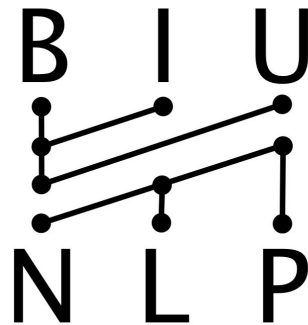# Adversarial Removal of Demographic Attributes from Text Data

**Yanai Elazar** and Yoav Goldberg

Bar-Ilan University / NLP Group

November 2, 2018

Text is used for predictions

- F                                                             e we predict:

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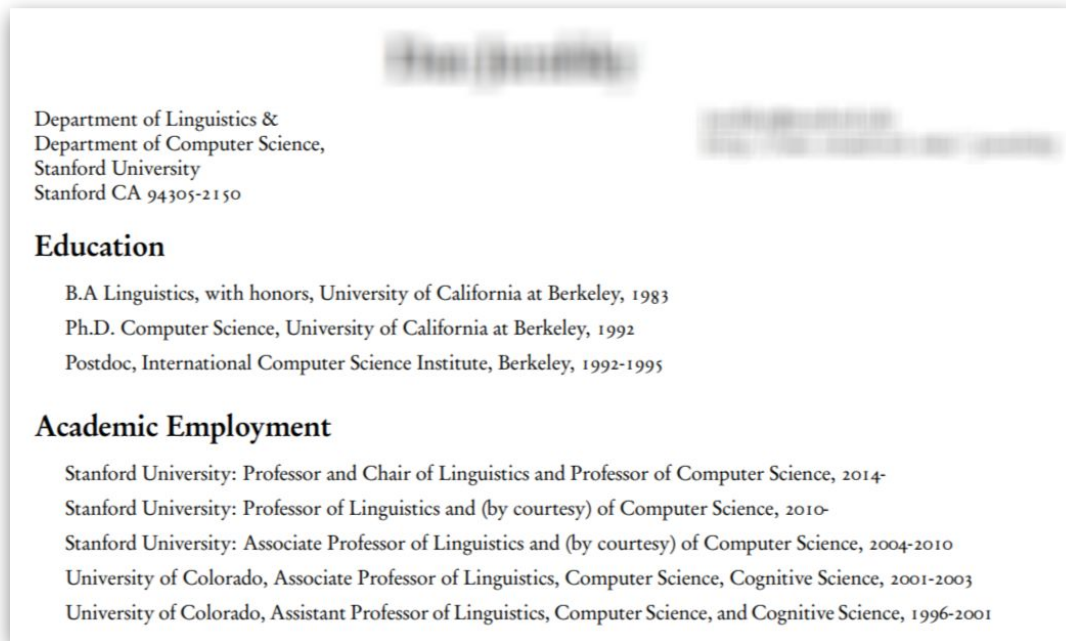Stanford University
Stanford CA 94305-2150

**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
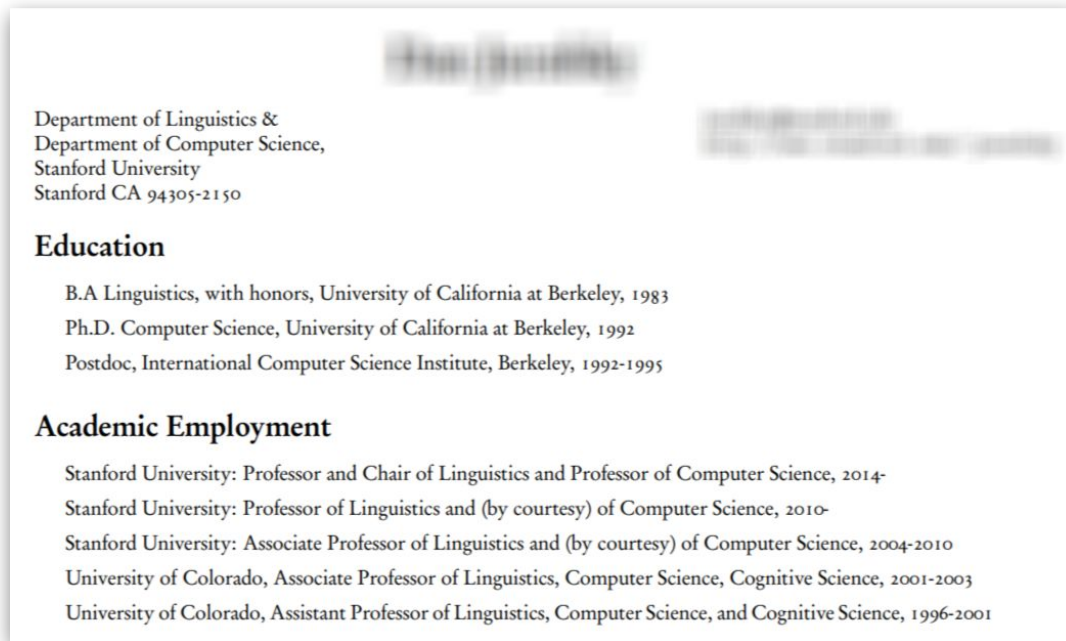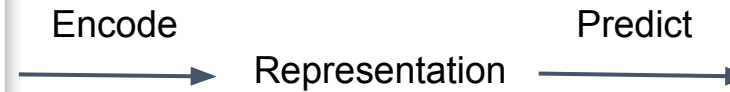
## The common implementation:



*Input CV*

ML
Model

Hire

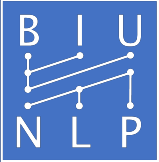Don't Hire

## The common implementation:

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

*Input CV*

Encode → Representation → Predict

Hire

Don't Hire

# Motivation



• But then we see this

- When deciding on recruiting an applicant from his/her 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 attributes (e.g. race, gender, …)

We do not have access to sensitive tasks like Resumes.

We will focus on other tasks, less sensitive

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 Felbo[1], Alan Mislove[2], Anders Søgaard[3], Iyad Rahwan[1], Sune Lehmann[4]

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

Let's predict... EMOJIS

*Deep Moji (Felbo et al., 2017)*

I love mom's cooking

| 😜 | 😍 | ❤️ | ☺️ | 🤍 |
|------|------|------|------|------|
| 49.1% | 8.8% | 3.1% | 3.0% | 2.9% |

I love how you never reply back..

| 😒 | 😑 | 😠 | 😐 | 💔 |
|------|------|------|------|------|
| 14.0% | 8.3% | 6.3% | 5.4% | 5.1% |

I love cruising with my homies

| 😎 | 👌 | ✌️ | ☺️ | 💯 |
|------|------|------|------|------|
| 34.0% | 6.6% | 5.7% | 4.1% | 3.8% |

I love messing with yo mind!!

| 😜 | 😈 | 😏 | 😉 | 🙊 |
|------|------|------|------|------|
| 17.2% | 11.8% | 8.0% | 6.4% | 5.3% |

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

| 💔 | 😔 | 😞 | 😪 | 😢 |
|------|------|------|------|------|
| 39.1% | 11.0% | 7.3% | 5.3% | 4.5% |

This is shit

| 😠 | 😡 | 😞 | 😒 | 😤 |
|------|------|------|------|------|
| 7.0% | 6.4% | 6.0% | 6.0% | 5.8% |

This is the shit

| 🎧 | 🎶 | 👌 | 😎 | 😏 |
|------|------|------|------|------|
| 10.9% | 9.7% | 6.5% | 5.7% | 4.8% |

11

Let's predict... EMOJIS

Deep Moji (Felbo et al., 2017)

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

Encode

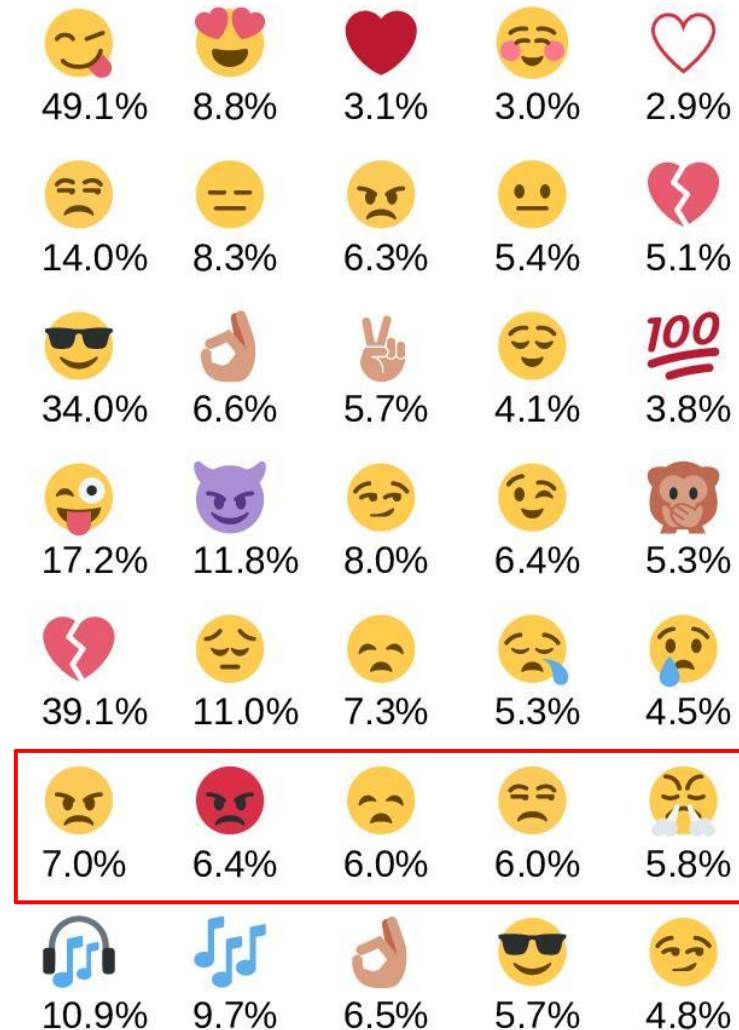Predict

Let's predict... EMOJIS

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

Encode

Predict

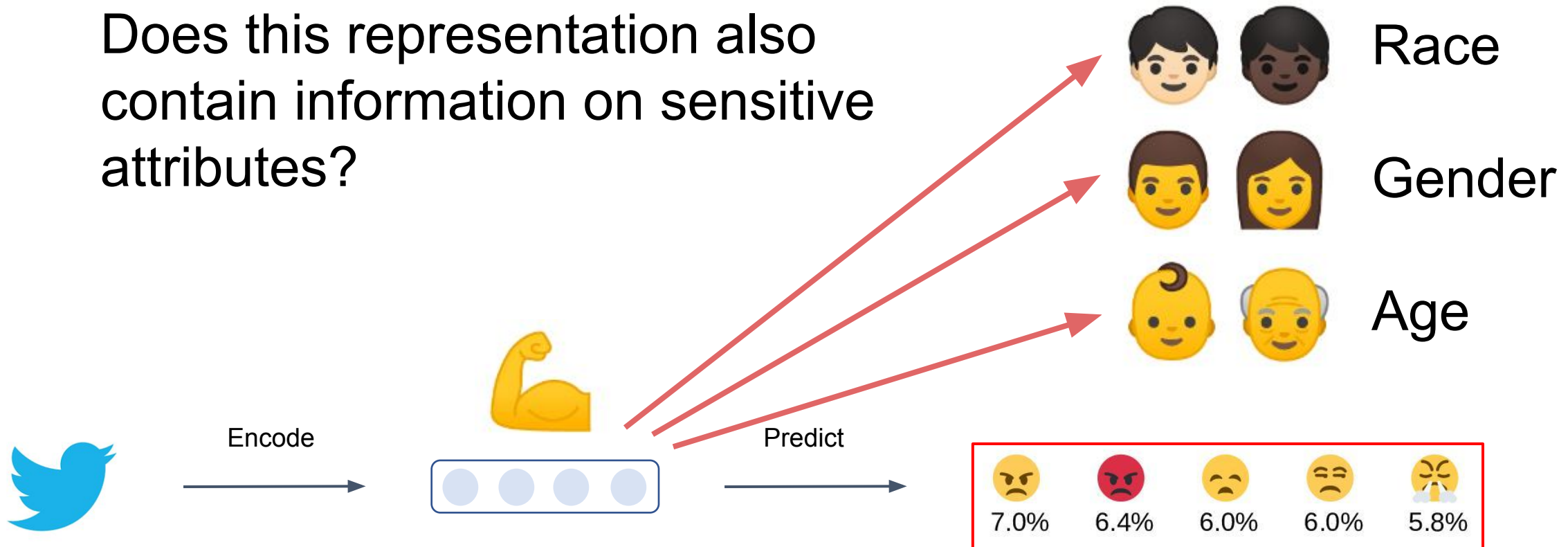- ● Achieved several SOTA results on text classification
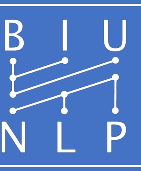
Deep Moji (Felbo et al., 2017)

| 😜 | 😍 | ❤️ | 🥰 | 🤍 |
|---|---|---|---|---|
| 49.1% | 8.8% | 3.1% | 3.0% | 2.9% |
| 😒 | 😑 | 😠 | 😐 | 💔 |
| 14.0% | 8.3% | 6.3% | 5.4% | 5.1% |
| 😎 | 👌 | ✌️ | 😌 | 💯 |
| 34.0% | 6.6% | 5.7% | 4.1% | 3.8% |
| 😜 | 😈 | 😏 | 😉 | 🙊 |
| 17.2% | 11.8% | 8.0% | 6.4% | 5.3% |
| 💔 | 😔 | 😞 | 😪 | 😢 |
| 39.1% | 11.0% | 7.3% | 5.3% | 4.5% |
| 😠 | 😡 | 😞 | 😒 | 😤 |
| 7.0% | 6.4% | 6.0% | 6.0% | 5.8% |
| 🎧 | 🎶 | 👌 | 😎 | 😏 |
| 10.9% | 9.7% | 6.5% | 5.7% | 4.8% |

Let's predict... EMOJIS

Does this representation also contain information on sensitive attributes?

Race

Gender

Age

Encode

Predict

7.0%   6.4%   6.0%   6.0%   5.8%

Task
(Emojis)

We use the
representation that
predict Emojis

😃 | ☹️

Deep Moji (Felbo et al., 2017)

Classifier

Representation  $h(x)$

Encoder  *DeepMoji Encoder*

Embeddings

I love messing with yo mind   $x$

# Setup

BIU
NLP

Task
(Emojis)

😀 | ☹️

Demographics
(Gender)

👩 | 👨

a.k.a. **Attacker**

Classifier

We use the
representation that
predict Emojis

Representation  ⬤⬤⬤⬤  $h(x)$

And use them to predict
demographics.

We define:
**leakage** = 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

a.k.a. **Attacker**

Demographics
(e.g. Gender)

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



Above Chance Scores of DeepMoji representation

DeepMoji

- Why do we get this major "help" in predicting other

  attributes than those we trained on?

- One option is the correlation between attributes in
  the data

Fair enough. Let's control it

# Controlled Setup

- ## We use Twitter data

- ## We focus on sentiment prediction, emoji based

- ## With *Race*, *Gender* and *Age* as protected attributes

*Blodgett et al., 2016*                    *Rangel et al., 2016*                    *Rangel et al., 2016*

**Balanced Dataset**

Demographics

50% Male          50% Female

Task (Sentiment)

50% Positive

50% Negative

$M_+ P_-$          $M_+ P_+$

$M_- P_-$          $M_- P_+$

# Training our own encoder on the balanced datasets

Main Task (sentiment)

Classifier

Representation

50% Male          50% Female

50%
Positive          $M_+P_-$          $M_+P_+$

50%
Negative          $M_-P_-$          $M_-P_+$

Encoder

Embeddings

I love messing with yo mind

And using the Attacker to check for leakage

Protected Attribute (gender)

$att(h(x))$

**Trainable**

Representation

$h(x)$

50% Male

50% Female

Encoder

**Freeze**

50% Positive

50% Negative

$M_+P_-$

$M_+P_+$

$M_-P_-$

$M_-P_+$

Embeddings

I love messing with yo mind

$x$

We wanted to see something like this:

But instead...



Above Chance Scores

Random Guess

The Attacker manages to extract a substantial amount of sensitive information

**Even in a balanced setup, leakage exists**



Above Chance Scores

Random Guess

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

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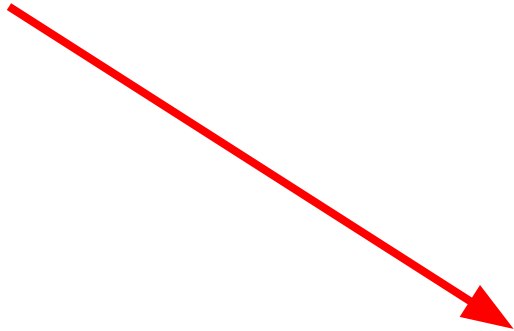  - Is predictive of the main task (e.g. sentiment)

and

  - Is **not** predictive of protected attribute (e.g. gender, race)

- Interesting technical problem – How to **unlearn** something?
- Interesting technical problem – <span style="color:red">Can</span> we **unlearn** something?

# Actively Reducing Leakage

# Adversarial Setup

- First introduced by Goodfellow et al., 2014

  - A very active line of research

  - We will go through the details

## Generative Adversarial Nets

Ian J. Goodfellow,  Jean Pouget-Abadie[*] Mehdi Mirza,  Bing Xu,  David Warde-Farley,
Sherjil Ozair[†] Aaron Courville,  Yoshua Bengio[‡]
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

- The motivation came from "Generative Models"

  - We would like to automatically create images

  - From… random input?

- 2 components:

  - Generator

  - Discriminator

$$\min_{G} \boxed{\max_{D}} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \boxed{[\log D(\boldsymbol{x})]} + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))].$$

A good Discriminator
(real data gets a high score,
meaning it's real)

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))].$$

A good Generator
(fake data gets a high score,
for maximizing *D*'s probability)

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \boxed{[\log D(\boldsymbol{x})]} + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} \boxed{[\log(1 - D(G(\boldsymbol{z})))]}.$$

- 2 competing objectives.
- We don't know how to solve this

$$\min_G \max_D V(D, G) = \boxed{\mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})]} + \boxed{\mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]}.$$

Goodfellow et al. solution:
iterate training between the Generator and Discriminator

- Update the discriminator by ascending its stochastic gradient:

- Update the generator by descending its stochastic gradient:

- The Adversarial setup was invented to create an "output"

- Which can't (or seem hard) to separate real from fake

- What if we want to create an intermediate representation?

# Adversarial Setup

- The Adversarial setup was invented to create an "output"

- Which can't (or seem hard) to separate real from fake

- What if we want to create an intermediate representation…

  - Which is indistinguishable for some feature or attribute?

# Adversarial Setup

- Ganin and Lempitsky, 2015

- Application: Domain Adaptation

- New trick for adversary train: Gradient Reversal Layer (GRL)

---

## Unsupervised Domain Adaptation by Backpropagation

---

**Yaroslav Ganin**                                    GANIN@SKOLTECH.RU
**Victor Lempitsky**                                  LEMPITSKY@SKOLTECH.RU
Skolkovo Institute of Science and Technology (Skoltech)

**Predict Sentiment**

$f(h(x))$

Classifier 1
(Main Task)

Representation    $h(x)$

Encoder

Embeddings

*I love messing with yo mind*    $x$

**Predict Sentiment**                    **Predict Race**

$f(h(x))$                    Classifier 1          Classifier 2 - Adv          $adv(h(x))$
                             (Main Task)          (Protected Attribute)

                                                  **try to interfere**

                             Representation  ⬜⬜⬜⬜          $h(x)$

                             Encoder

                             Embeddings  ⬜⬜ ⬜⬜ ⬜⬜

                             *I love messing with yo mind*          $x$

**3 different sub-objectives**

$f(h(x))$

Classifier 1
(Main Task)

$h(x)$

$x$

classify well

$adv(h(x))$

Classifier 2 - Adv
(Protected Attribute)

$h(x)$

$x$

adversary should
succeed

$-adv(h(x))$

Classifier 2 - Adv
(Protected Attribute)

$h(x)$

$x$

encoder should
make adversary
fail

$f(h(x))$

$adv(h(x))$

$-adv(h(x))$

Classifier 1
(Main Task)

Classifier 2 - Adv
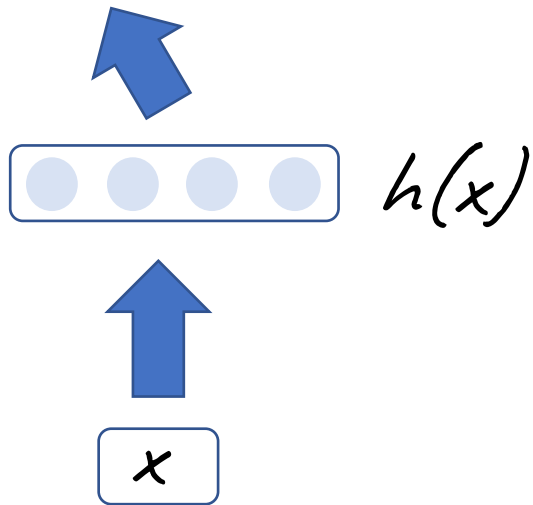(Protected Attribute)

Classifier 2 - Adv
(Protected Attribute)

$h(x)$

$h(x)$

$h(x)$

$x$

$x$

$x$

**blue**: update parameters
**white**: don't update

$f(h(x))$

$adv(h(x))$

$-adv(h(x))$

Classifier 1
(Main Task)

Classifier 2 - Adv
(Protected Attribute)

Classifier 2 - Adv
(Protected Attribute)

$h(x)$

$h(x)$

$h(x)$

$x$

$x$

$x$

**blue**: update parameters
**white**: don't update

$grad(-adv(h(x)))$

$f(h(x))$               $adv(h(x))$           $-adv(h(x))$

Classifier 1              Classifier 2 - Adv        Classifier 2 - Adv
(Main Task)           (Protected Attribute)     (Protected Attribute)

$h(x)$                $h(x)$               $h(x)$

$x$                 $x$                 $x$

**blue**: update parameters
**white**: don't update

$grad(-adv(h(x)))=-grad(adv(h(x)))$

46

Classifier 2 - Adv
(Protected Attribute)

$adv(h(x))$

$f(h(x))$

Classifier 1
(Main Task)

gradient reversal layer

Remove stuff
from
representation

Representation

$h(x)$

Encoder

$-\lambda \dfrac{\partial L_{adv}}{adv(h(x))}$

Embeddings

*I love messing with yo mind*
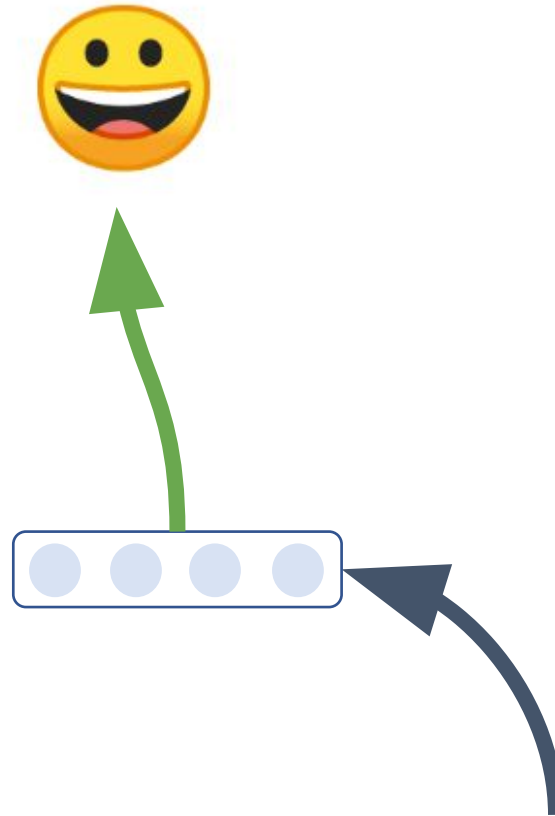
$x$

- In their paper, the representation after the adversarial training seems invariant to the domain



*before*                              *after*
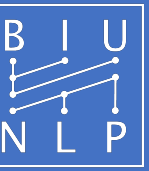
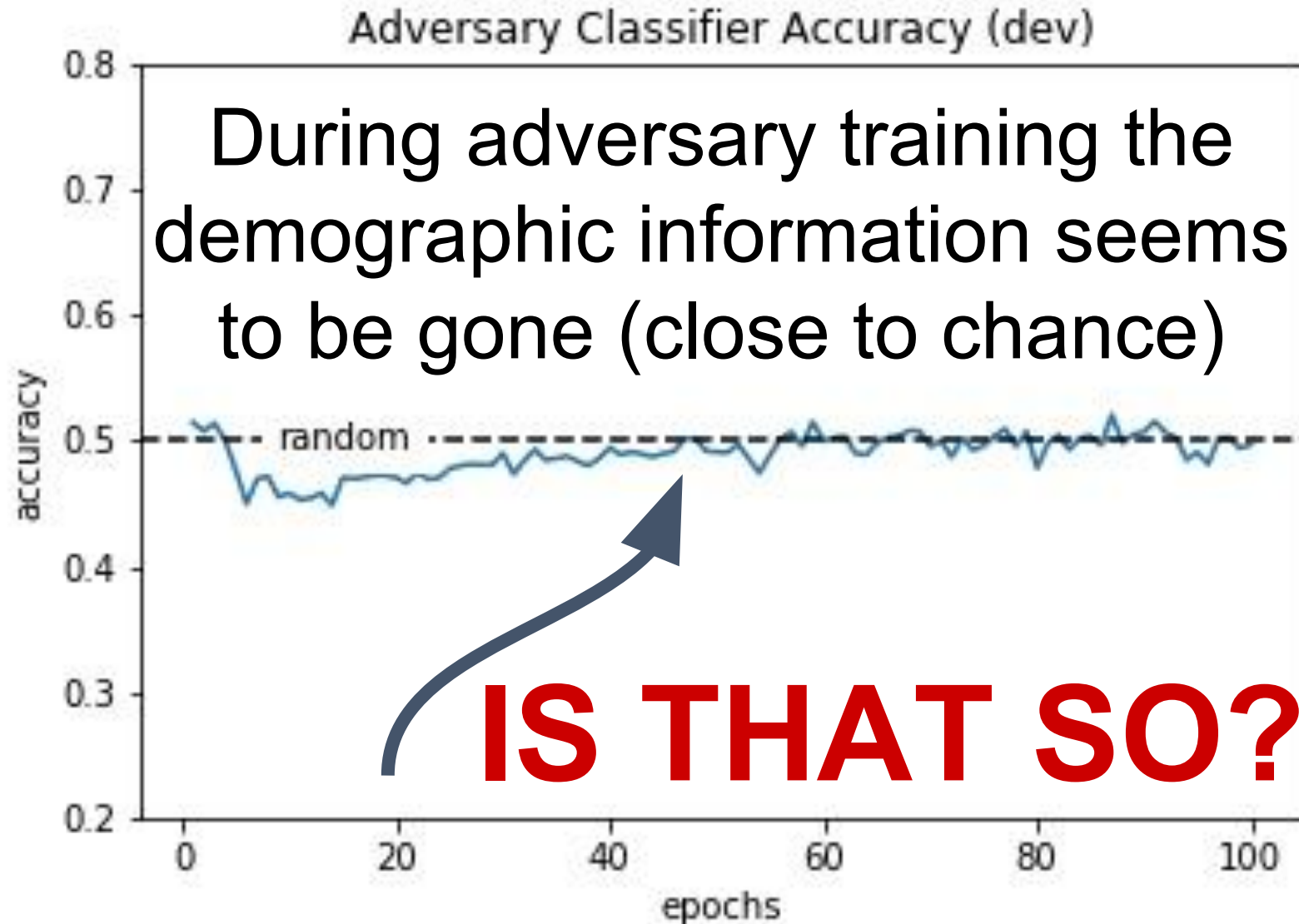**Successfully
predicting sentiment**
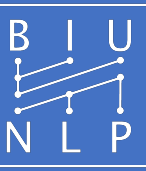
"I love mom's cooking"

Successfully removed demographics?

"I love mom's cooking"

Adversary Classifier Accuracy (dev)

During adversary training the demographic information seems to be gone (close to chance)

**IS THAT SO?**

BIU
NLP

When training the Attacker

We can still recover a considerable amount of information



Attacker Classifier Accuracy (dev)

**Consistent across tasks and protected attributes**

## Above Chance Scores of Attacker



| | Leakage |

10

Leakage

0

Sentiment-Race · Mention1-Race · Mention2-Gender · Mention2-Age

Random Guess

53

Well, the adversarial method does help.
But not enough



**Adversarial Contribution**

Legend: Adv, No-Adv

Categories: Sentiment-Race, Mention2-Gender, Mention2-Age

Random Guess

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

# Can we make stronger adversaries?

**More Parameters!** ~~Baselines!~~

Classifier 2 - Adv
(Protected Attribute)

$adv(h(x))$

$f(h(x))$

Classifier 1
(Main Task)

**gradient reversal layer**

Representation

$h(x)$

$$-\lambda \frac{\partial L_{adv}}{adv(h(x))}$$

Encoder

Embeddings

*I love messing with yo mind*

$x$

# Stronger, Better, Bigger???

Bigger Weight!
~~Baseline~~

$f(h(x))$

Classifier 1
(Main Task)

Classifier 2 - Adv
(Protected Attribute)

$adv(h(x))$

gradient reversal layer

Representation

$h(x)$

$$-\lambda \frac{\partial L_{adv}}{adv(h(x))}$$

Scale the reverse gradients

Encoder

Embeddings

*I love messing with yo mind*

$x$

58

BIU NLP

**More Adversaries!**

Classifier 3 - Adv
(Protected Attribute)

Classifier 2 - Adv
(Protected Attribute)

$adv(h(x))$

$f(h(x))$

Classifier 1
(Main Task)

**gradient reversal layer**

**gradient reversal layer**

Representation

$h(x)$

Encoder

$$-\lambda \sum_i \frac{\partial L_{adv_i}}{adv(h(x))}$$

Embeddings

*I love messing with yo mind*

$x$

# Stronger, Better, Bigger???

**Better, but still not perfect**

# Error Analysis

- 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
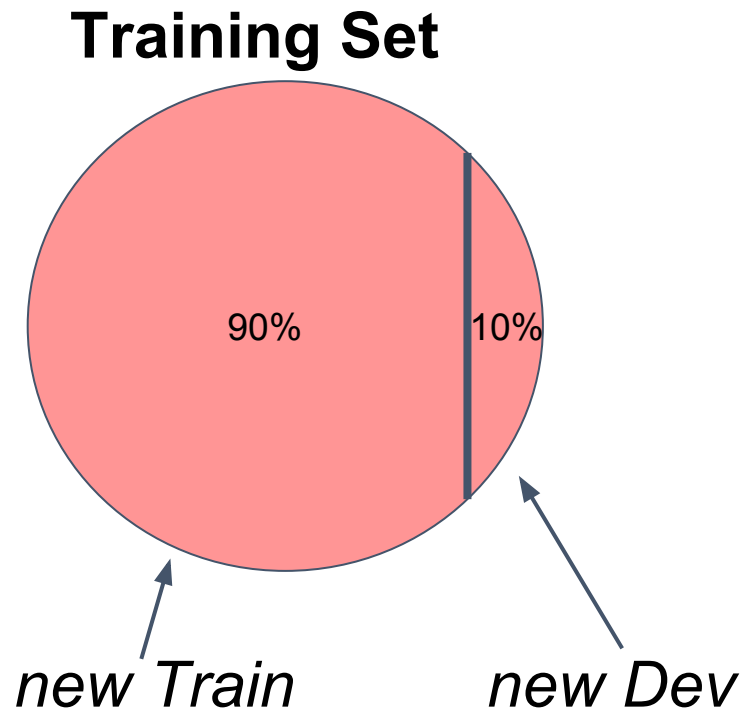- Is it possible that we just overfitted on the training-set?

- "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**

90%　10%

*new Train*　*new Dev*

Above Chance Scores of Attacker Training

Leakage

12.2 — Sentiment-Race
14.3 — Mention1-Race
8.1 — Mention2-Gender
9.7 — Mention2-Age

**It is more than that**

- 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

# Persistent Examples

AAE("non-hispanic blacks")

Enoy yall day

_ Naw im cool

My Brew Eatting

My momma Bestfrand died

Tonoght was cool

SAE ("non-hispanic whites")

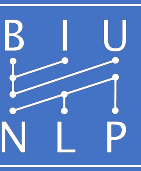I want to be tan again

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*

- 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 out setup (balanced data) these definitions are identical

# Summary

- When training a text encoder for some task
  - Encoded vectors are still useful for predicting various 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**
  - **Recheck your model using an "Attacker"**

Thank you