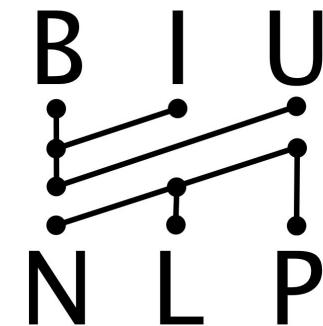


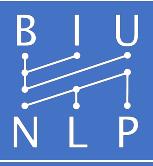
Adversarial Removal of Demographic Attributes from Texts and Beyond

Yanai Elazar

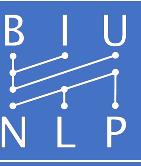
Bar-Ilan University / NLP Group

Textkernel, February 28, 2018





Our data is used for predictions



Motivation

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

• we predict:

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

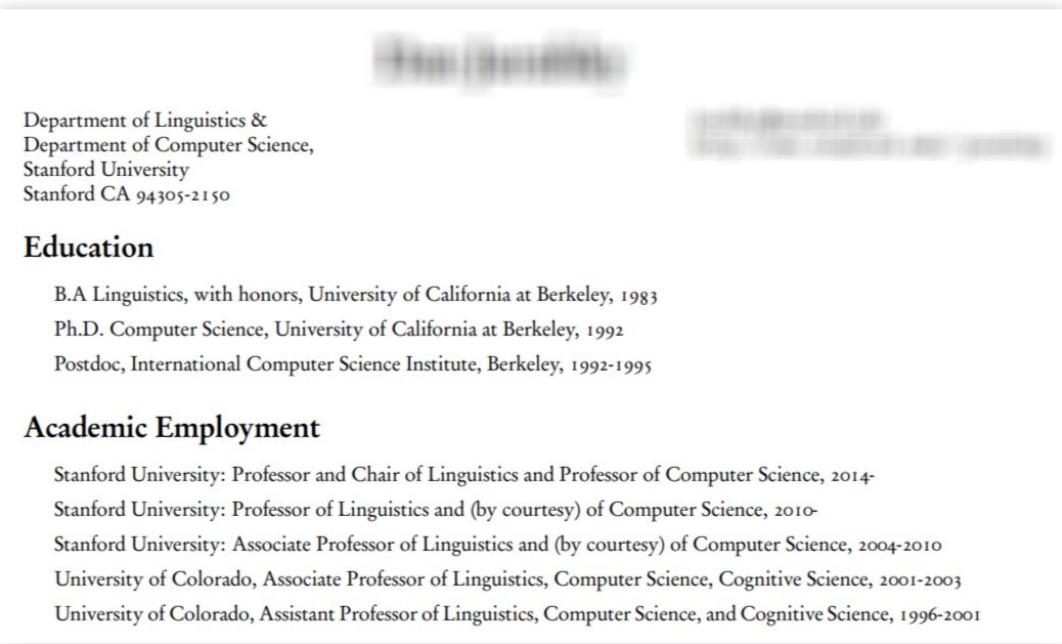
University of Colorado, Associate Professor of Linguistics, Computer Science, Cognitive Science, 2001-2003

University of Colorado, Assistant Professor of Linguistics, Computer Science, and Cognitive Science, 1996-2001

This applicant would easily get any NLP job

Motivation

The common implementation:



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

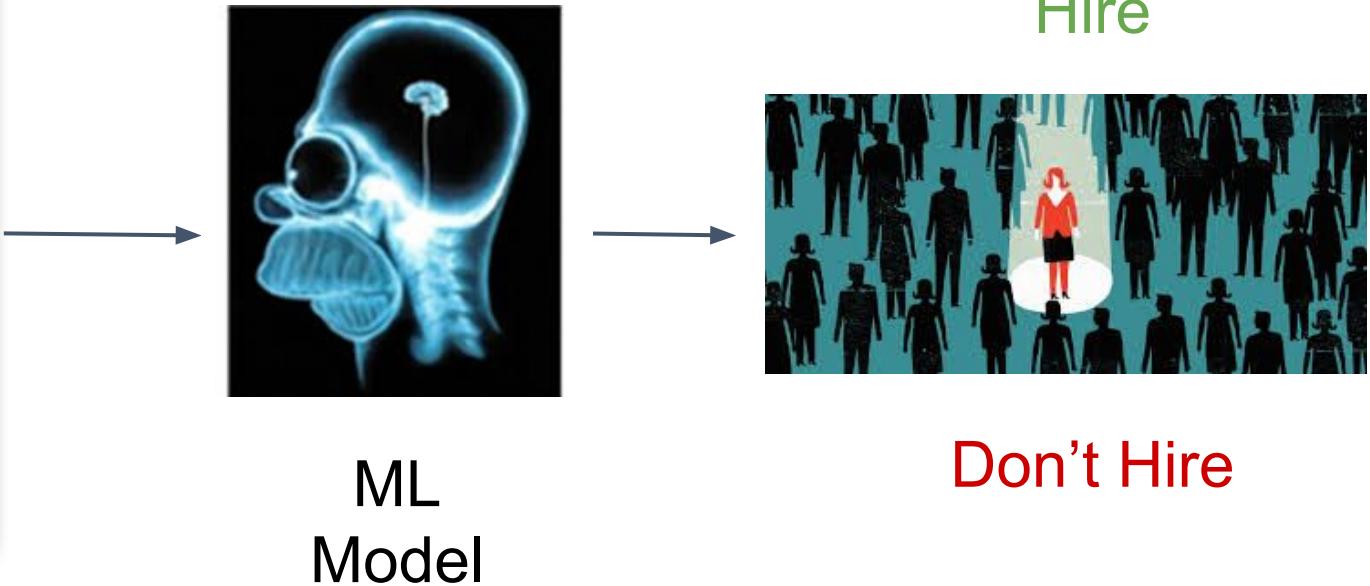
Education

- 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



Motivation

The common implementation:

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

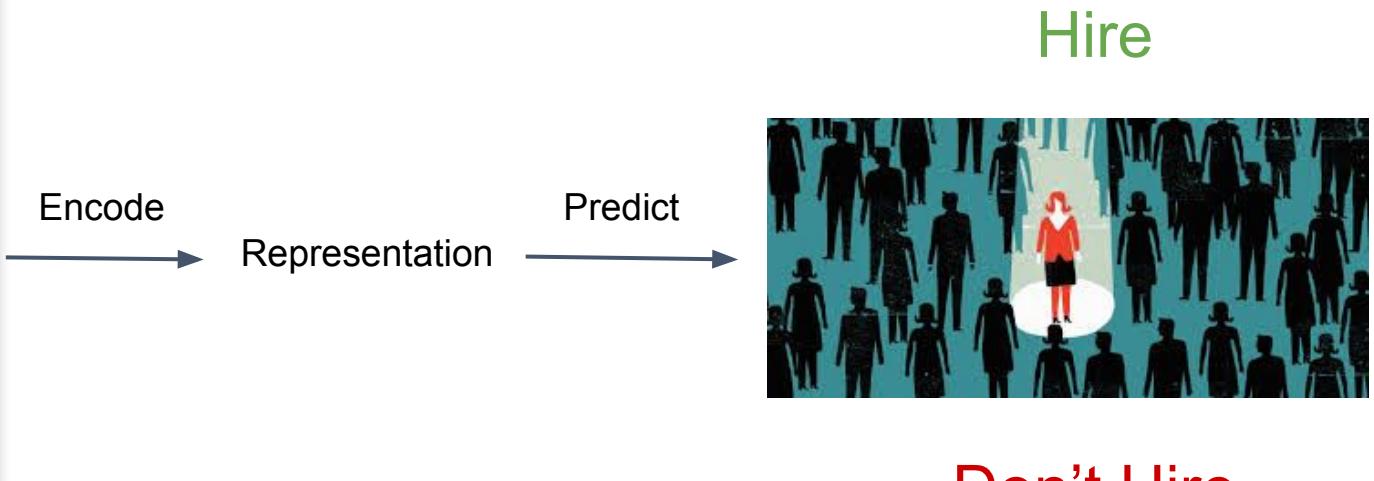
Education

- 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



Motivation

BUSINESS
INSIDER

TECH | FINANCE | POLITICS | STRATEGY | LIFE | ALL

PRIME | INTELLIGENCE



Amazon built an AI tool to hire people but had to shut it down because it was discriminating against women

Isobel Asher Hamilton 19h

But then we see this



BUSINESS NEWS

OCTOBER 10, 2018 / 6:12 AM / UPDATED 16 HOURS AGO

Amazon scraps secret AI recruiting tool that showed bias against women



CNBC

MENU

MARKETS

BUSINESS NEWS

INVESTING

TECH

POLITICS

CNBC TV

RETAIL

APPAREL | DISCOUNTERS | DEPARTMENT STORES | E-COMMERCE | FOOD AND BEVERAGE

Amazon scraps a secret A.I. recruiting tool that showed bias against women

had a big problem: their new

e 2014 to review job
search for top talent, fiveintelligence to give job
uch like shoppers rate

REUTERS



Motivation

- 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

Motivation

- 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, ...)



Motivation

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

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

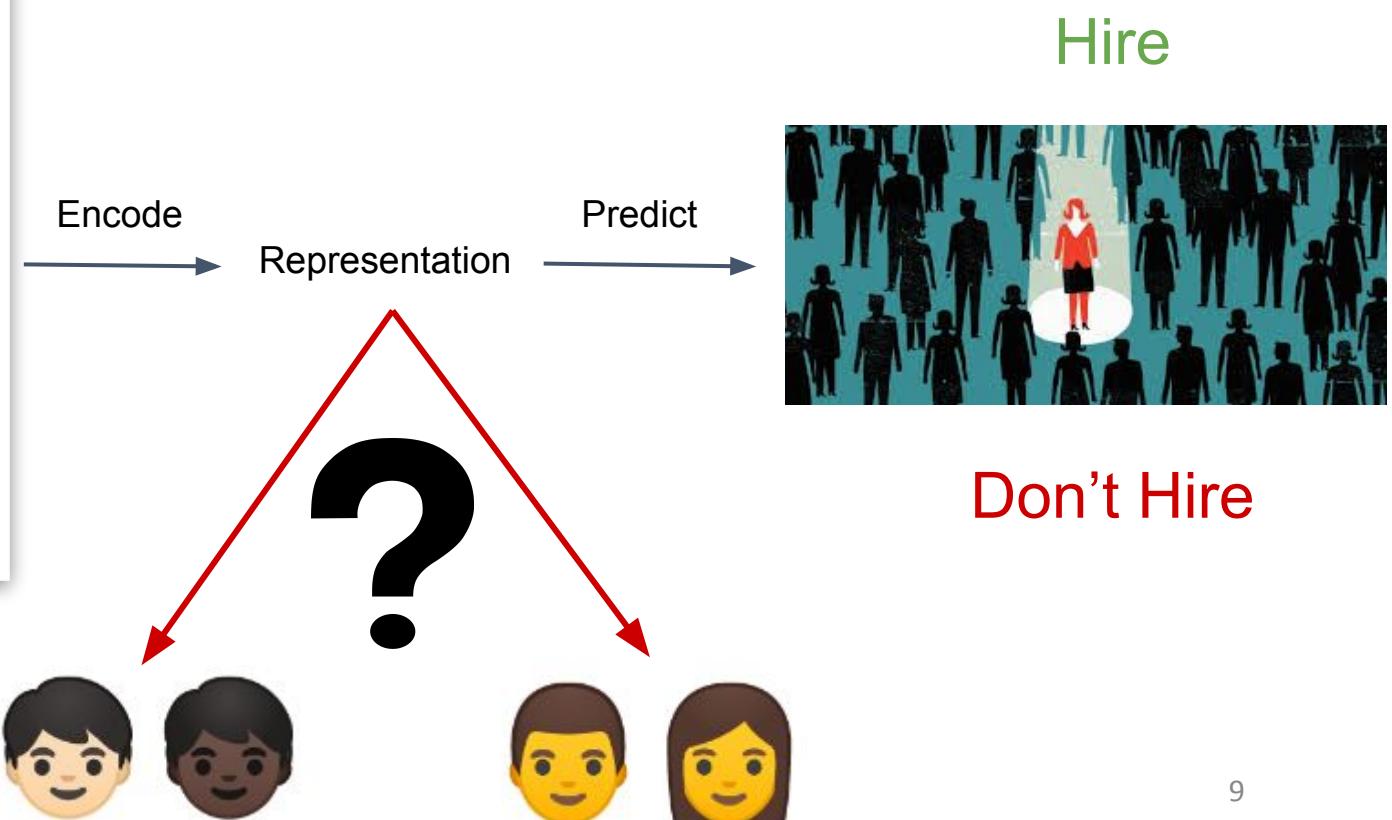
Education

- 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

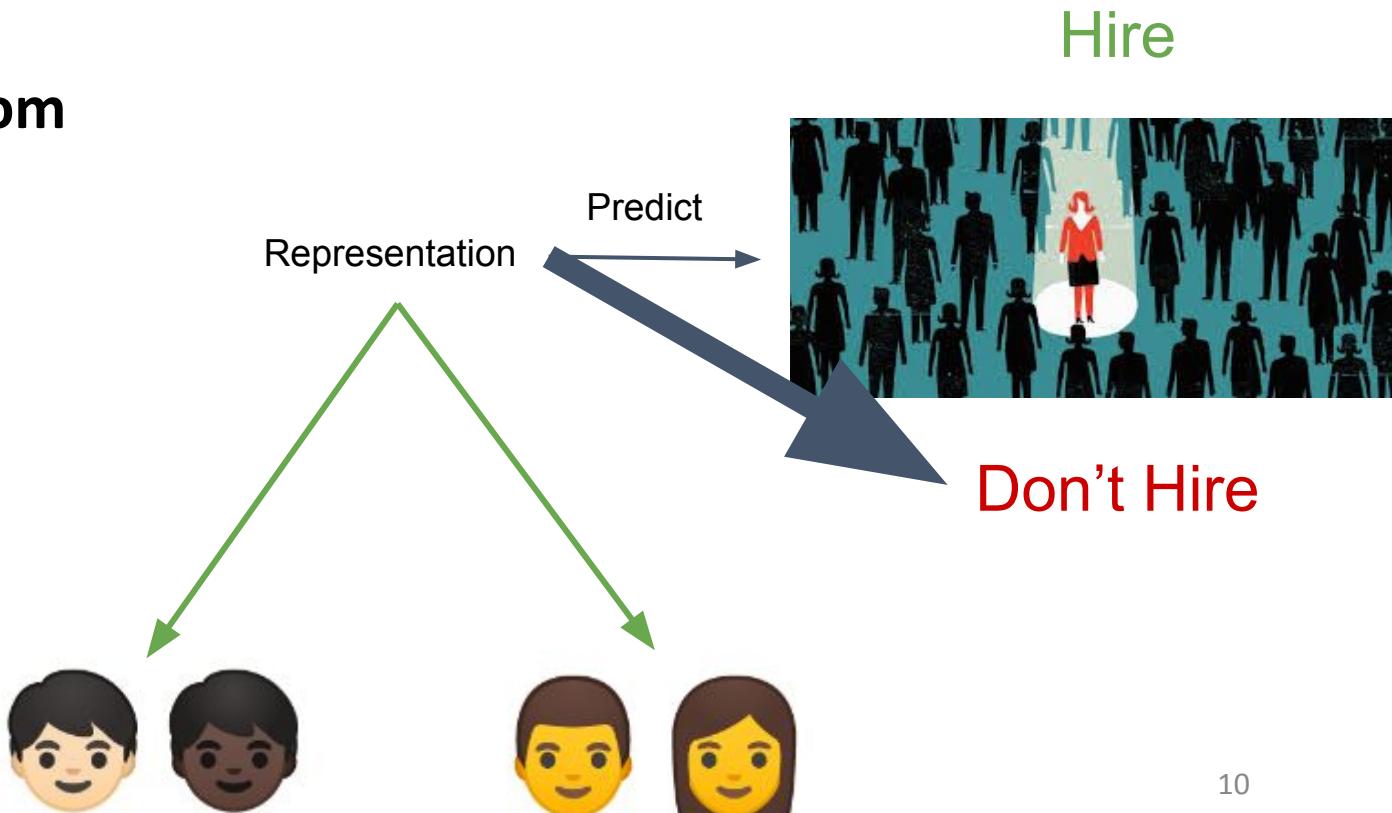
Input CV



Motivation

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

A talented candidate might suffer from demographic discrimination

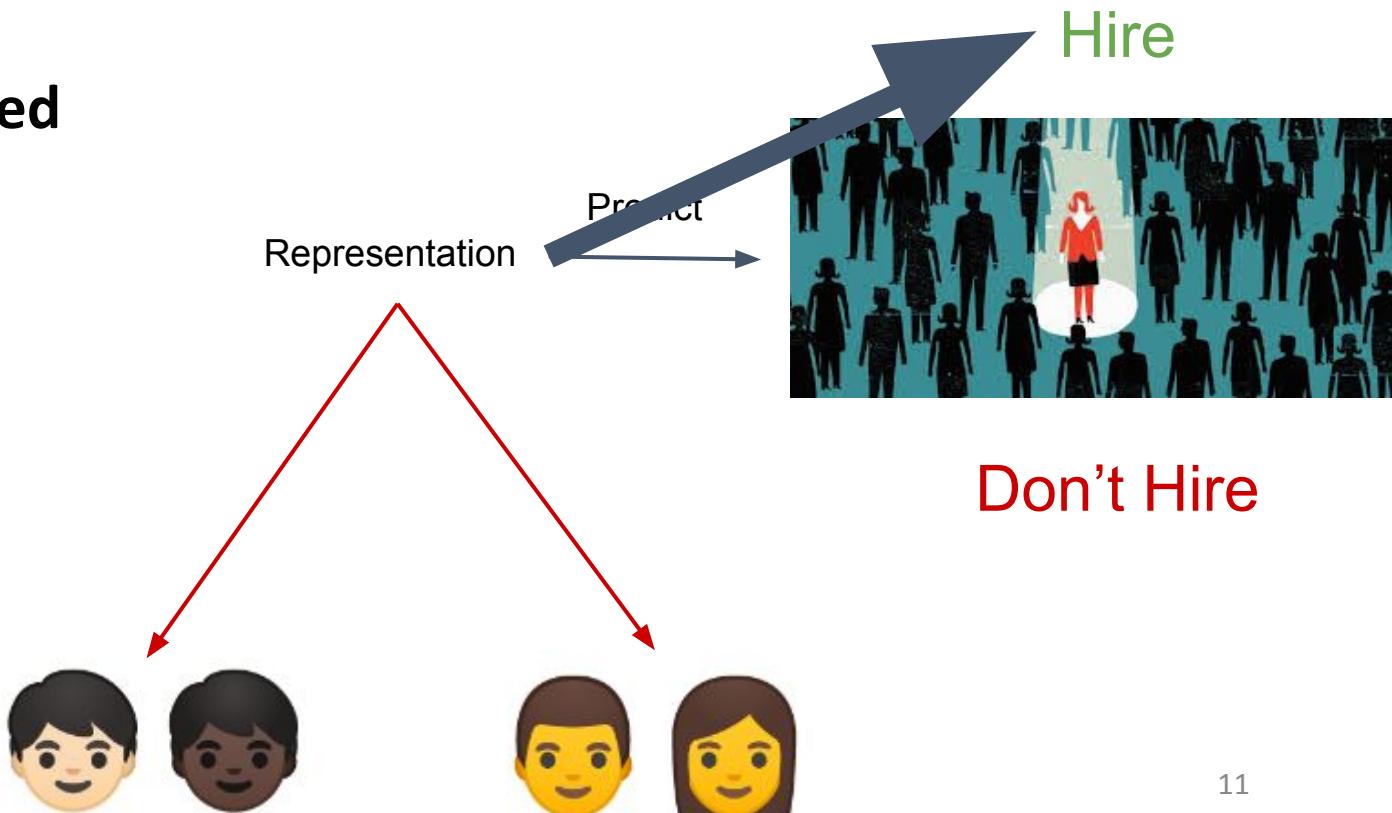


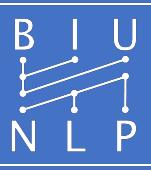
Motivation

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





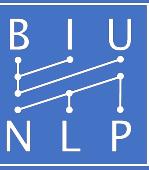
Adversarial Removal of Demographic Attributes from Text Data

Yanai Elazar[†] and Yoav Goldberg^{†*}

[†]Computer Science Department, Bar-Ilan University, Israel

^{*}Allen Institute for Artificial Intelligence

{yanaiela, yoav.goldberg}@gmail.com

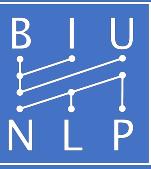


Text classification - Example

In this work:

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

we focus on other tasks, less sensitive



Text classification - Example

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

Text classification - Example

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

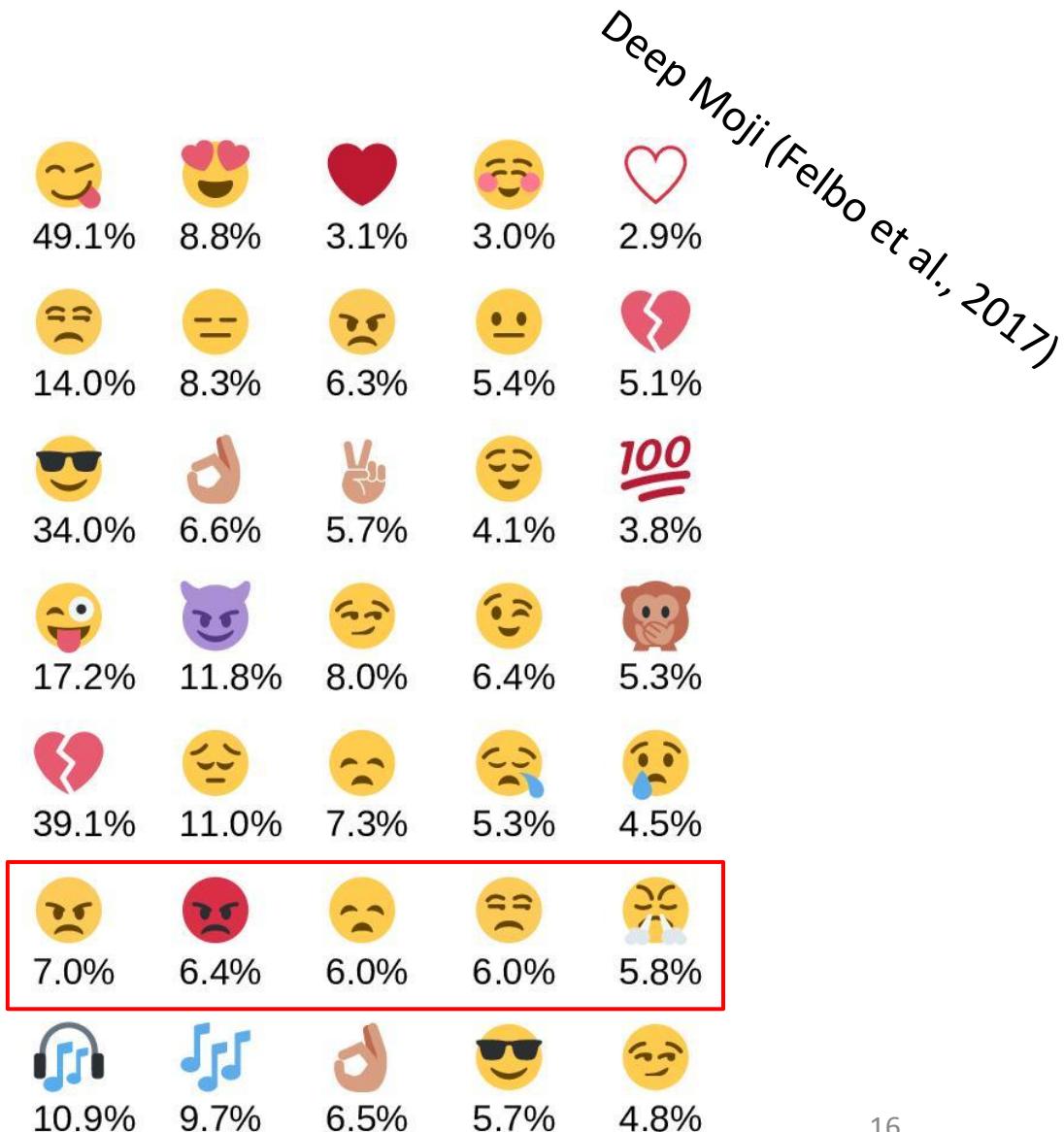
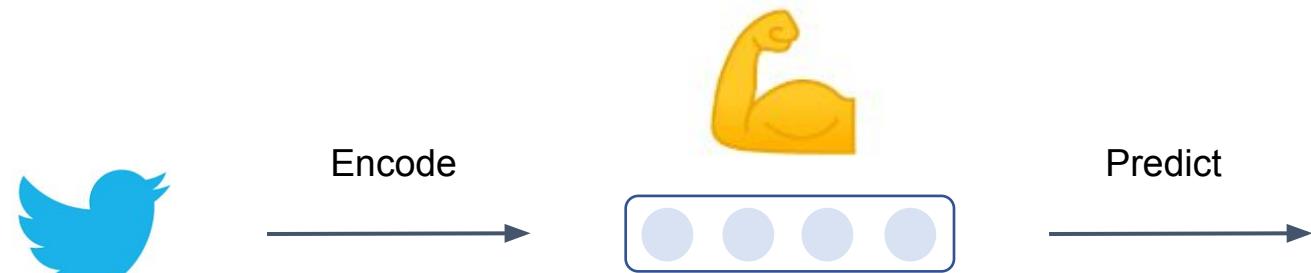


Deep Moji (Felbo et al., 2017)

Text classification - Example

Let's predict... EMOJIS

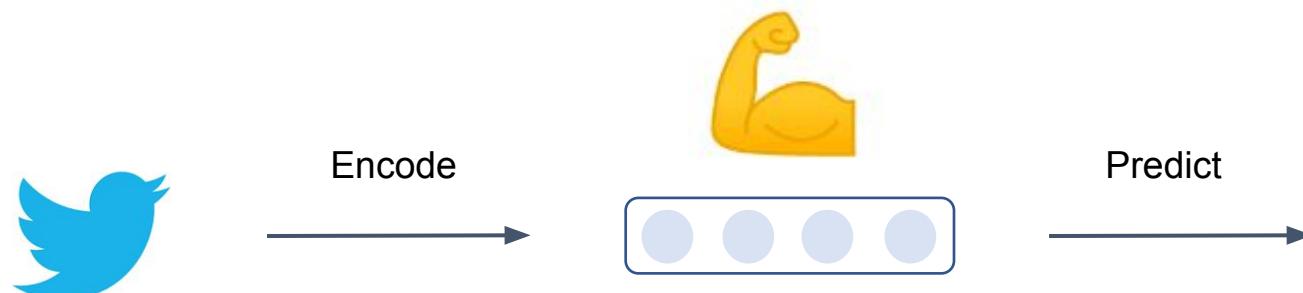
- DeepMoji is a strong and expressive model
- It also create powerful representations



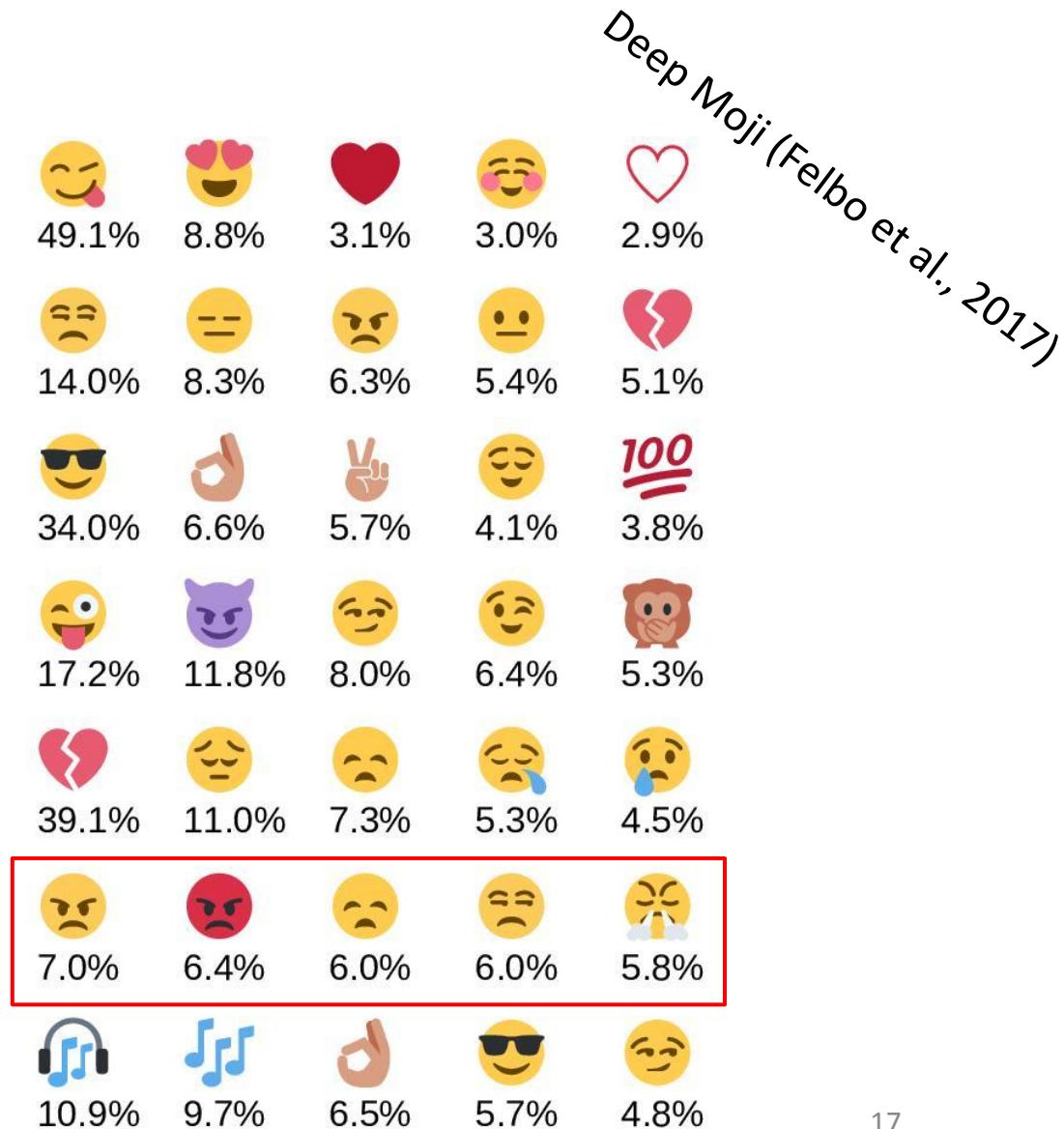
Text classification - Example

Let's predict... EMOJIS

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



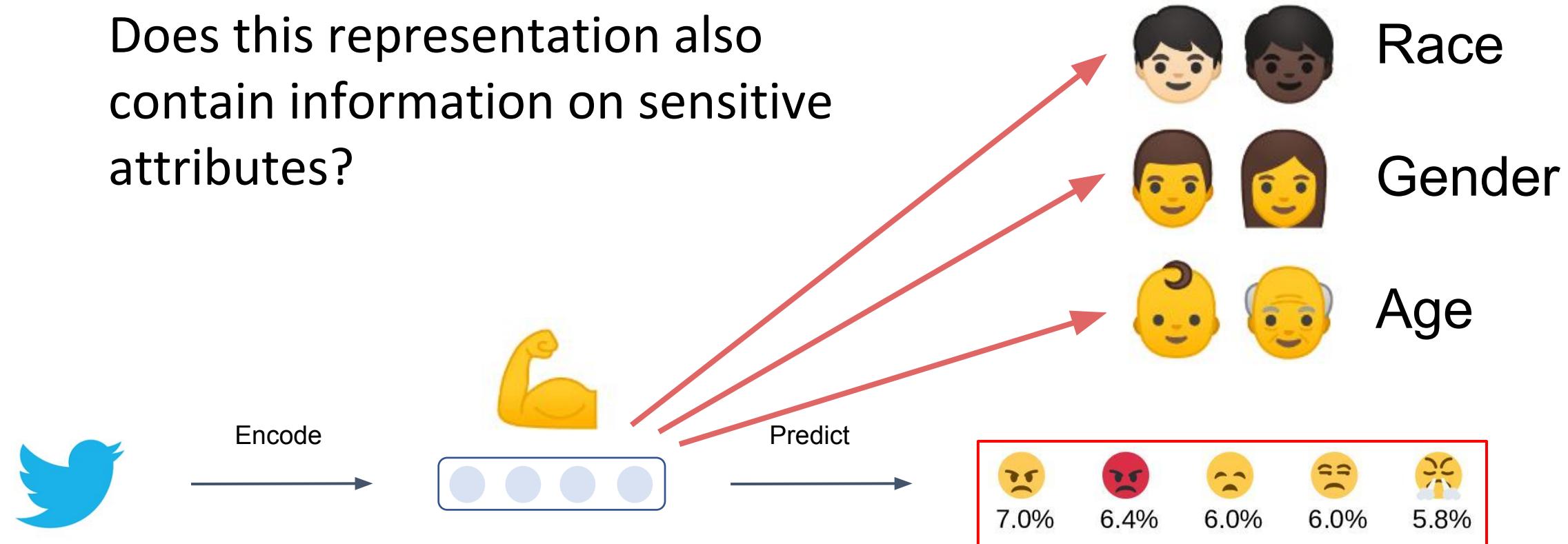
- Achieved several SOTA results on text classification



Text classification - Example

Let's predict... EMOJIS

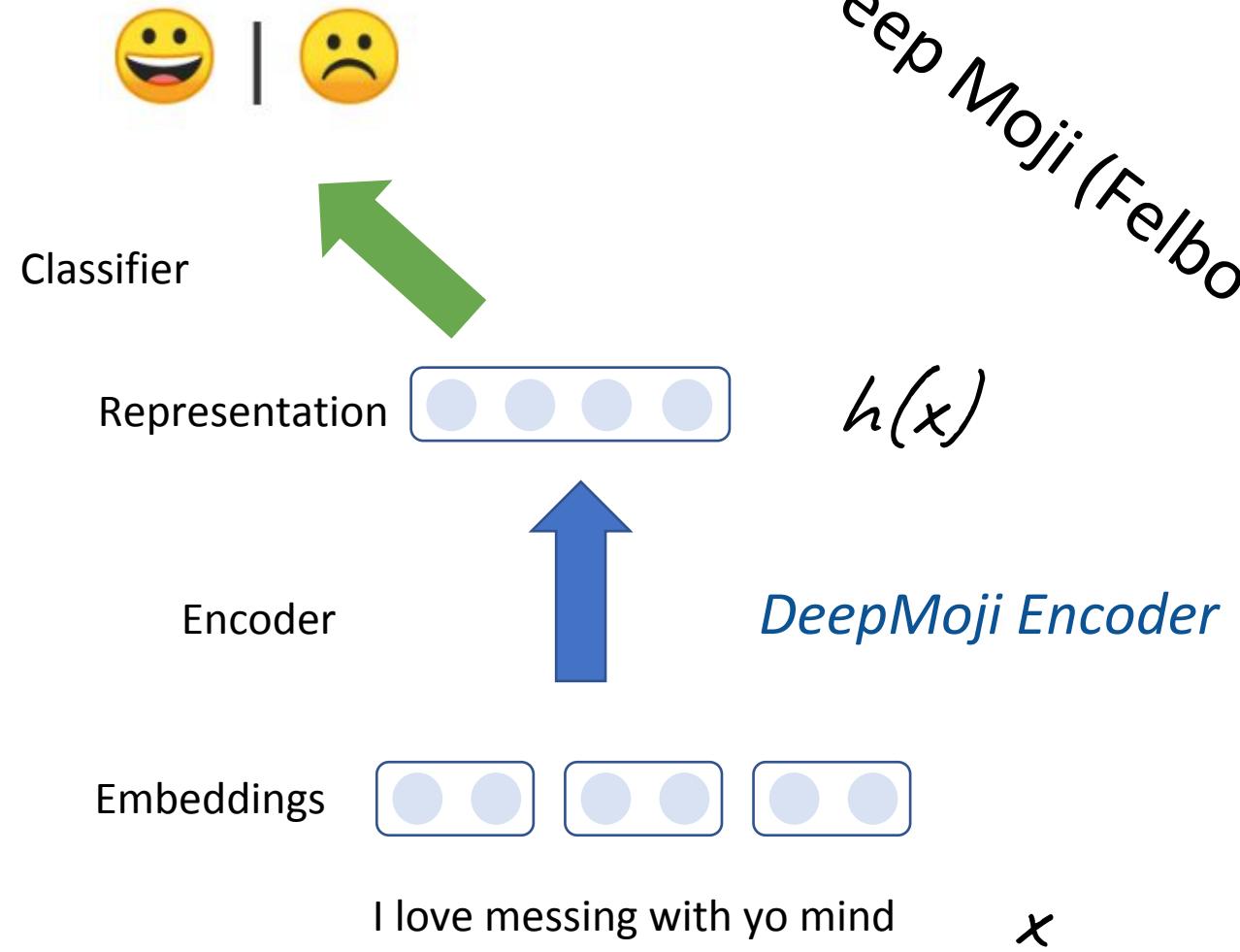
Does this representation also contain information on sensitive attributes?



Setup

Task
(Emojis)

We take the representation that predict Emojis



Setup

Task
(Emojis)

We take the representation that predict Emojis

Classifier

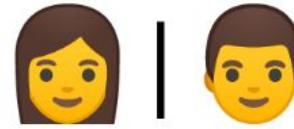


Representation



a.k.a. Attacker

$h(x)$



Demographics
(Gender)

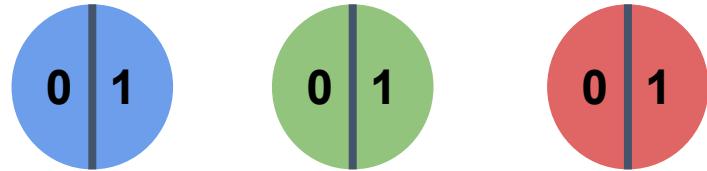
And use them to predict demographics.

We define:

leakage = score above a random guess an “Attacker” achieves

Text Leakage – Case Study

- 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



Demographics
(e.g. Gender)

Text Leakage – Case Study

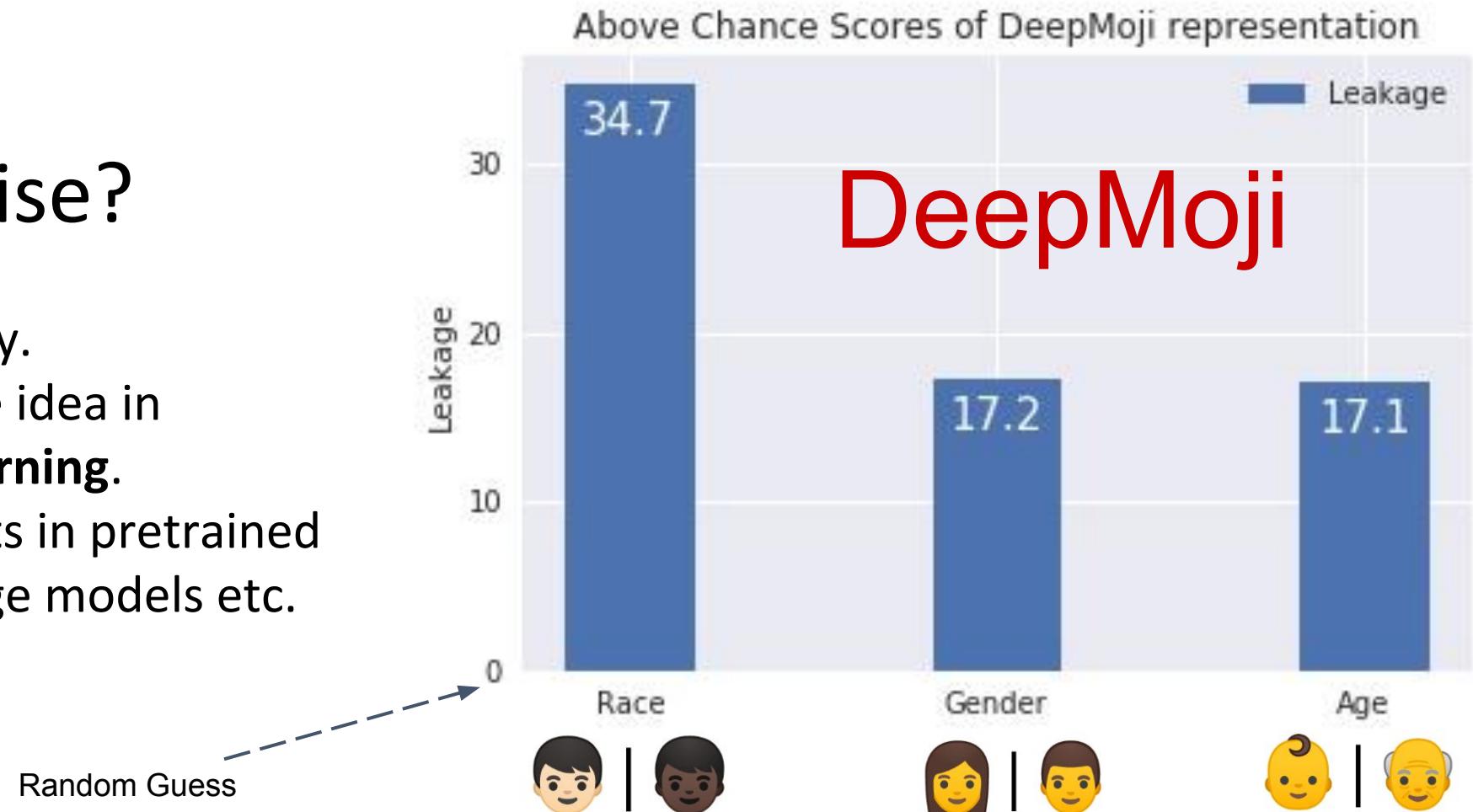
The dev-set scores above chance level are quite high

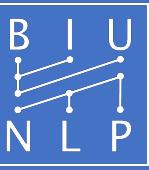
Big Surprise?

Not really.

This is the core idea in **Transfer-Learning**.

We've seen its benefits in pretrained embeddings, language models etc.

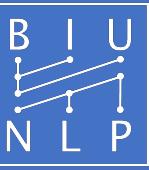




Text Leakage – Case Study

- 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.

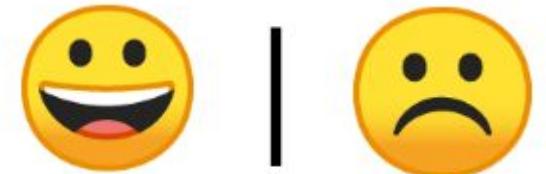
Fair enough. Let's control for it.



Controlled Setup

New 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

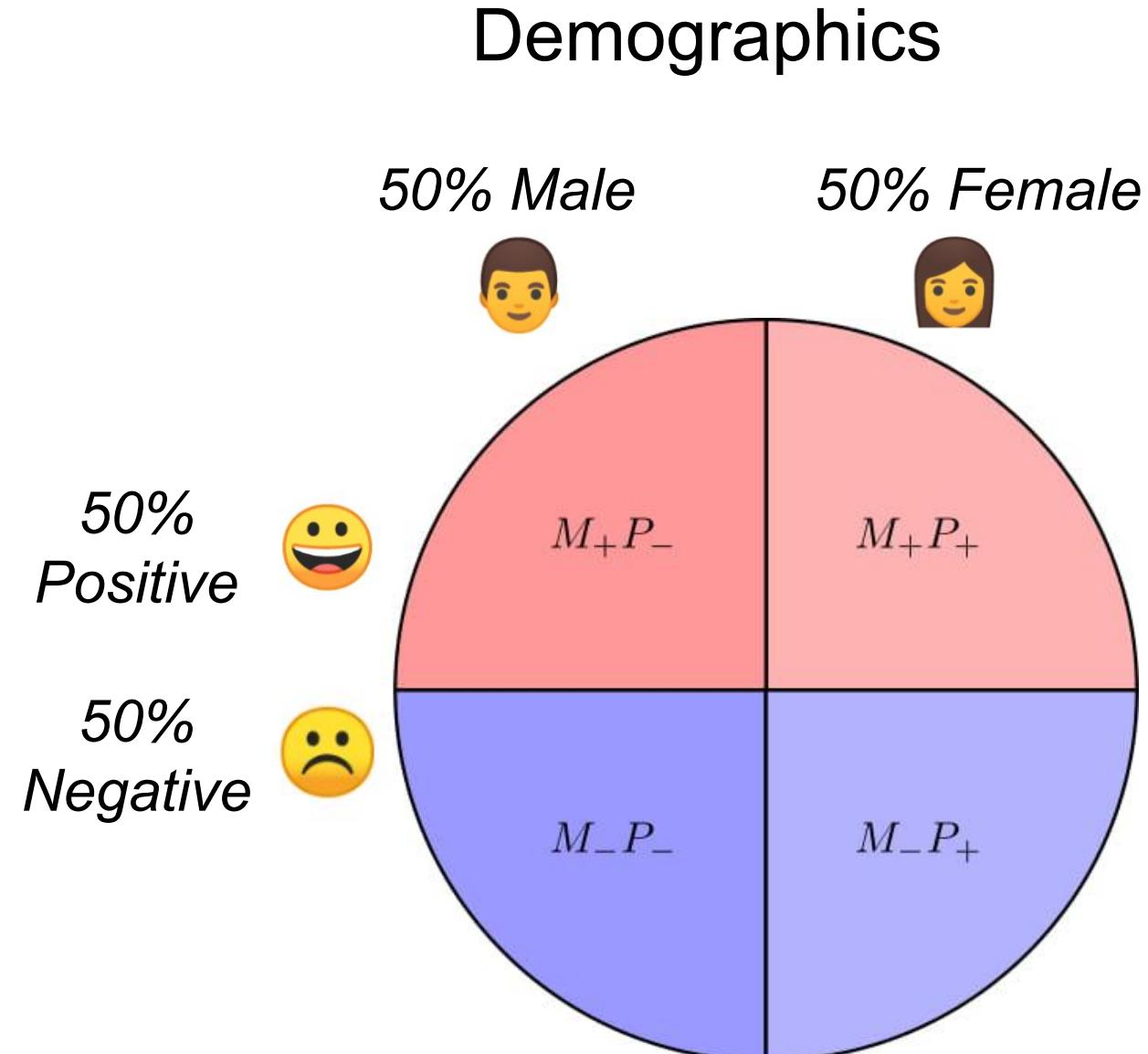
Rangel et al., 2016

New setup

Balanced Dataset

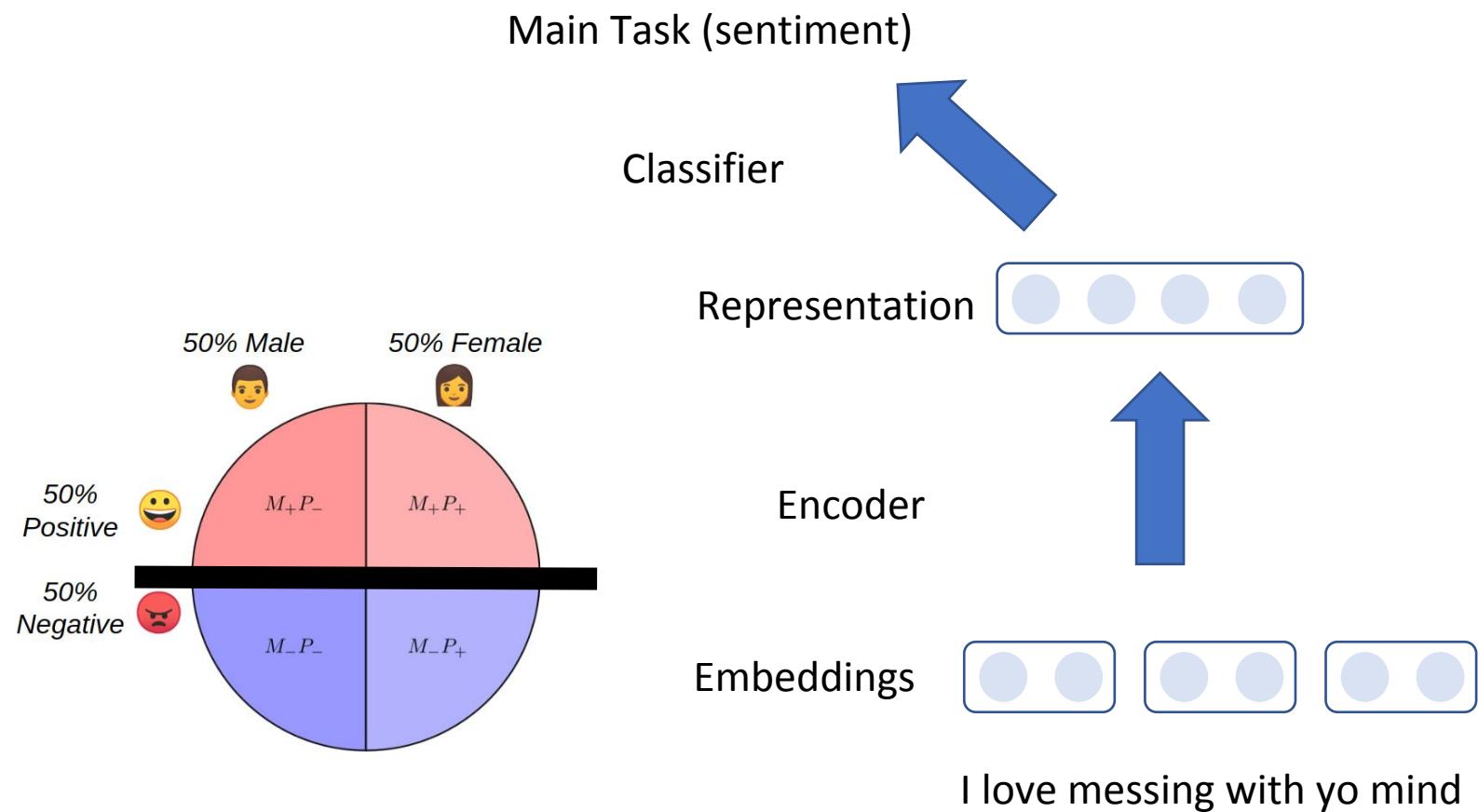
Task
(Sentiment)

50%
Positive
50%
Negative



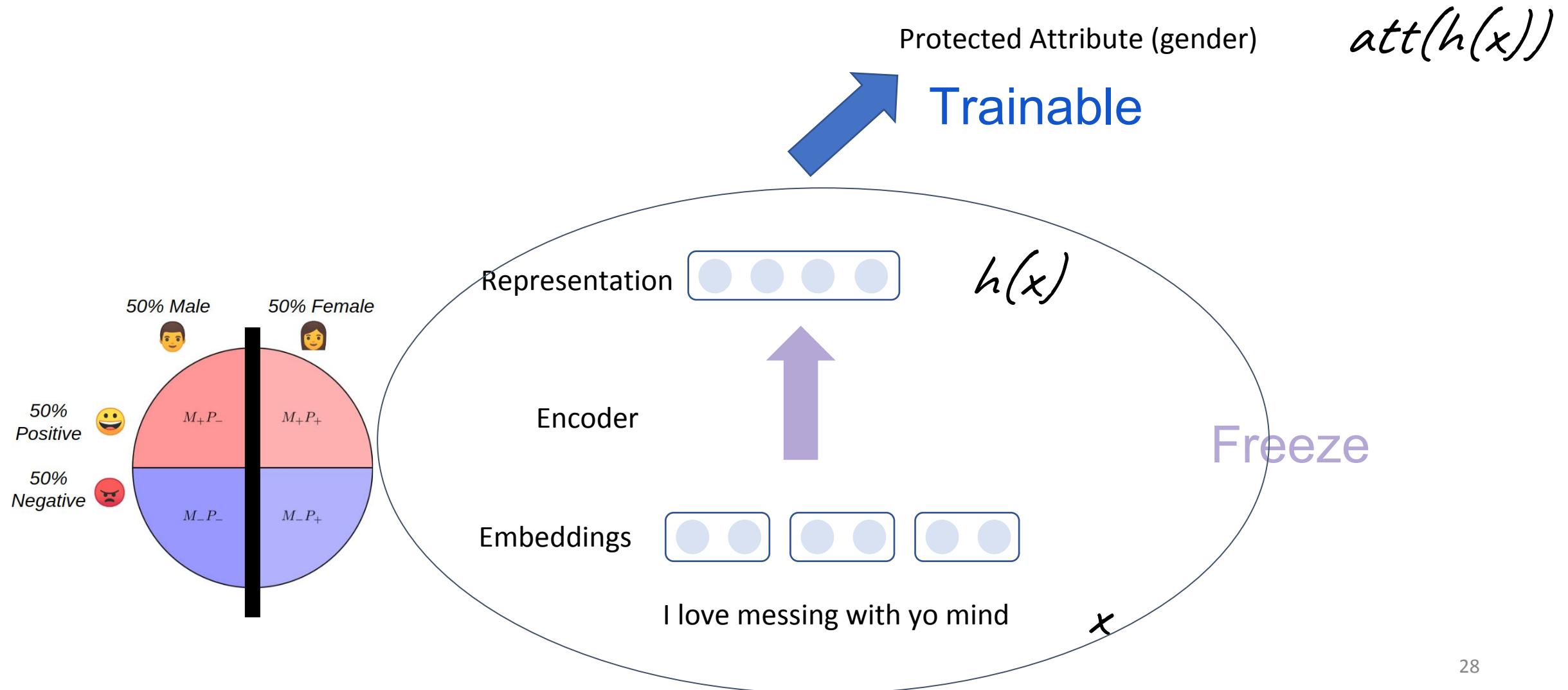
Balanced Training

Training our own encoder on the balanced datasets



Balanced Training

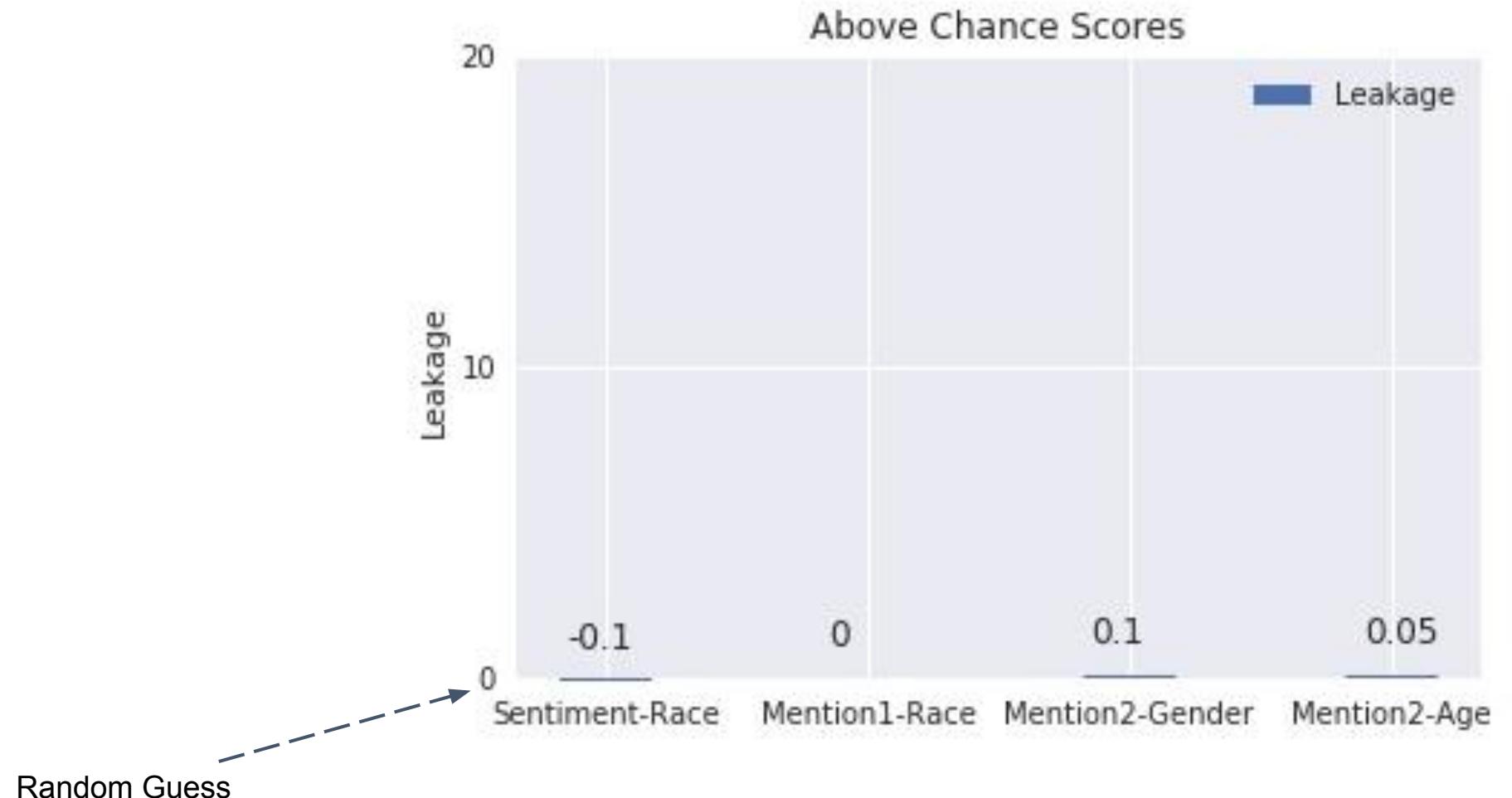
And using the Attacker to check for leakage



Balanced Training - Leakage

We wanted to see something like this:

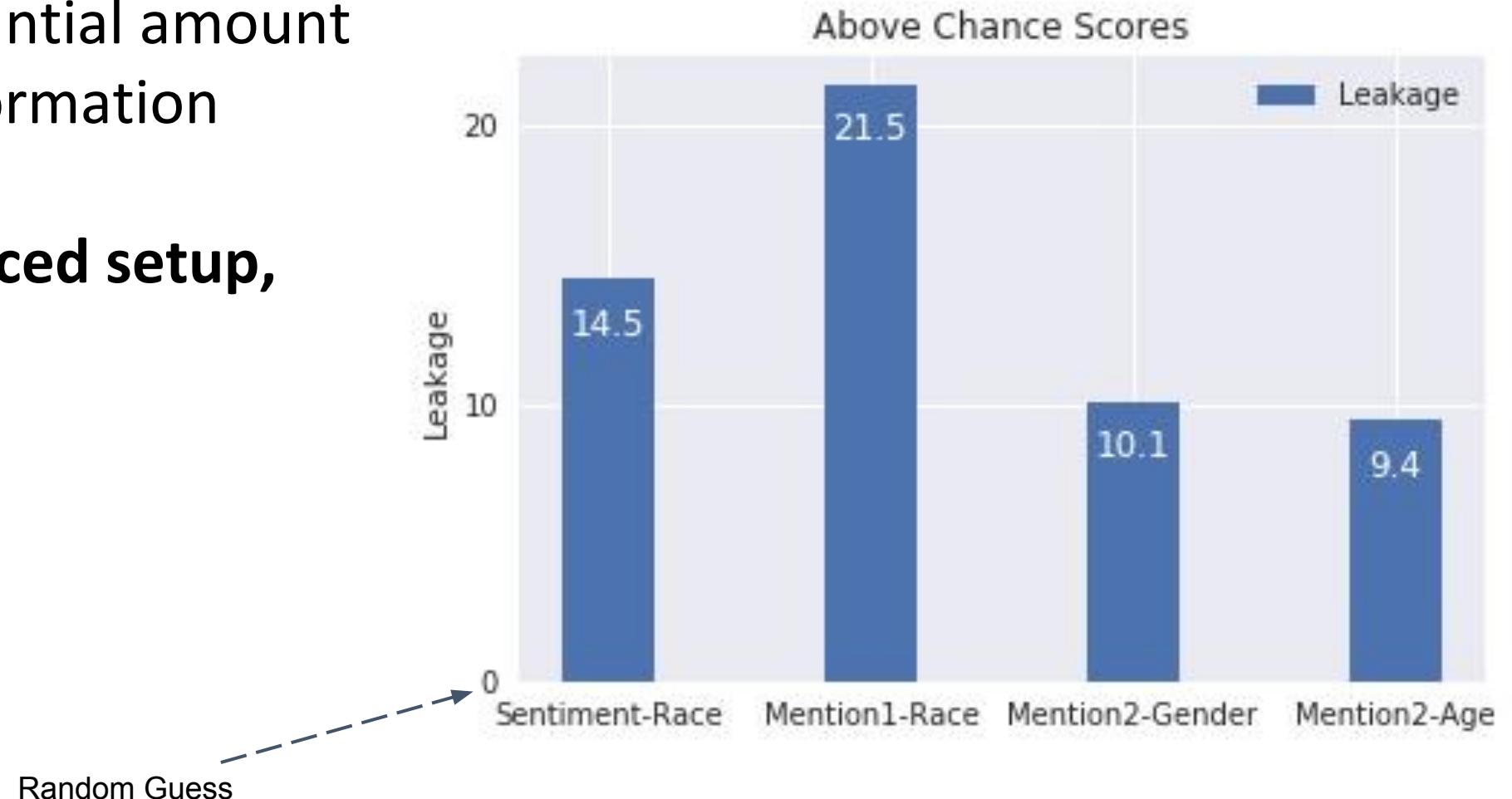
But instead...



Balanced Training - Leakage

The Attacker manages to extract a substantial amount of sensitive information

Even in a balanced setup, leakage exists



Our objective

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



Our objective

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



|



and

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

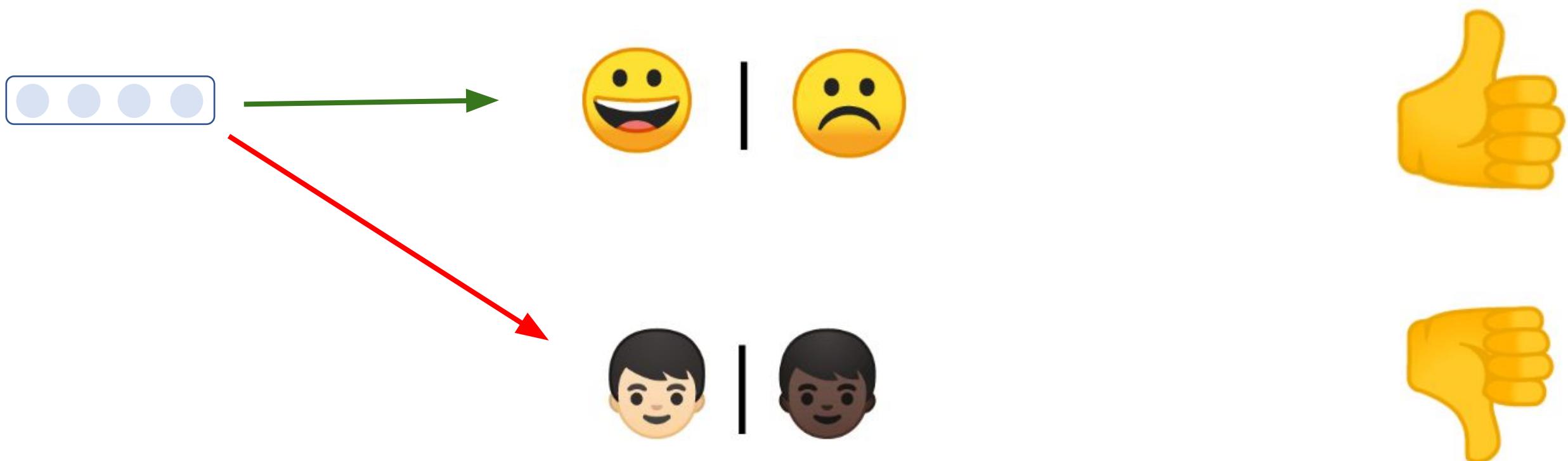


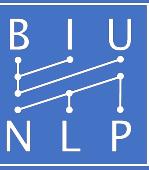
|



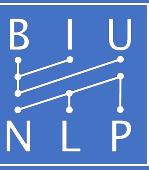
Our objective

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

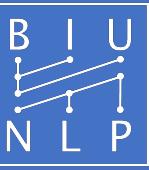
Generative Adversarial Nets

**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

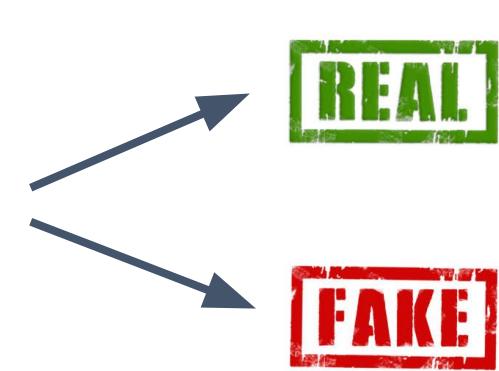


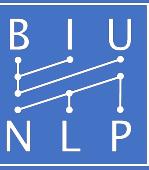
Adversarial Setup

- The motivation came from “Generative Models”
 - We would like to automatically create images
 - From... random input?

Adversarial Setup

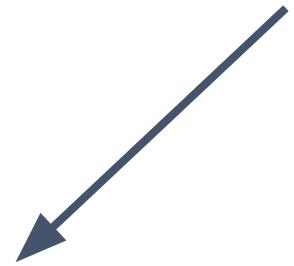
- 2 components:
 - Generator
 - Discriminator



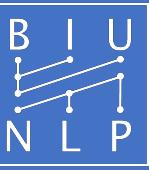


Adversarial Setup

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



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

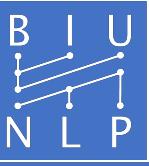


Adversarial Setup

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

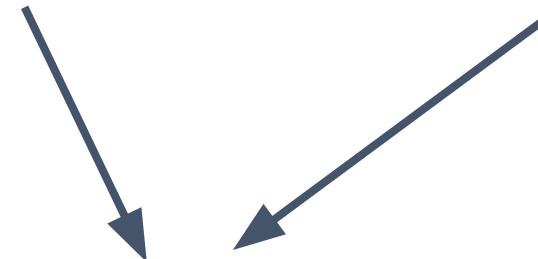


A good Generator
(fake data gets a high score, for
maximizing D 's probability)

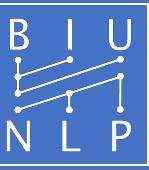


Adversarial Setup

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



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

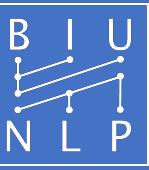


Adversarial Setup

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

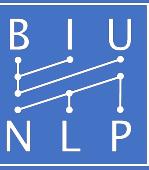
Goodfellow et al. solution:
iterate training between the Generator and Discriminator

- Update the discriminator by ascending its stochastic gradient:
- Update the generator by descending its stochastic gradient:



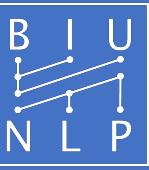
Adversarial Setup

- 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?



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

Victor Lempitsky

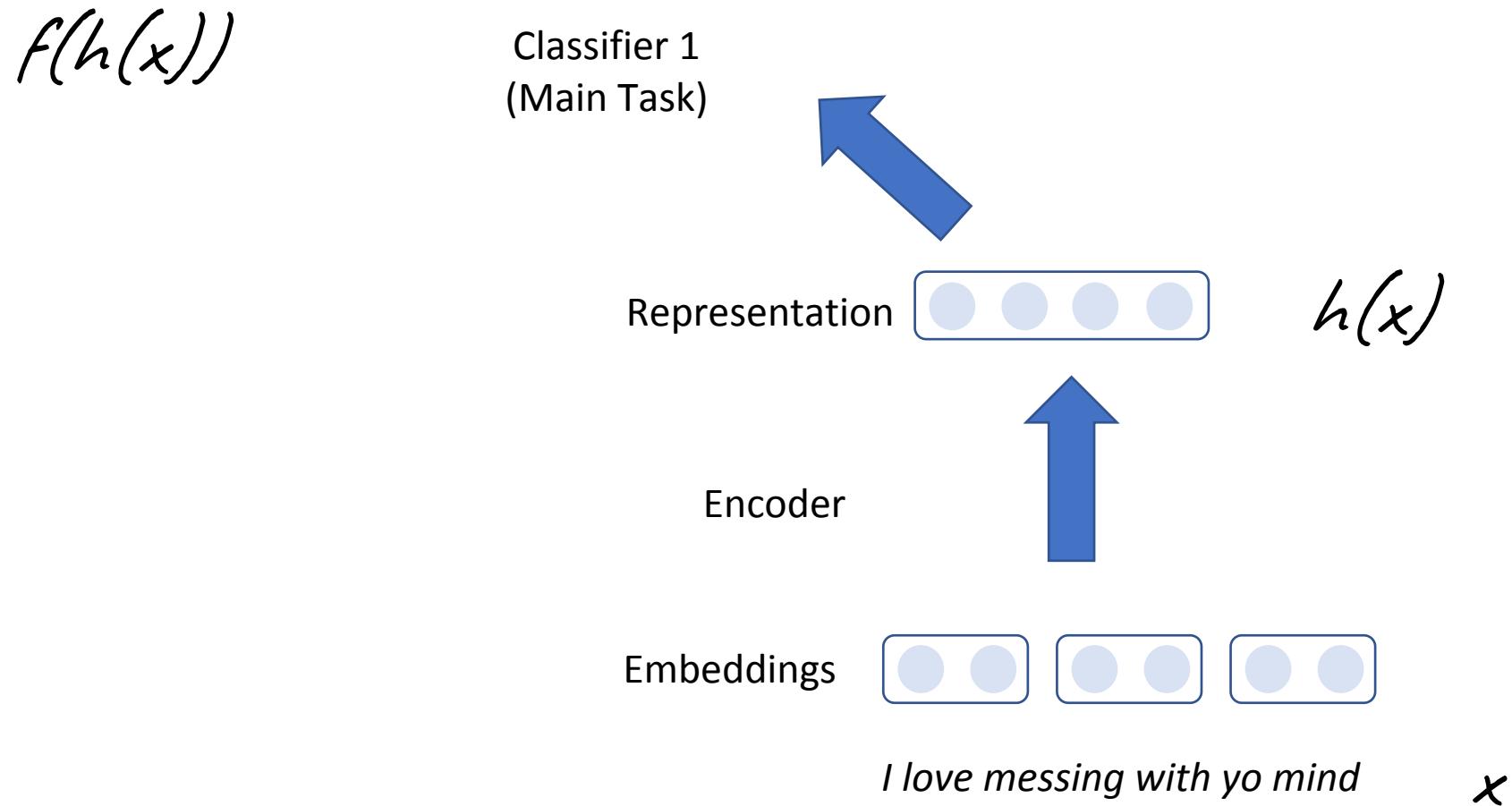
Skolkovo Institute of Science and Technology (Skoltech)

GANIN@SKOLTECH.RU

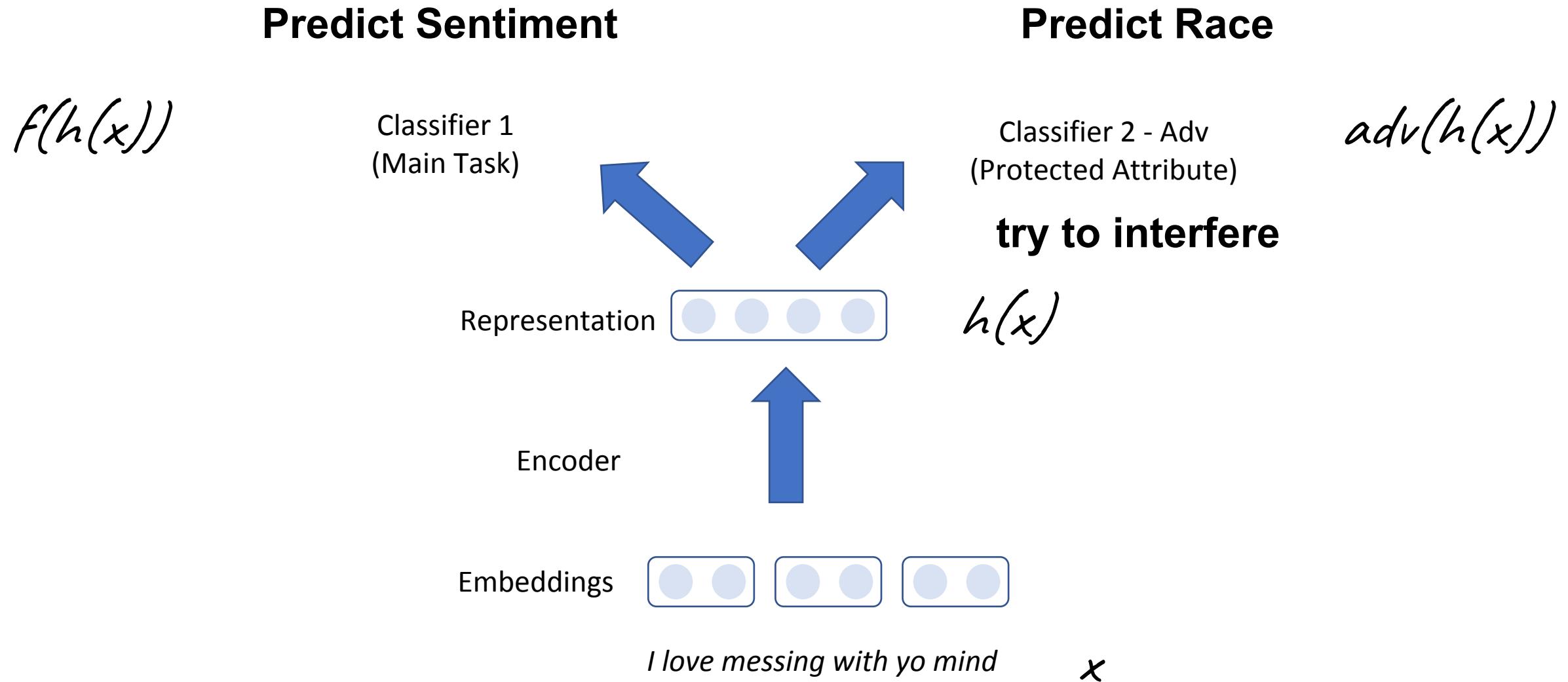
LEMPITSKY@SKOLTECH.RU

Adversarial Setup (Ganin and Lempitsky, 2015)

Predict Sentiment



Adversarial Setup (Ganin and Lempitsky, 2015)

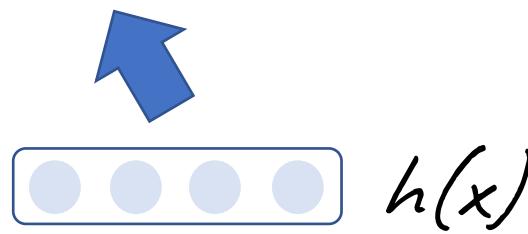


Adversarial Setup (Ganin and Lempitsky, 2015)

3 different sub-objectives

$f(h(x))$

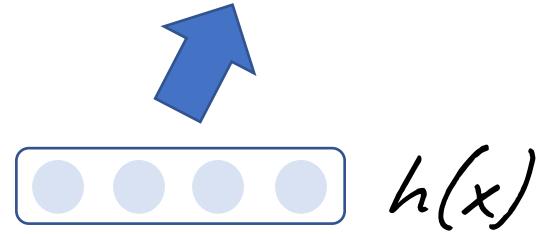
Classifier 1
(Main Task)



classify well

$adv(h(x))$

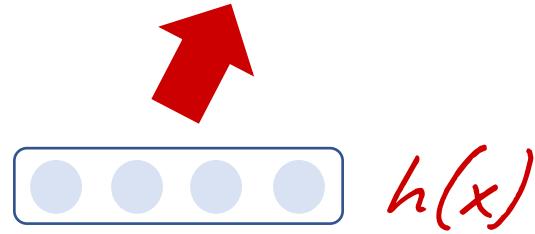
Classifier 2 - Adv
(Protected Attribute)



adversary should succeed

$-adv(h(x))$

Classifier 2 - Adv
(Protected Attribute)



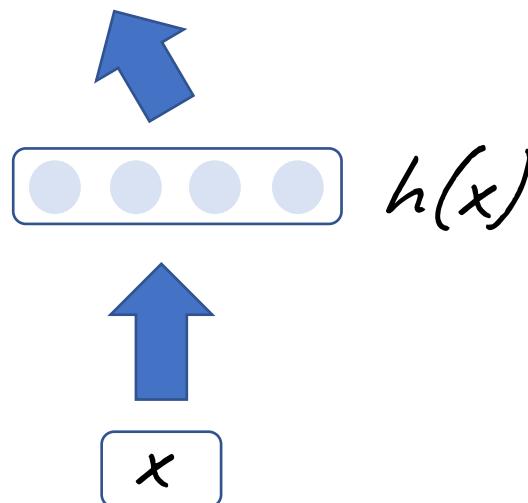
encoder should
make adversary
fail

Adversarial Setup (Ganin and Lempitsky, 2015)

3 different sub-objectives

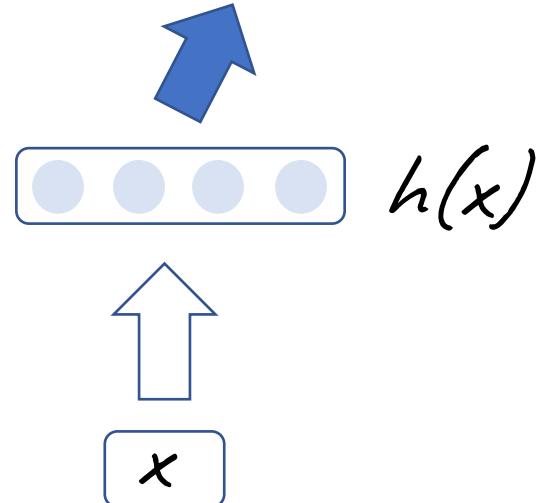
$f(h(x))$

Classifier 1
(Main Task)



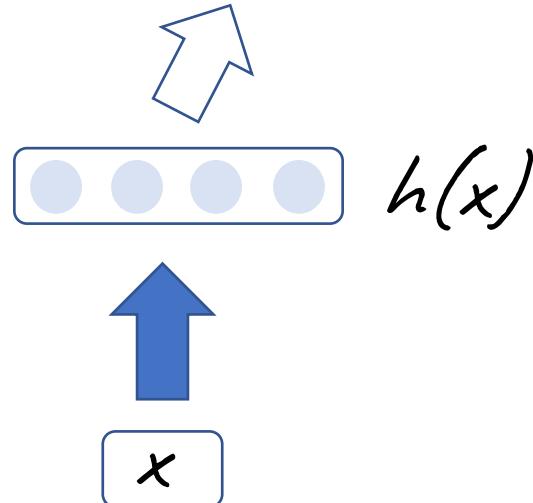
$adv(h(x))$

Classifier 2 - Adv
(Protected Attribute)

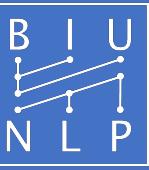


$-adv(h(x))$

Classifier 2 - Adv
(Protected Attribute)



blue: update parameters
white: don't update

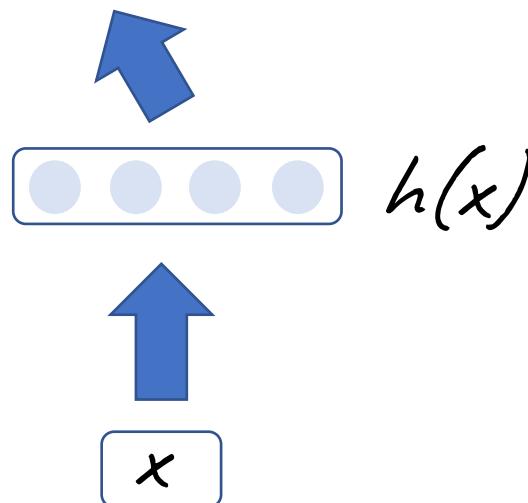


Adversarial Setup (Ganin and Lempitsky, 2015)

3 different sub-objectives

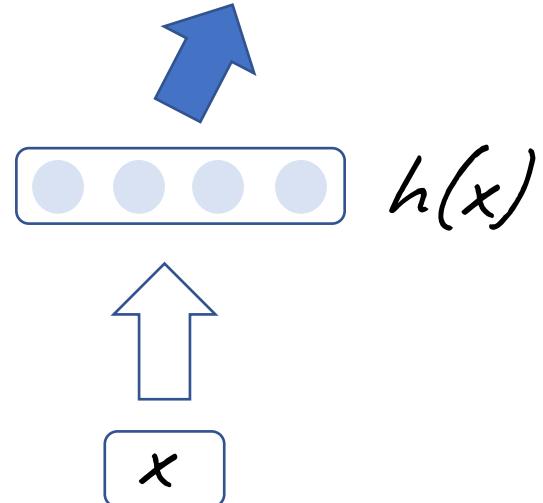
$$f(h(x))$$

Classifier 1
(Main Task)



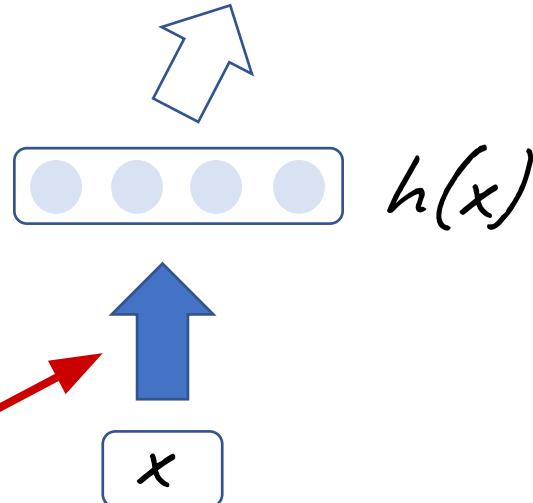
$$\text{adv}(h(x))$$

Classifier 2 - Adv
(Protected Attribute)



$$-\text{adv}(h(x))$$

Classifier 2 - Adv
(Protected Attribute)



blue: update parameters
white: don't update

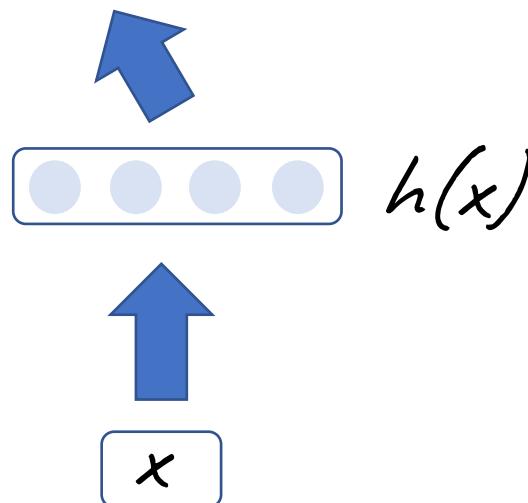
$$\text{grad}(-\text{adv}(h(x)))$$

Adversarial Setup (Ganin and Lempitsky, 2015)

3 different sub-objectives

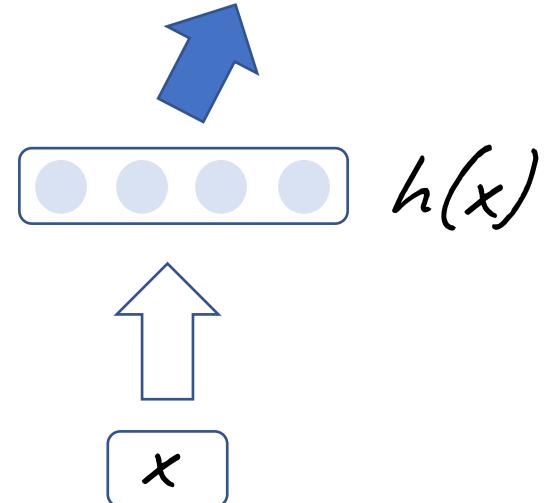
$$f(h(x))$$

Classifier 1
(Main Task)



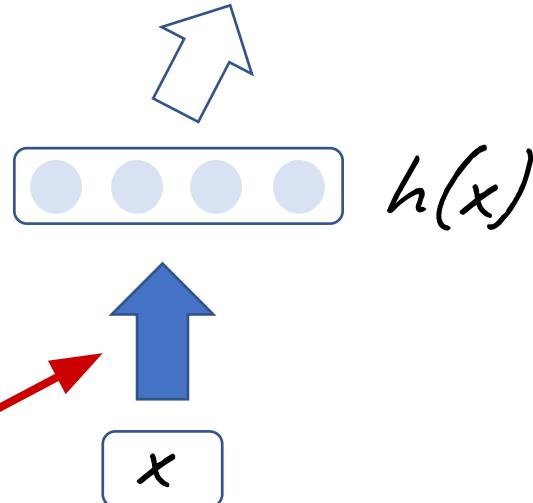
$$\text{adv}(h(x))$$

Classifier 2 - Adv
(Protected Attribute)



$$-\text{adv}(h(x))$$

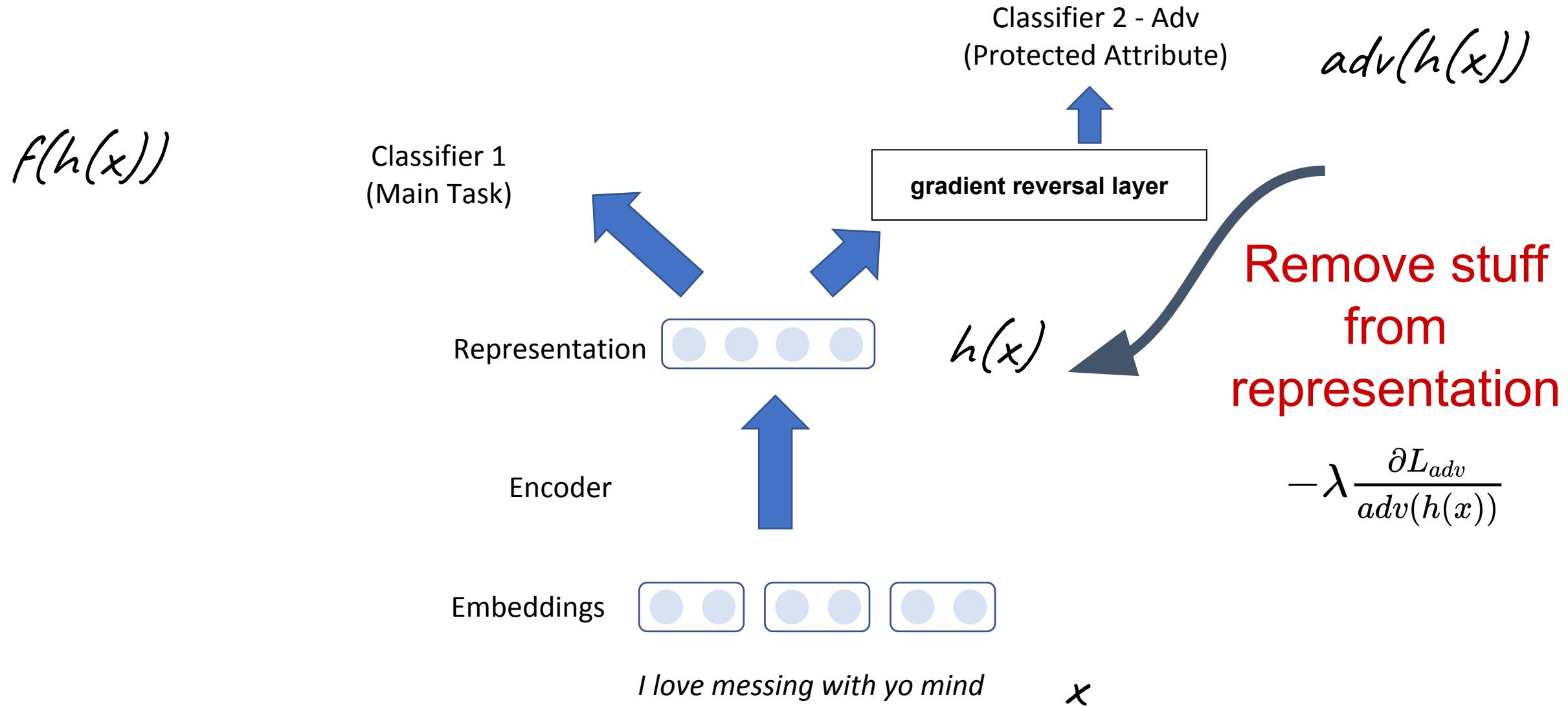
Classifier 2 - Adv
(Protected Attribute)



blue: update parameters
white: don't update

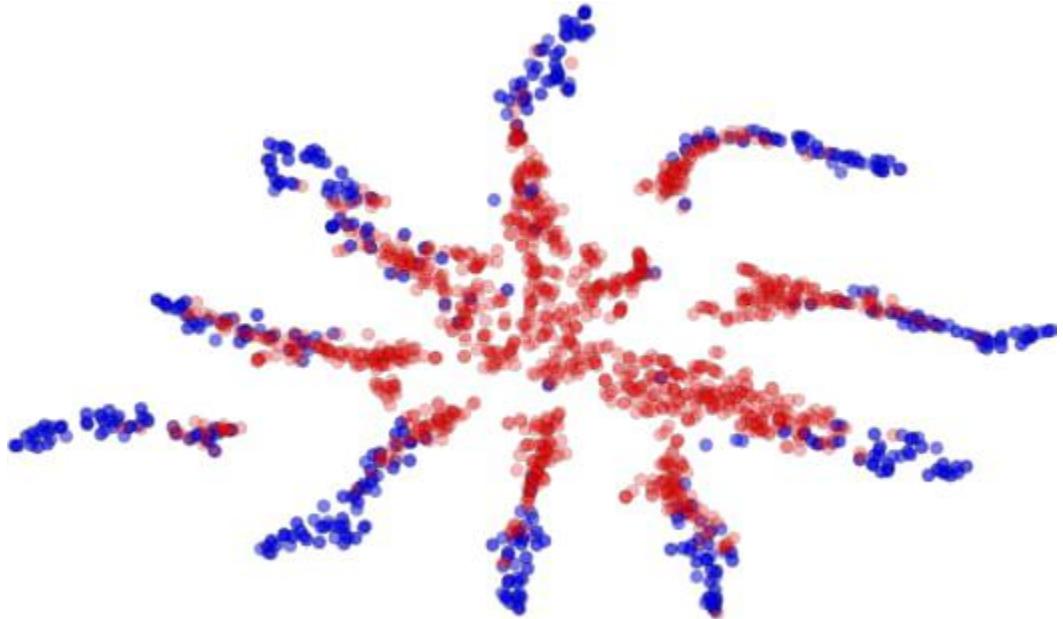
$$\text{grad}(-\text{adv}(h(x))) = -\text{grad}(\text{adv}(h(x)))$$

Adversarial Setup (Ganin and Lempitsky, 2015)

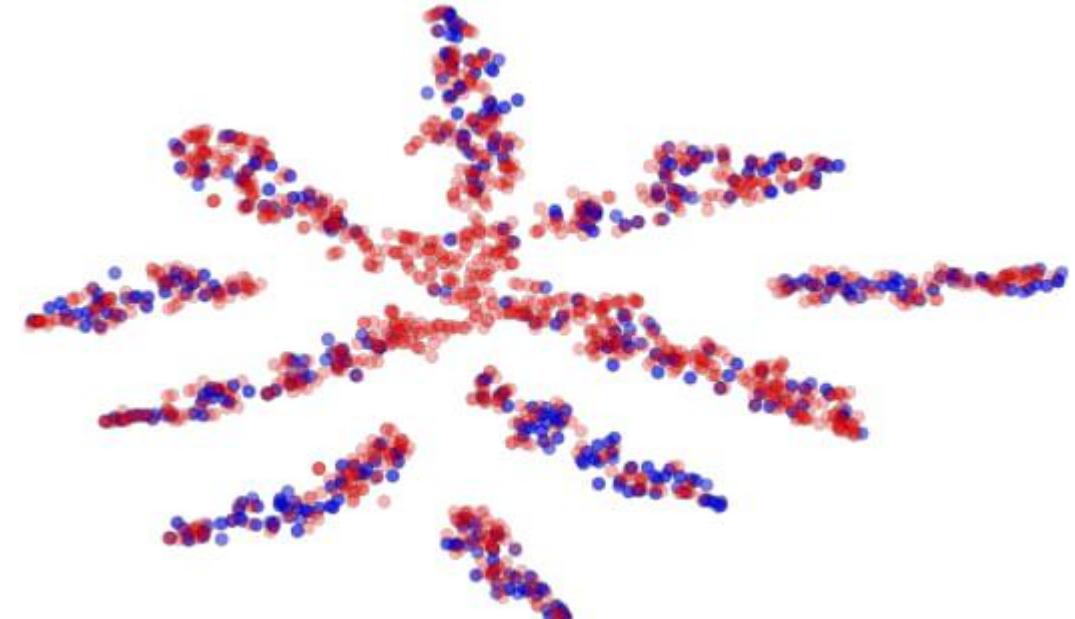


Adversarial Setup (Ganin and Lempitsky, 2015)

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

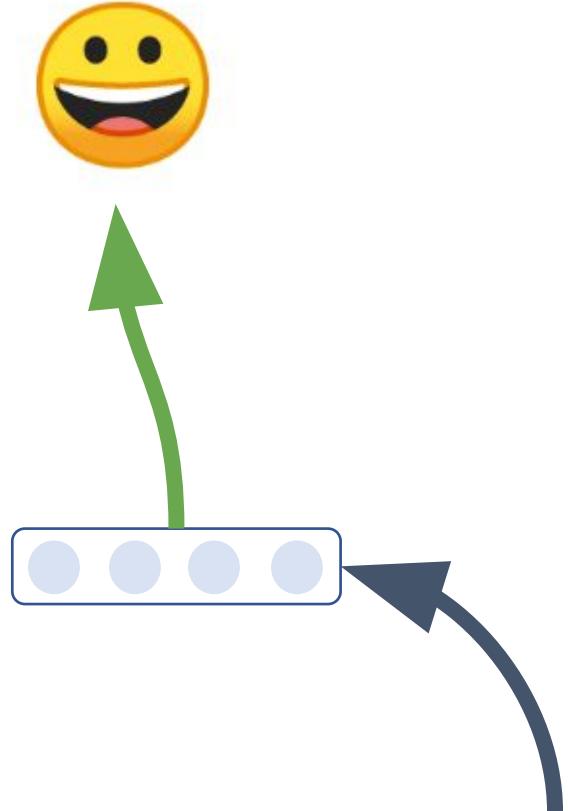


before



after

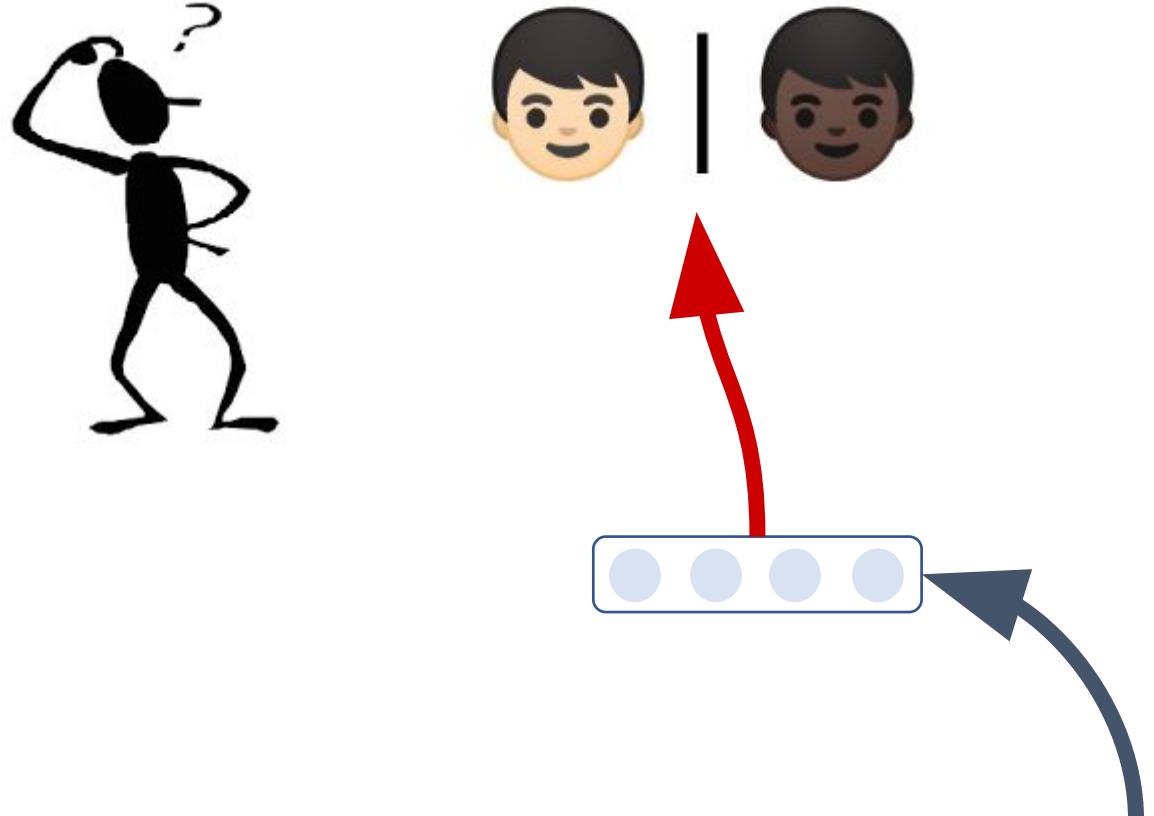
Does it work?



**Successfully predicting
sentiment**

"I love mom's cooking"

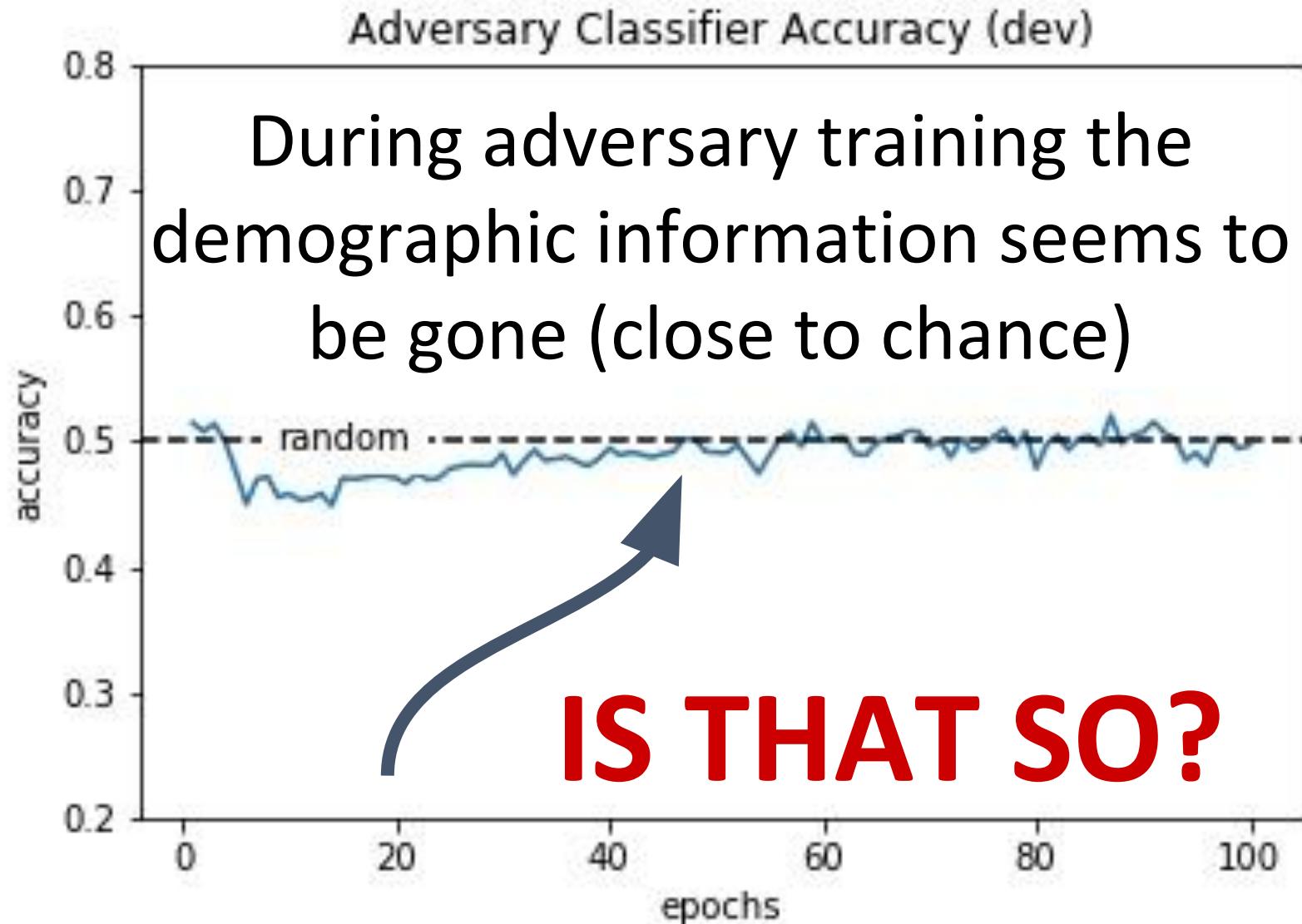
Does it work?



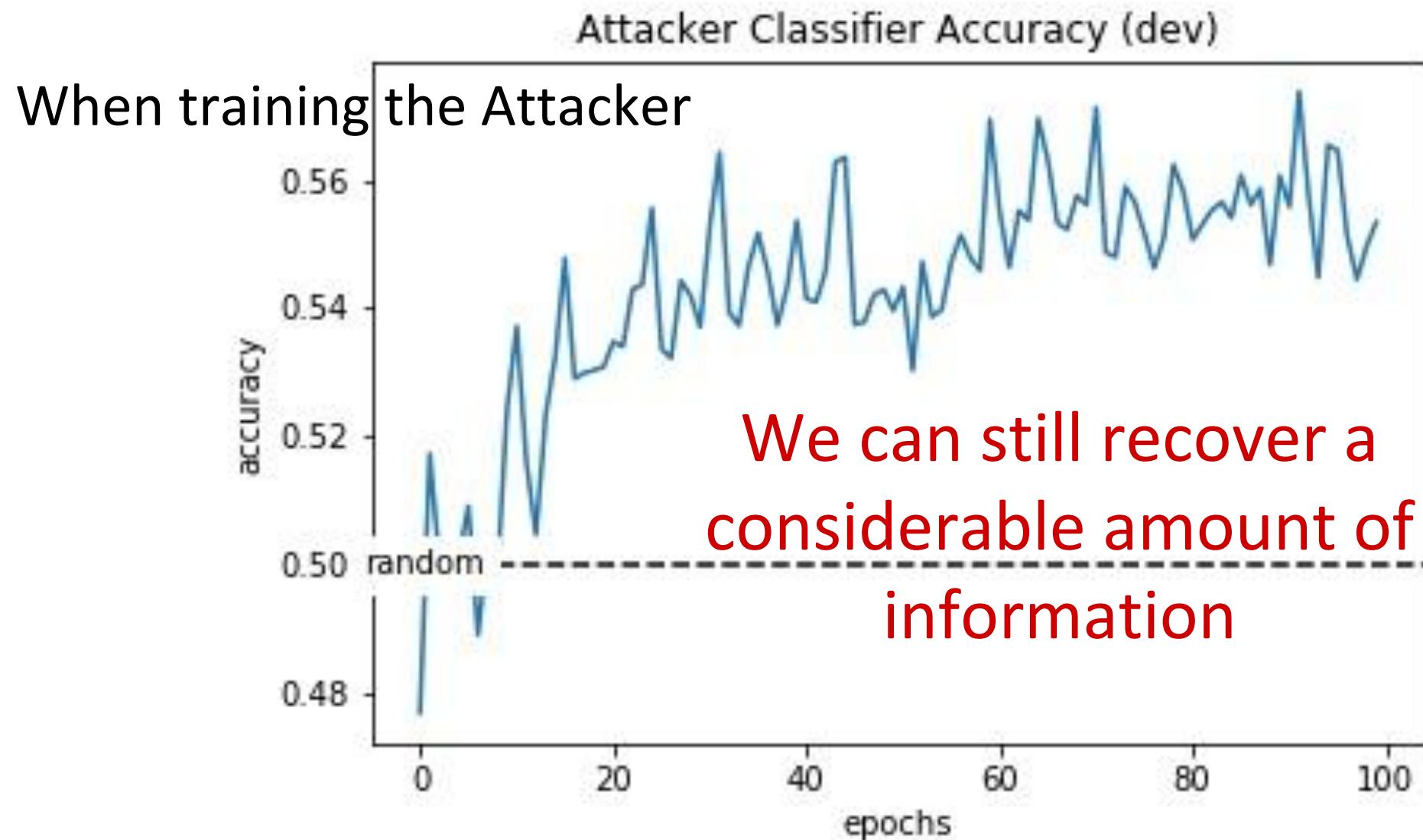
**Successfully removed
demographics?**

“I love mom’s cooking”

Does it work?



Does it work? Not so quickly...



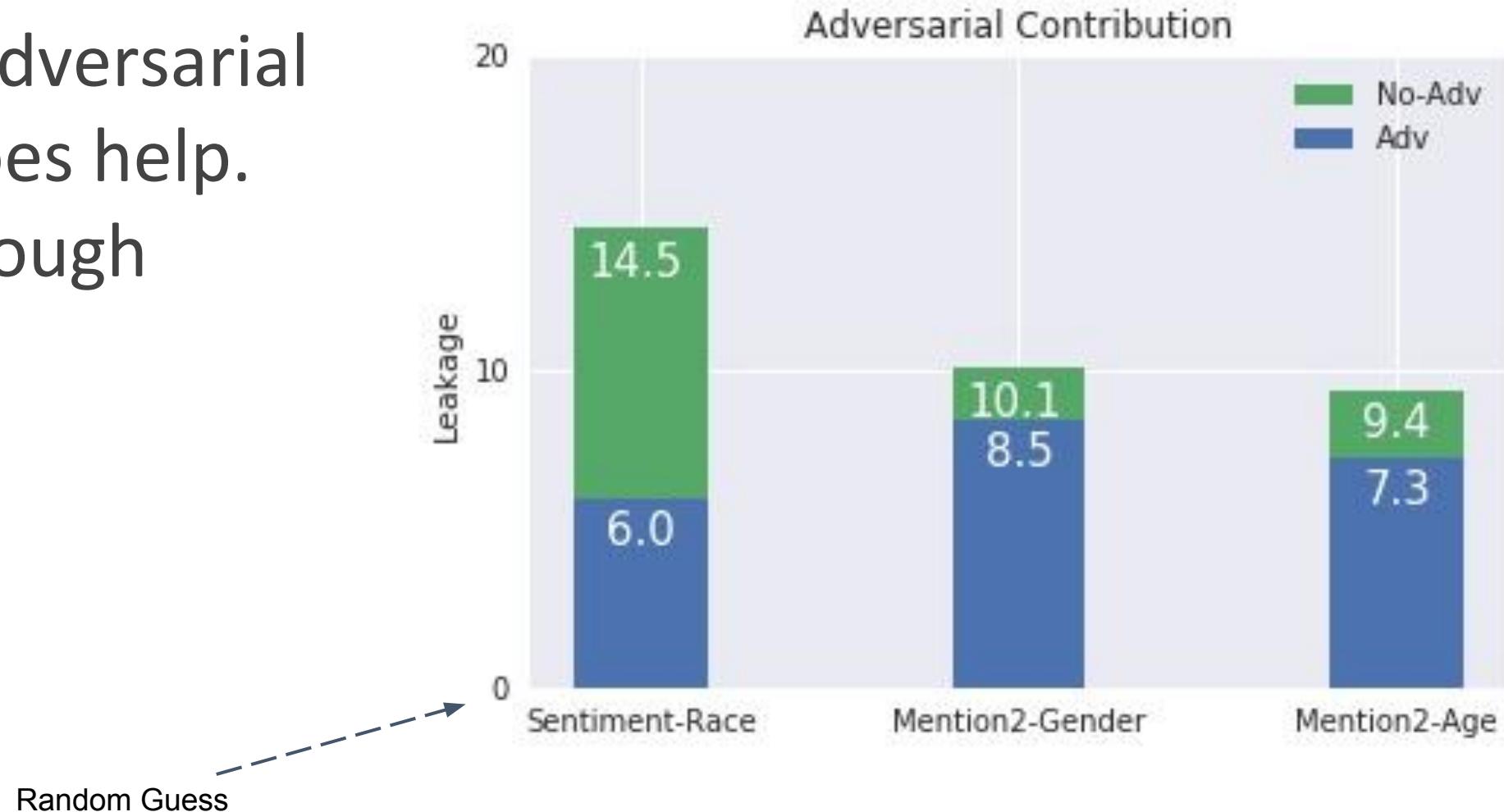
Does it work? Not so quickly...

Consistent across tasks and protected attributes

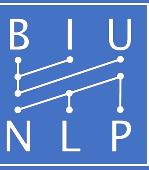


Does it work? more or less

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



Stronger, Better, Bigger???

Can we make stronger
adversaries?

Stronger, Better, Bigger???

More Baselineers!

$f(h(x))$

Classifier 1
(Main Task)

Representation

Encoder

Embeddings



I love messing with yo mind

Classifier 2 - Adv
(Protected Attribute)

gradient reversal layer

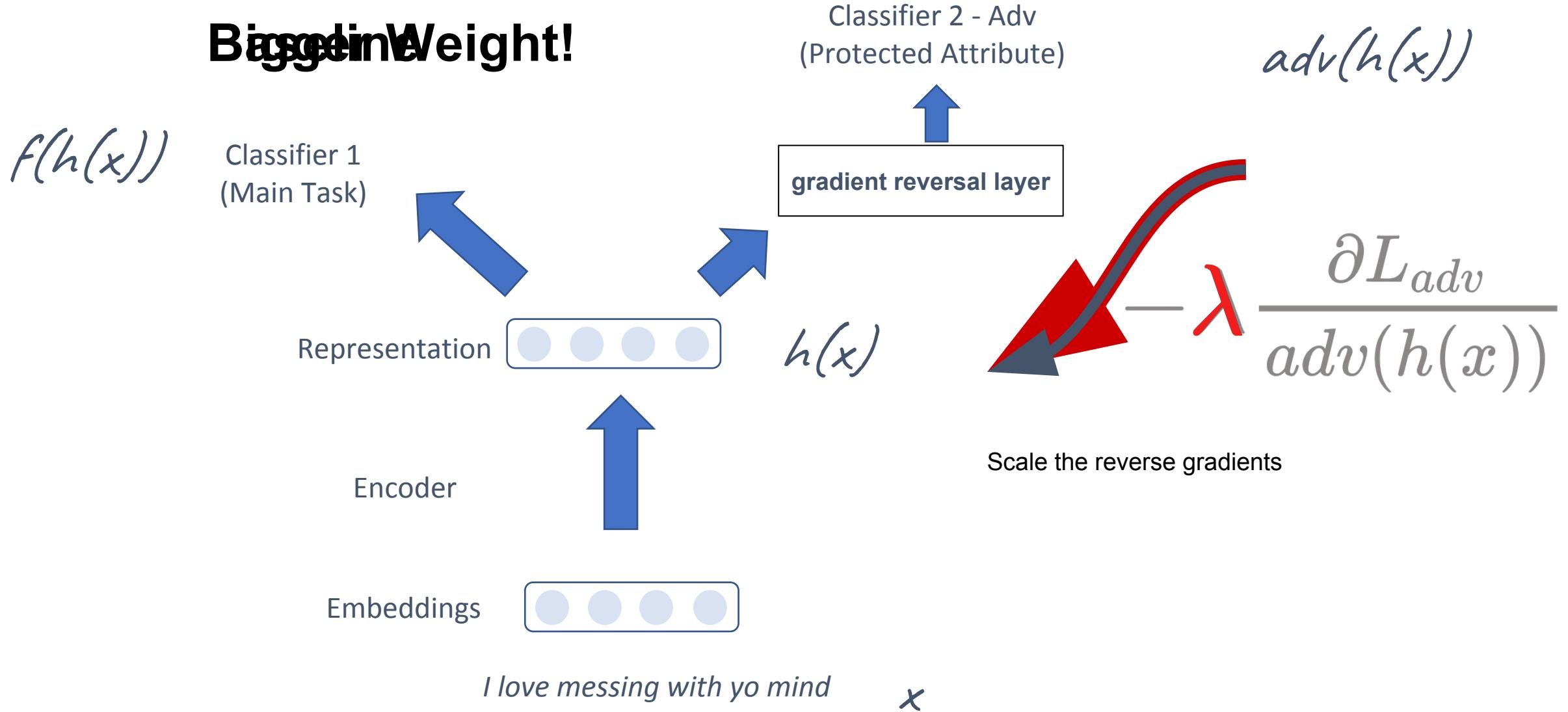
$h(x)$

x

$adv(h(x))$

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

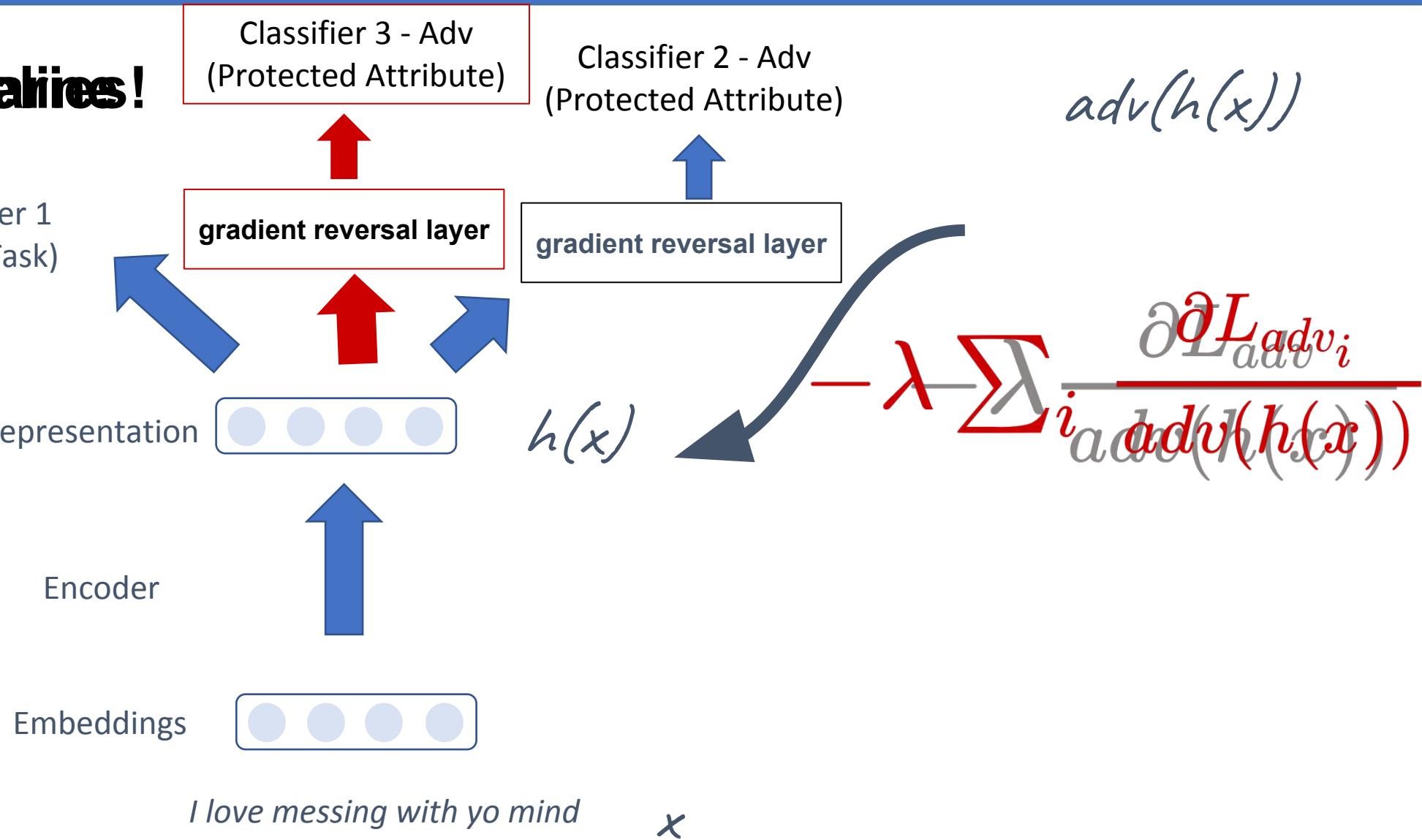
Stronger, Better, Bigger???



Stronger, Better, Bigger???

More Adversaries!

$f(h(x))$

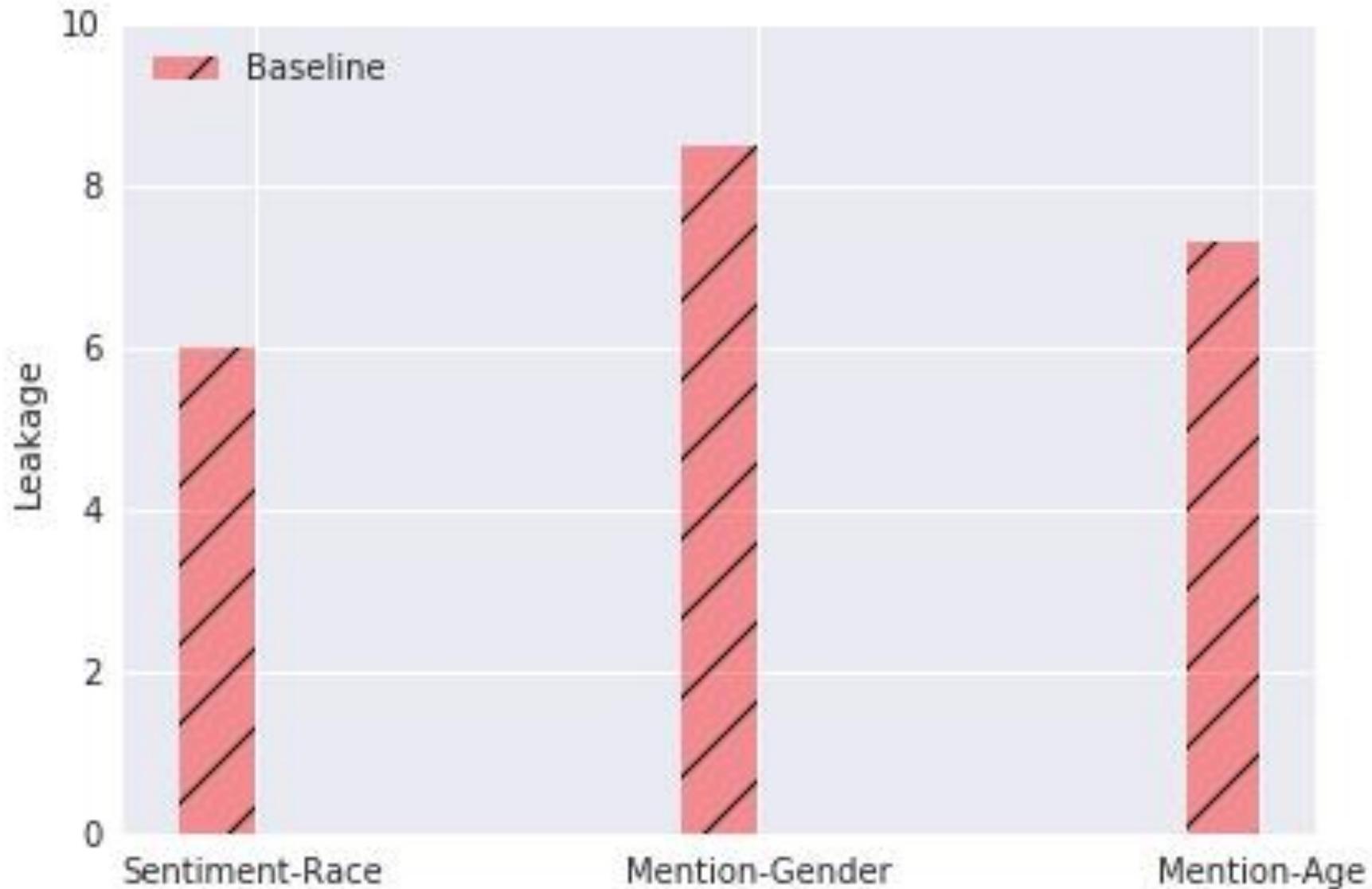


I love messing with yo mind

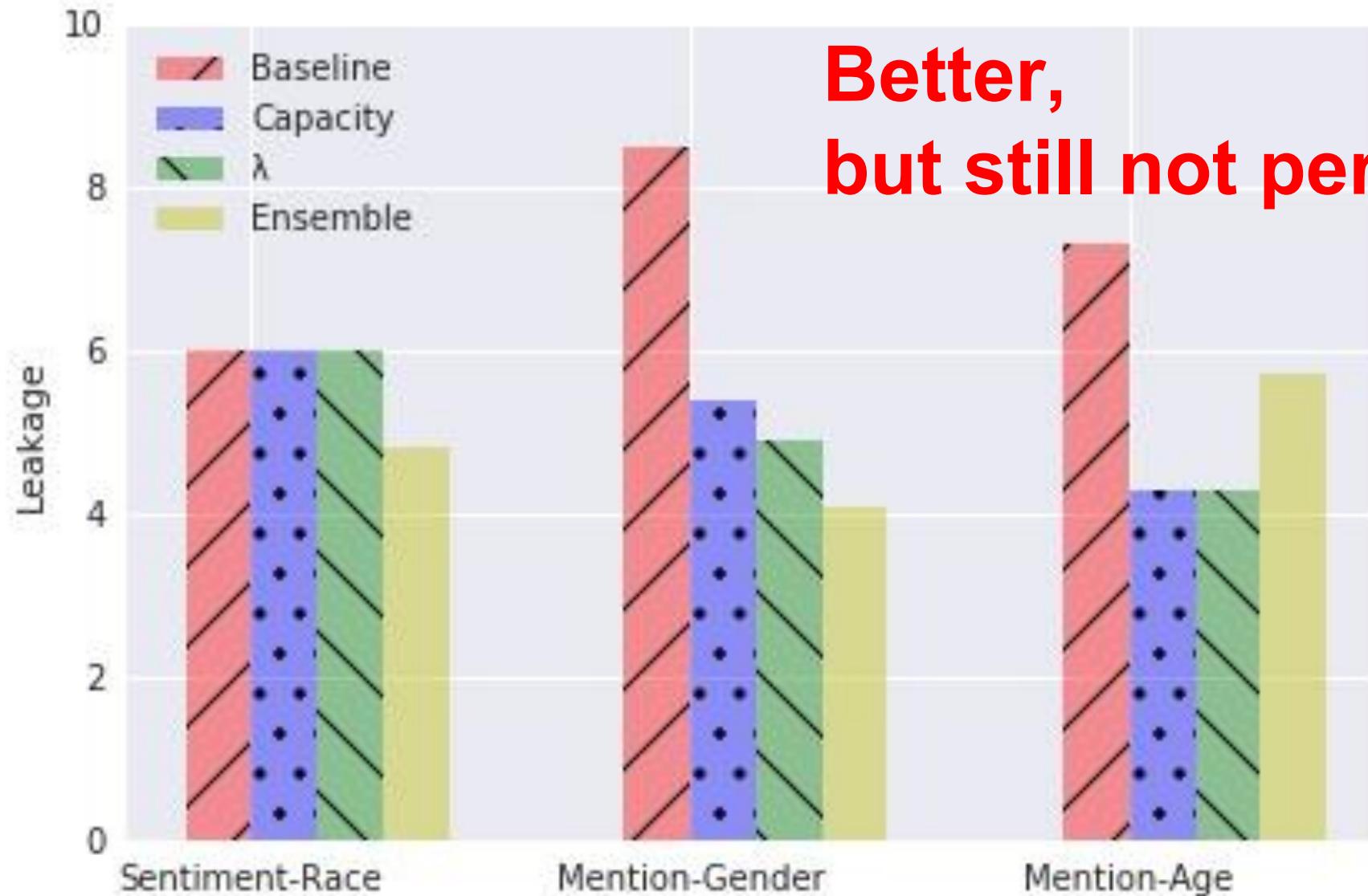
x

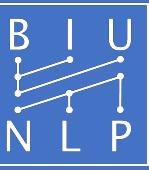


Stronger, Better, Bigger???

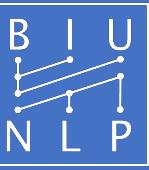


Stronger, Better, Bigger???



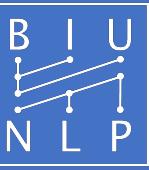


Error Analysis



Wait. I remember this thing called Overfitting

- 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?

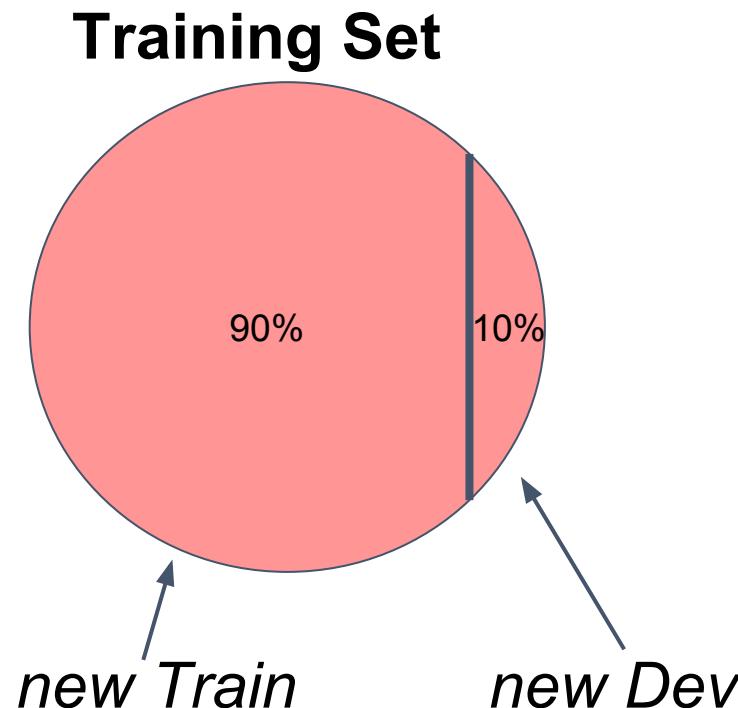


Wait. I remember this thing called Overfitting

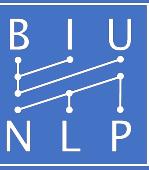
- “Adversary overfitting”:
 - Memorizing the training data
 - By removing all its sensitive information
 - While leaking in test time

Wait. I remember this thing called Overfitting

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

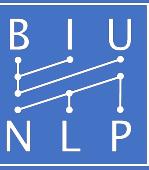


It is more than that



Persistent Examples

- 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.

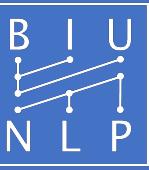
More about the leakage origin can be found in the paper

Few words about fairness

- 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



What happens in recommendation systems?

Privacy and Fairness in Recommender Systems via Adversarial Training of User Representations

Yehezkel S. Resheff ¹, Yanai Elazar ^{2,3}, Moni Shahar ¹, and Oren Sar Shalom ³

¹*Intuit Tech Futures, Israel*

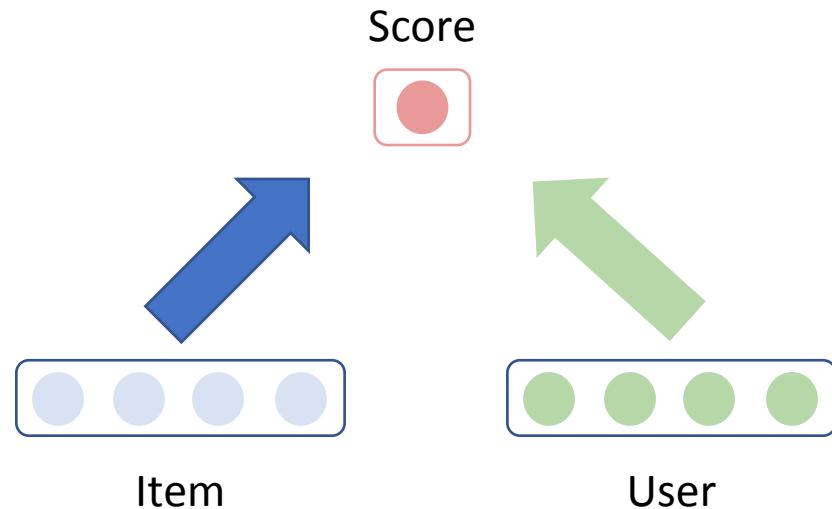
²*Bar Ilan University, Israel*

³*Intuit, Israel*

{*hezi.resheff, yanaiela, monishahar, oren.sarshalom*}@gmail.com

What happens in recommendation systems?

An Abstractive overview...

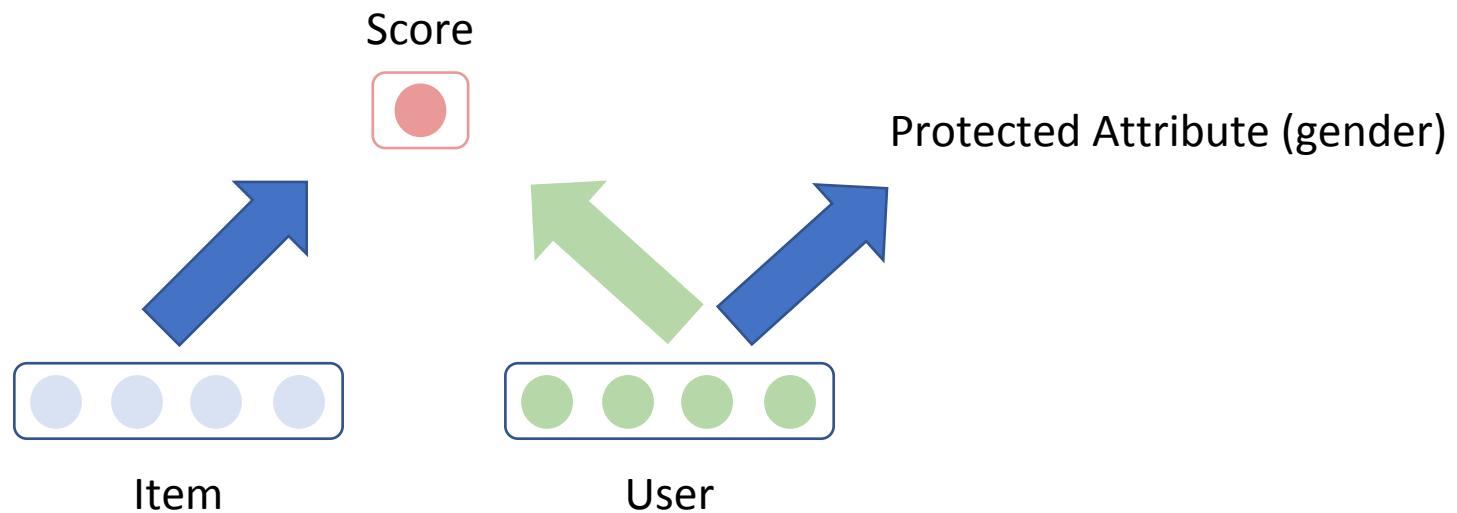


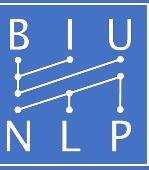
We train it such:

- Item-User pair which occurred get a higher score than the ones who did not

What happens in recommendation systems?

When adding the adversarial component:





What happens in recommendation systems?

Claiming that a representation is not leaky is hard

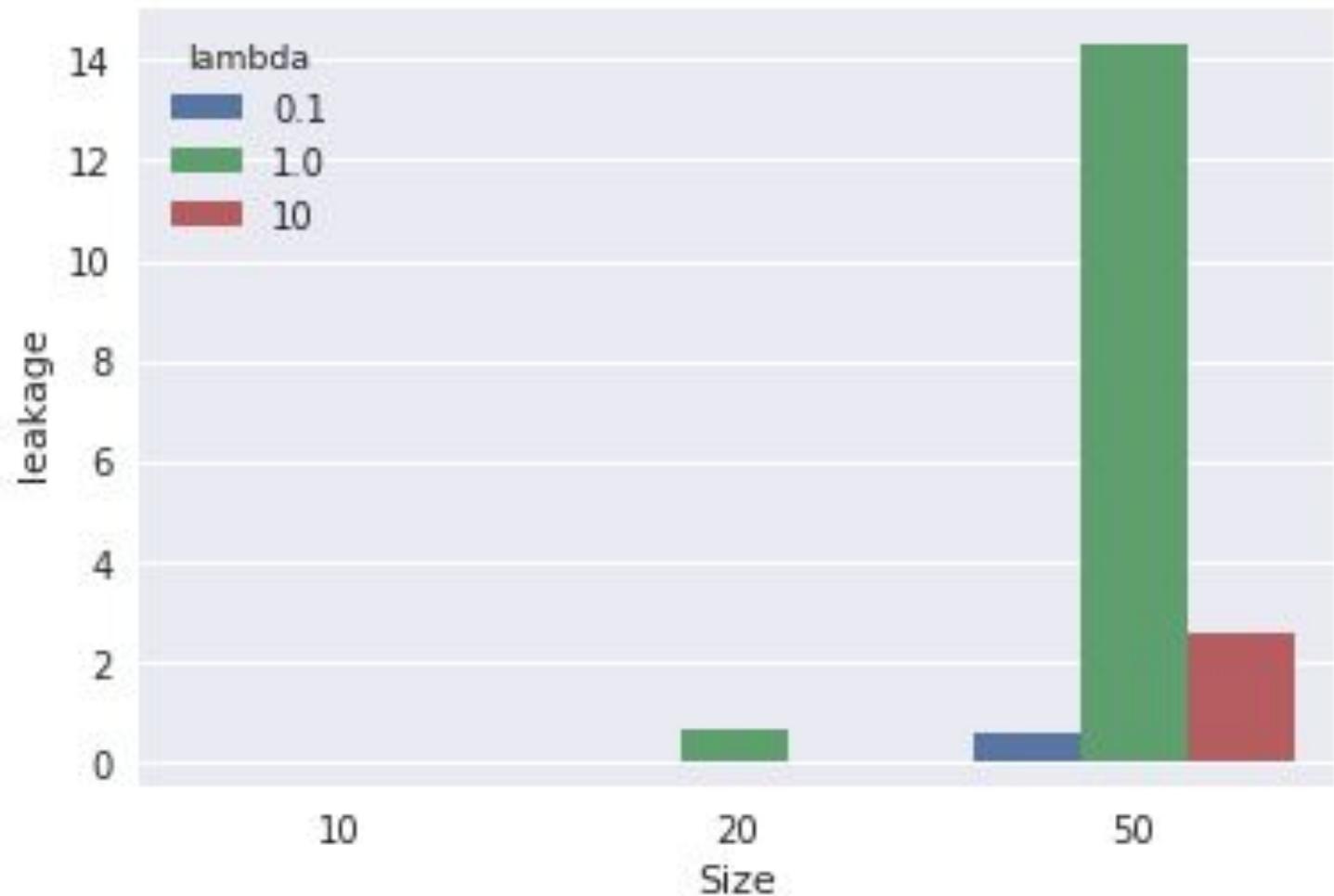
We tested for leakage with several classifiers:

- NN
- SVM
- Decision Tree
- Random Forest
- Gradient Boosting

What happens in recommendation systems?

Amount of leakage:

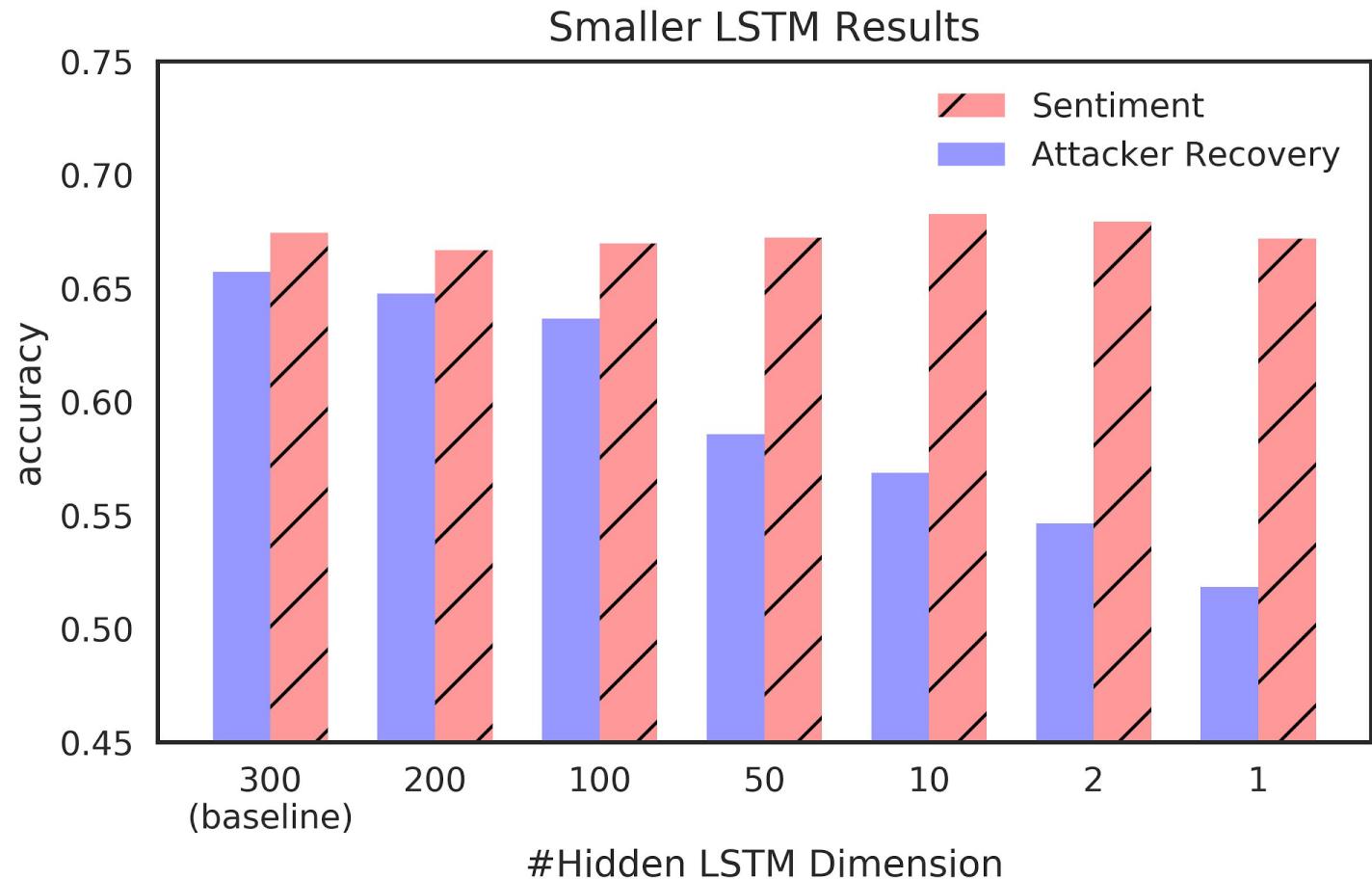
- Highly sensitive to the hyperparameters!
- Lower leakage with smaller dimensions!

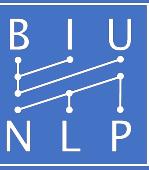


What happens in recommendation systems?

Amount of leakage:

- Highly sensitive to the hyperparameters!
- Lower leakage with smaller dimensions!
- Also in Texts!



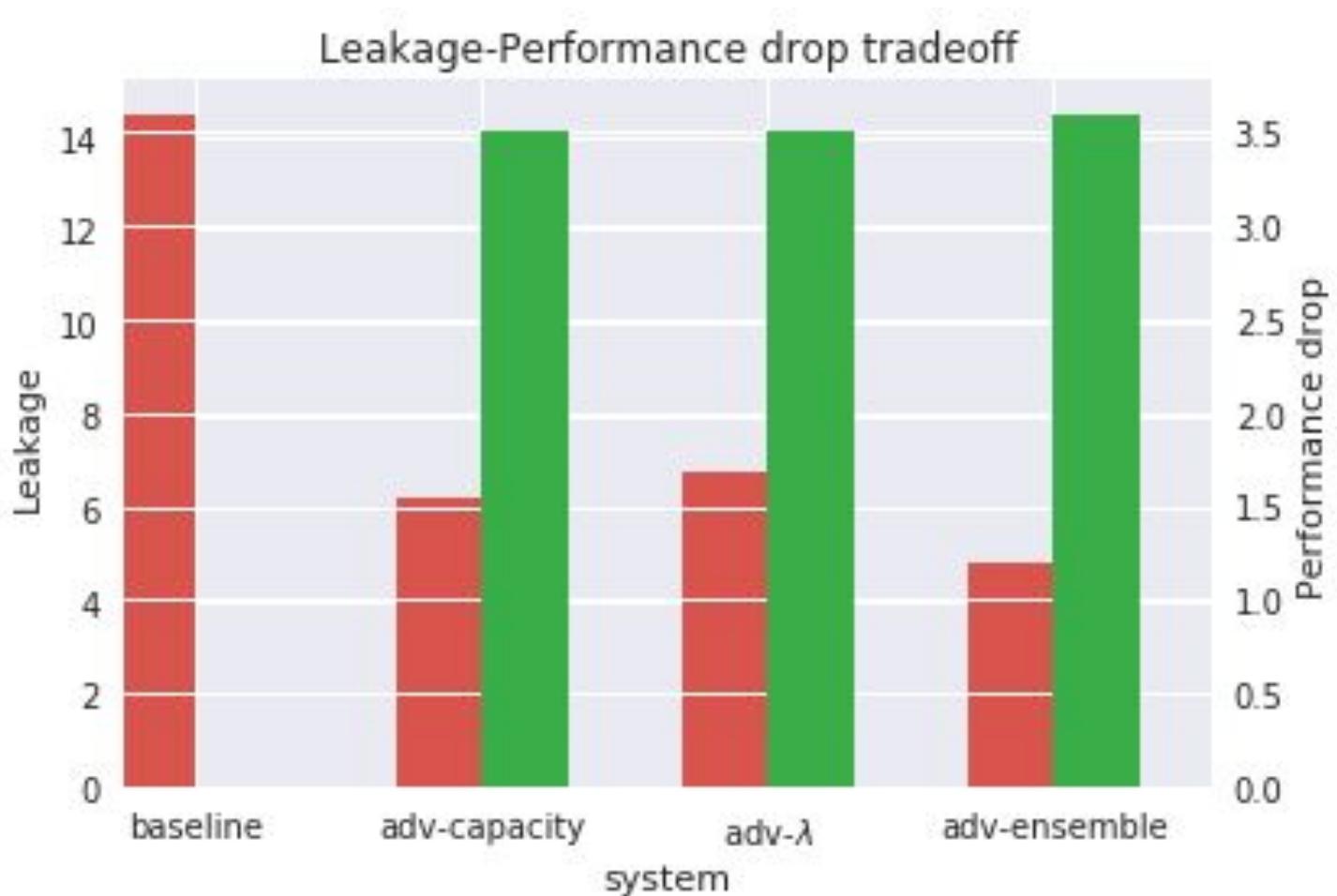


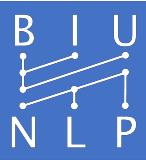
What happens in recommendation systems?

- A small representation size might not provide satisfying performance
- But tuning can help reduce leakage!

Bias is important... but what about performance?

- Performance drop is consistent across tasks
- Need to find the tradeoff
- But...
 - The numbers are not everything!





Bias is important... but what about performance?



Vered Shwartz

March 8, 2018 ·

...

Dear advertisers, happy women's day! Could you please stop targeting me with ads about babies? I know my gender, age, and marital status fit your profile, but trust me, you're wasting your time. And frankly, I'm a bit offended that you think this is all I can care about, while my husband sees ads on many other cooler and more relevant topics.



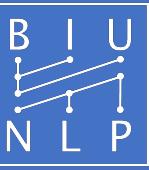
and 41 others

3 Comments

Like

Comment

Share



Bias is important... but what about performance?

- Look at the data
- Make user studies
- Find your trade-off

Summary

- 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”**
 - **Tune your models for the best trade-off**

Thank you