is it possible that we just overfitted on the training-set?

Wait. I remember this thing called Overfitting

- "Adversary overfitting":
 - Memorizing the training data
 - By removing all its sensitive information
 - While leaking in test time

We trained on 90% on the "overfitted" training set, and tested the remaining 10%

Training Set

Above Chance Scores of Attacker Training

It is more than that

• Throughout this work, we aimed in achieving zero leakage, or in

other words: fairness by blindness

- Many other definitions for "fairness" (>20)
- With 3 popular
 - Demographic parity
 - Equality of Odds
 - Equality of Opportunity

In the paper, we prove that in our setup (balanced data) these definitions are identical

