

# Linear Adversarial Concept Erasure

Shauli Ravfogel, Michael Twiton, Yoav Goldberg and Ryan Cotterell

**ETH** zürich

Bar-Ilan  
University

אוניברסיטת בר-אילן



**Ai2**

# Motivation

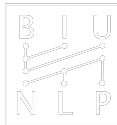
Neural models learn rich representations

# Motivation

Neural models learn rich representations

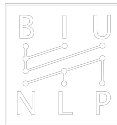
But can we control their content?

# Controlled Representation Learning



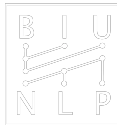
- Often, we want to make sure some concept is *not* encoded.

# Controlled Representation Learning



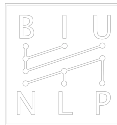
- Often, we want to make sure some concept is *not* encoded.
  - Word embeddings without tense distinctions

# Controlled Representation Learning



- Often, we want to make sure some concept is *not* encoded.
  - Word embeddings without tense distinctions
  - Sensitivity to content, but not to style

# Controlled Representation Learning



- Often, we want to make sure some concept is *not* encoded.
  - Word embeddings without tense distinctions
  - Sensitivity to content, but not to style
  - **Representations that do not leak protected attributes**

# The linear concept subspace hypothesis

Useful use case: the concept lives in low-dimensional **subspace** within the representation space.



# The linear concept subspace hypothesis

Useful use case: the concept lives in low-dimensional **subspace** within the representation space.

**How can we identify the concept subspace?**

# Formulation

We formulate an **adversarial game** between a projection matrix  $P \in \mathcal{P}_k$  that tries to remove the information, and a predictor  $\theta \in \Theta$  that tries to recover it.

# Formulation

We formulate an **adversarial game** between a projection matrix  $P \in \mathcal{P}_k$  that tries to remove the information, and a predictor  $\theta \in \Theta$  that tries to recover it.

$$\min_{\theta \in \Theta} \max_{P \in \mathcal{P}_k} \sum_{n=1}^N \ell \left( y_n, g^{-1} \left( \theta^\top P \mathbf{x}_n \right) \right)$$

# Formulation

We formulate an **adversarial game** between a projection matrix  $P \in \mathcal{P}_k$  that tries to remove the information, and a predictor  $\theta \in \Theta$  that tries to recover it.

$$\min_{\theta \in \Theta} \max_{P \in \mathcal{P}_k} \sum_{n=1}^N \ell \left( y_n, g^{-1} \left( \theta^\top \boxed{P} x_n \right) \right)$$

Projection matrix that tries to maximize the loss

# Formulation

We formulate an **adversarial game** between a projection matrix  $P \in \mathcal{P}_k$  that tries to remove the information, and a predictor  $\theta \in \Theta$  that tries to recover it.

$$\min_{\theta \in \Theta} \max_{P \in \mathcal{P}_k} \sum_{n=1}^N \ell \left( y_n, g^{-1} \left( \boxed{\theta}^\top \boxed{P} x_n \right) \right)$$

Predictor that tries to minimize it

Projection matrix that tries to maximize the loss

# Side trip: Iterative Nullspace Projection (INLP)

## **Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection**

**Shauli Ravfogel<sup>1,2</sup> Yanai Elazar<sup>1,2</sup> Hila Gonen<sup>1</sup> Michael Twiton<sup>3</sup> Yoav Goldberg<sup>1,2</sup>**

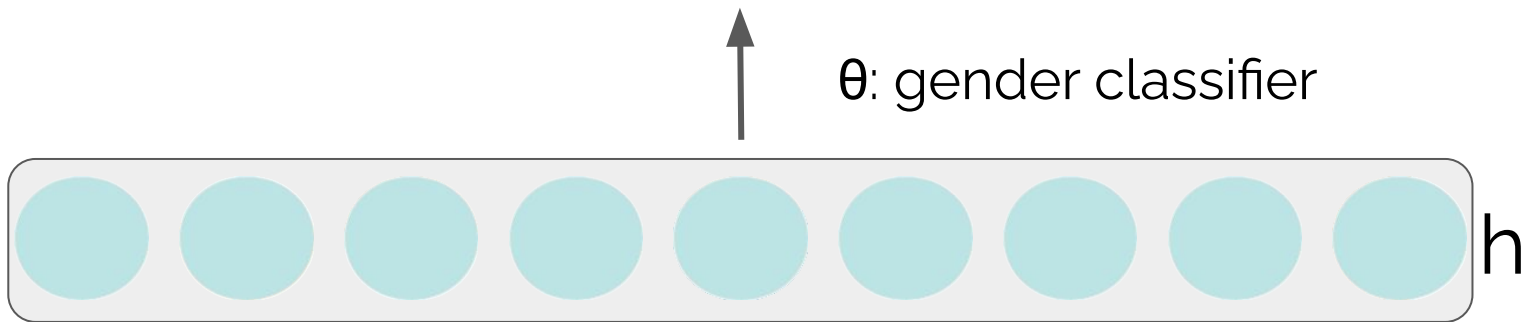
<sup>1</sup>Computer Science Department, Bar Ilan University

<sup>2</sup>Allen Institute for Artificial Intelligence

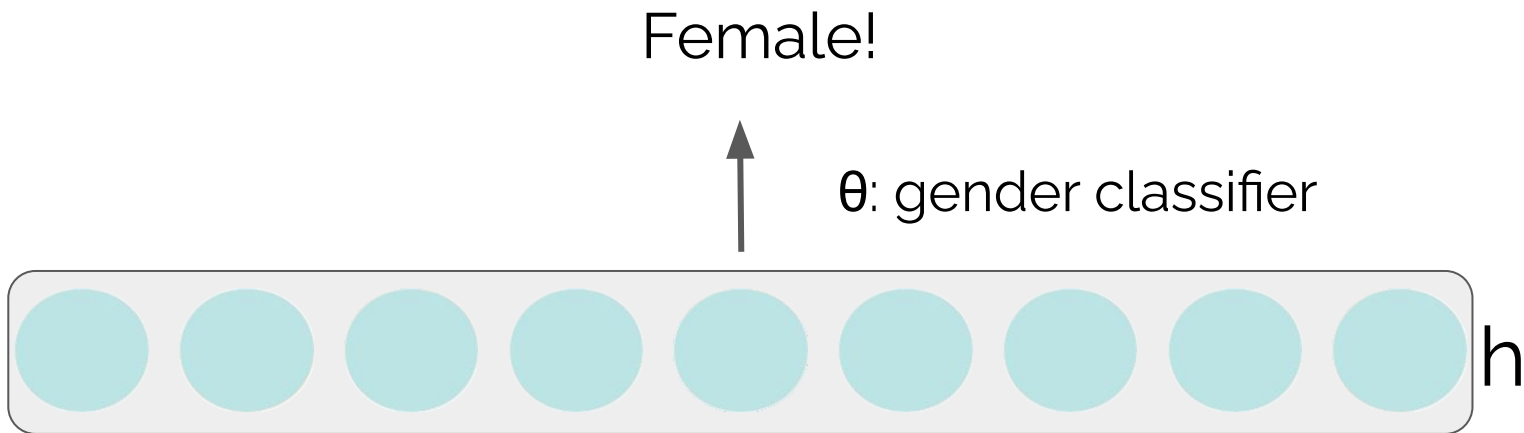
<sup>3</sup>Independent researcher

~~[shauli.ravfogel@bar-ilan.ac.il](#) [yanai.elazar@bar-ilan.ac.il](#) [hila.gonen@bar-ilan.ac.il](#) [mtwiton101@gmail.com](#) [yoav.goldberg@gmail.com](#)~~

# Nullspace projections



# Nullspace projections





# Nullspace projections

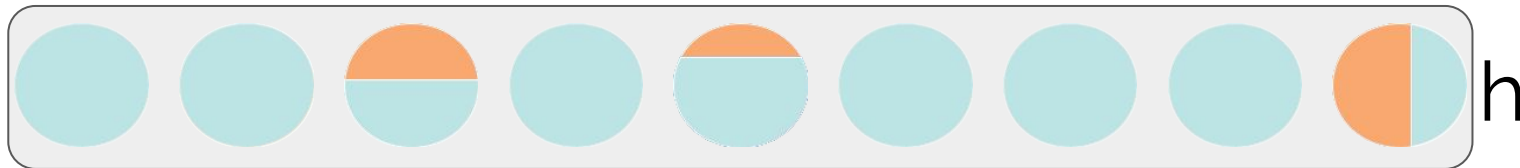


Features that  $\theta$  finds indicative of gender

Female!



$\theta$ : gender classifier



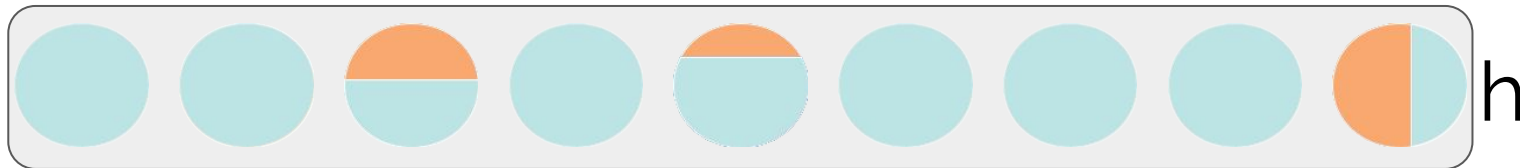
# Nullspace projections



Features that  $\theta$  finds indicative of gender



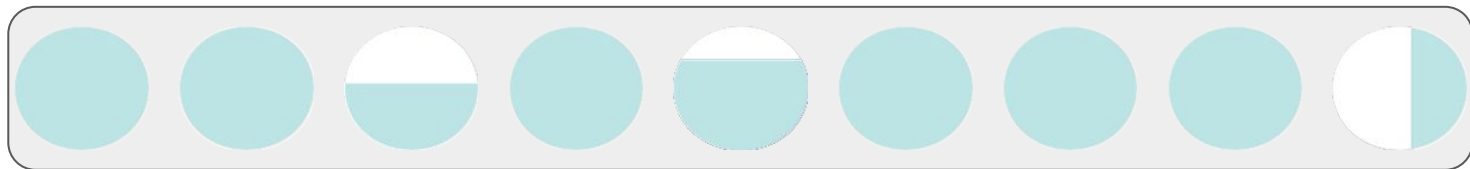
Remove gender features



# Nullspace projections



Features that  $\theta$  finds indicative of gender



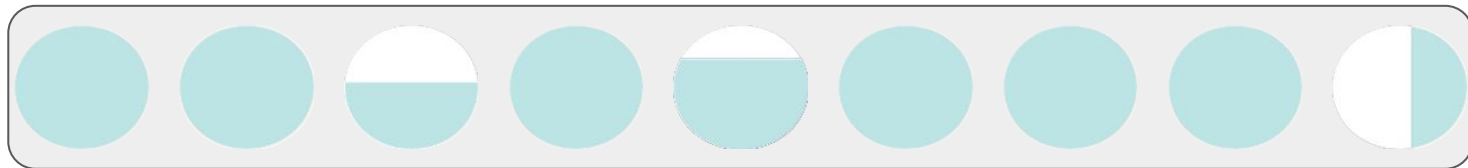
Remove gender features



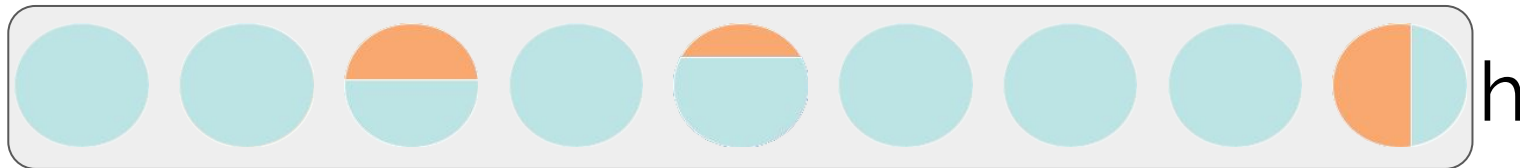
# Nullspace projections



Features that  $\theta$  finds indicative of gender



Remove gender features **how?**

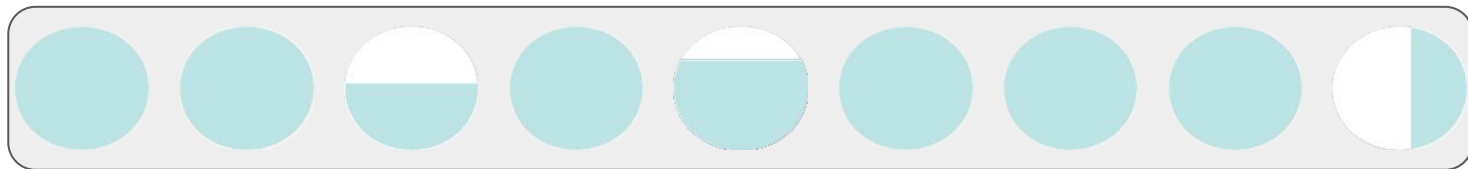


# Nullspace projections



Features that  $\theta$  finds indicative of gender

**Project**  $h$  to the orthogonal complement of  $\theta$



Remove gender features **how?**

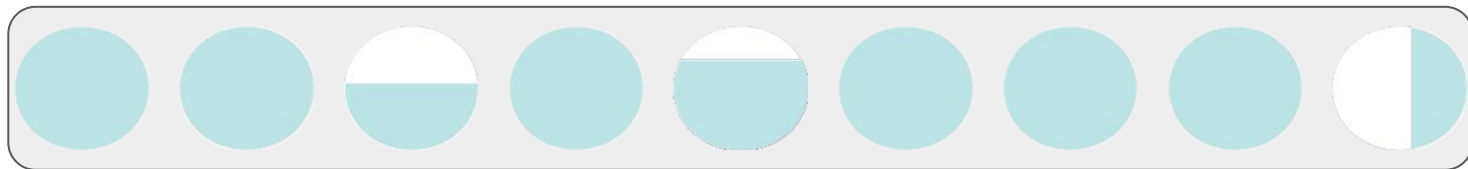


# Nullspace projections



Features that  $\theta$  finds indicative of gender

**Project**  $h$  to the orthogonal complement of  $\theta$   
**Iteratively**



Remove gender features **how?**



# Special Cases

We treat several important special cases of this objective:

# Special Cases

We treat several important special cases of this objective:

1. Linear regression



# Special Cases

We treat several important special cases of this objective:

1. Linear regression
2. Rayleigh quotient losses

# Special Cases

We treat several important special cases of this objective:

1. Linear regression
2. Rayleigh quotient losses
- 3. Classification**

# Special Cases

We treat several important special cases of this objective:

1. Linear regression
2. Rayleigh quotient losses

Closed-form solution  
(details in the paper)

## 3. **Classification**

# Special Cases

We treat several important special cases of this objective:

- |                             |   |  |
|-----------------------------|---|--|
| 1. Linear regression        | ] | Closed-form solution<br>(details in the paper) |
| 2. Rayleigh quotient losses |   |  |
| 3. <b>Classification</b>    | ] | Gradient-based optimization                    |

# Classification Case

The loss is an arbitrary classification loss (hinge, logistic, etc).

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^D} \max_{P \in \mathcal{P}_k} \sum_{n=1}^N y_n \log \frac{\exp \boldsymbol{\theta}^\top P \mathbf{x}_n}{1 + \exp \boldsymbol{\theta}^\top P \mathbf{x}_n}$$

# Convex Relaxation (RLACE)

The problem is nonconvex in  $P$ .

# Convex Relaxation (RLACE)

The problem is nonconvex in  $P$ .

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^D} \max_{\substack{P \in \mathcal{P}_k \\ P \in \mathcal{F}_k}} \sum_{n=1}^N y_n \log \frac{\exp \boldsymbol{\theta}^\top P \mathbf{x}_n}{1 + \exp \boldsymbol{\theta}^\top P \mathbf{x}_n}$$

# Convex Relaxation (RLACE)

The problem is nonconvex in  $P$ .

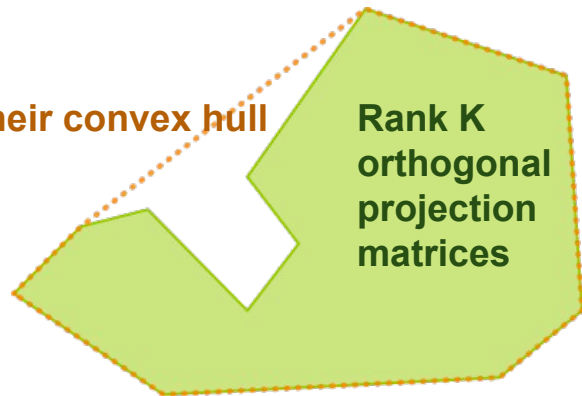
$$\min_{\boldsymbol{\theta} \in \mathbb{R}^D} \max_{\substack{P \in \mathcal{P}_k \\ P \in \mathcal{F}_k}} \sum_{n=1}^N y_n \log \frac{\exp \boldsymbol{\theta}^\top P \mathbf{x}_n}{1 + \exp \boldsymbol{\theta}^\top P \mathbf{x}_n}$$

Where we define:

$$\mathcal{F}_k = \text{conv}(\mathcal{P}_k)$$

Their convex hull

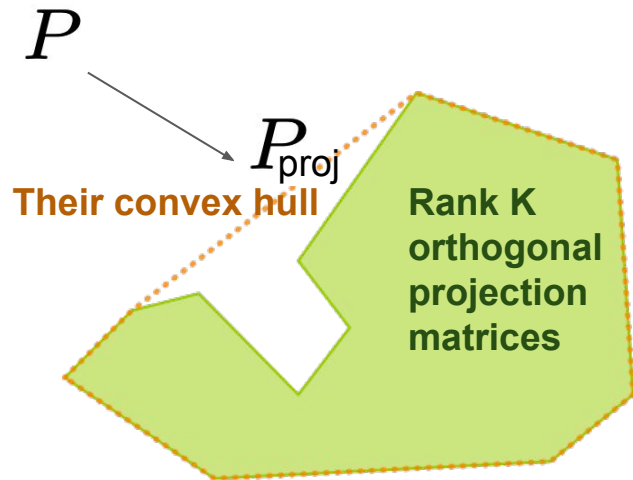
Rank K  
orthogonal  
projection  
matrices





# Convex Relaxation (RLACE)

In training, we optimize over an arbitrary matrix  $P$ , and we project it to the convex hull of orthogonal projection matrices at each step.

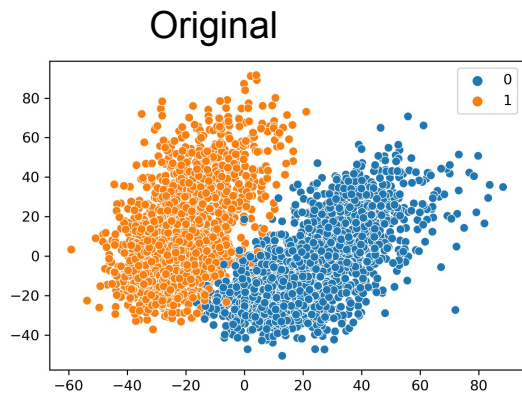


# Experimental Evaluation

We conduct experiments on GloVe embeddings and on contextualized representations of short biographies (annotated for both gender and profession).

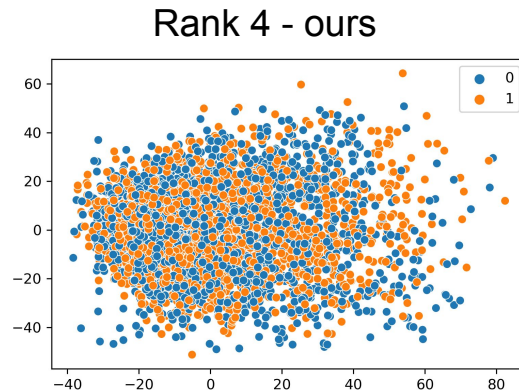
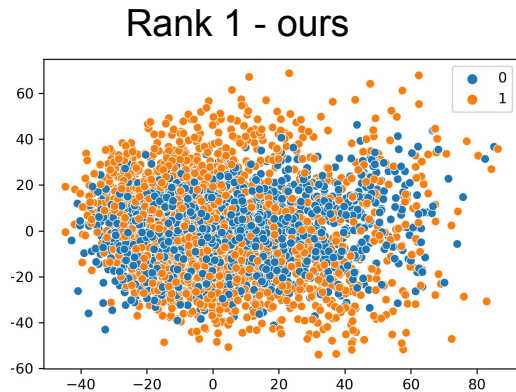
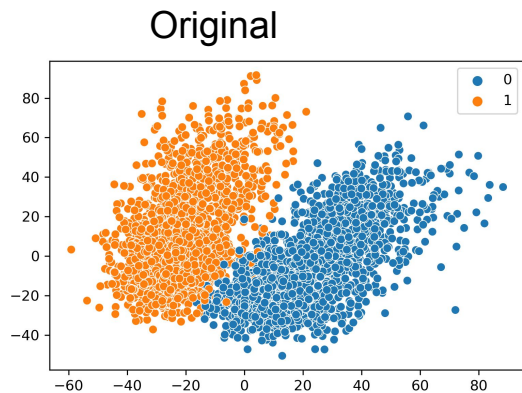
# Results

PCA: less clustering of representations by gender



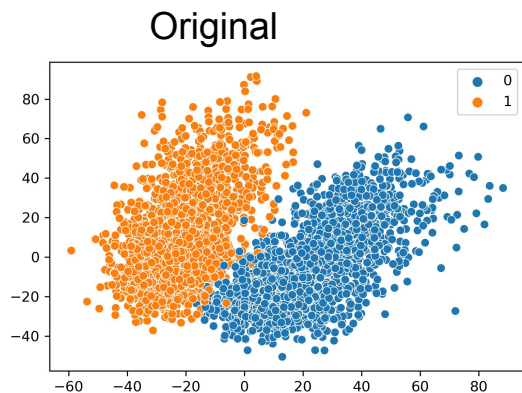
# Results

PCA: less clustering of representations by gender

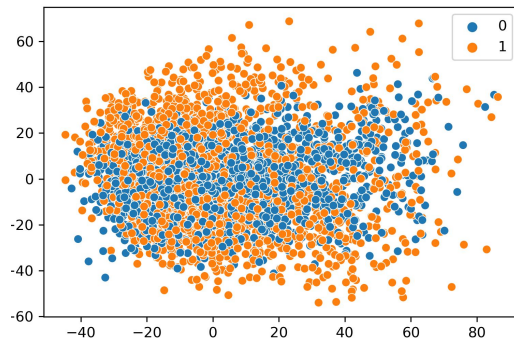


# Results

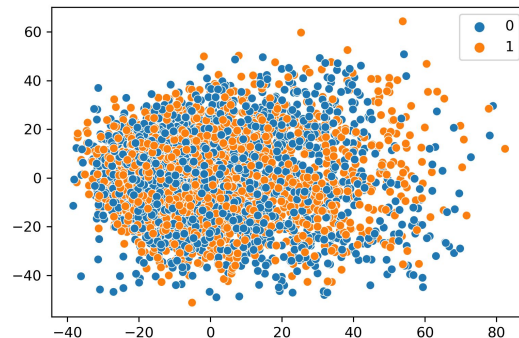
PCA: less clustering of representations by gender



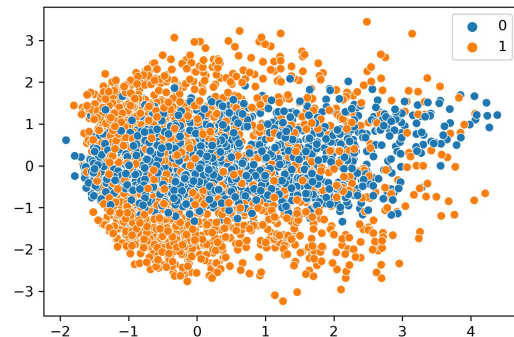
Rank 1 - ours



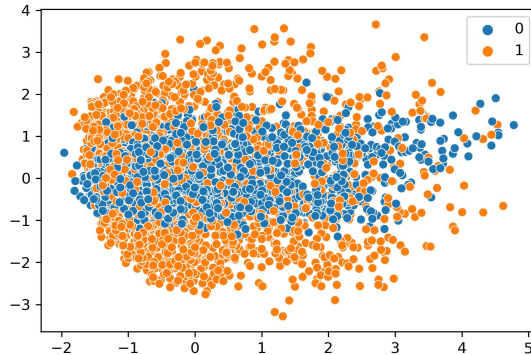
Rank 4 - ours



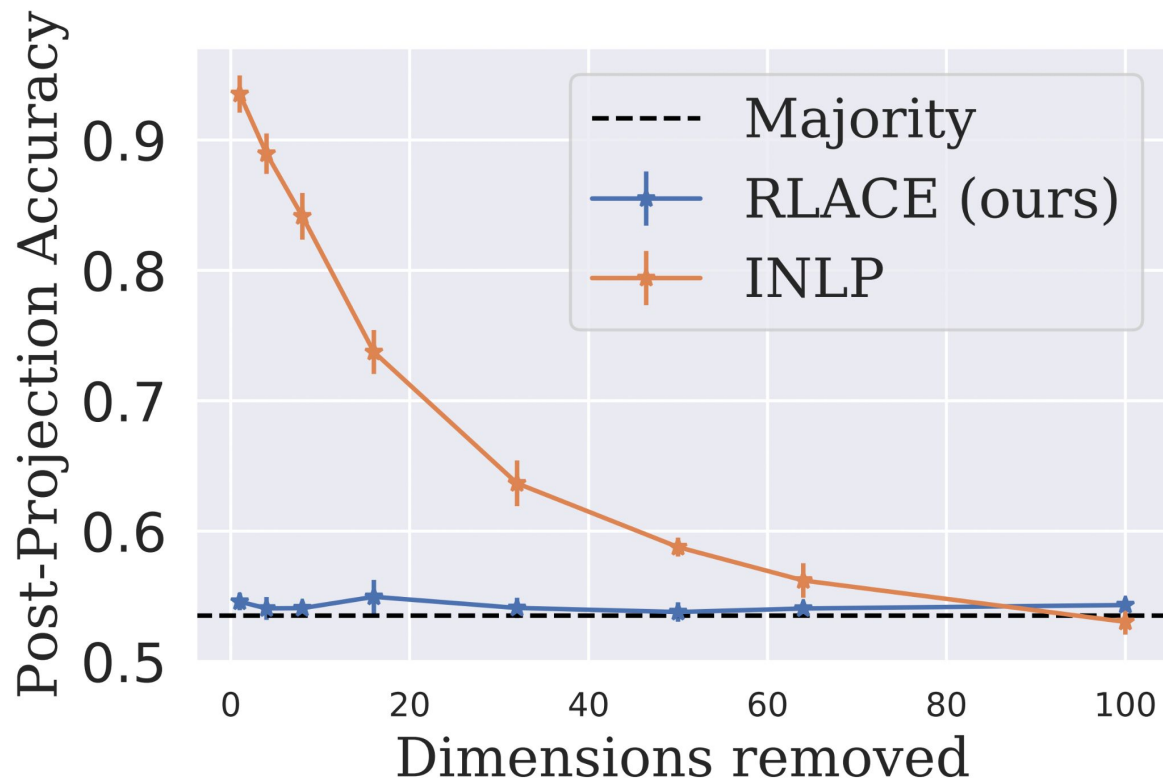
Rank 1 - INLP



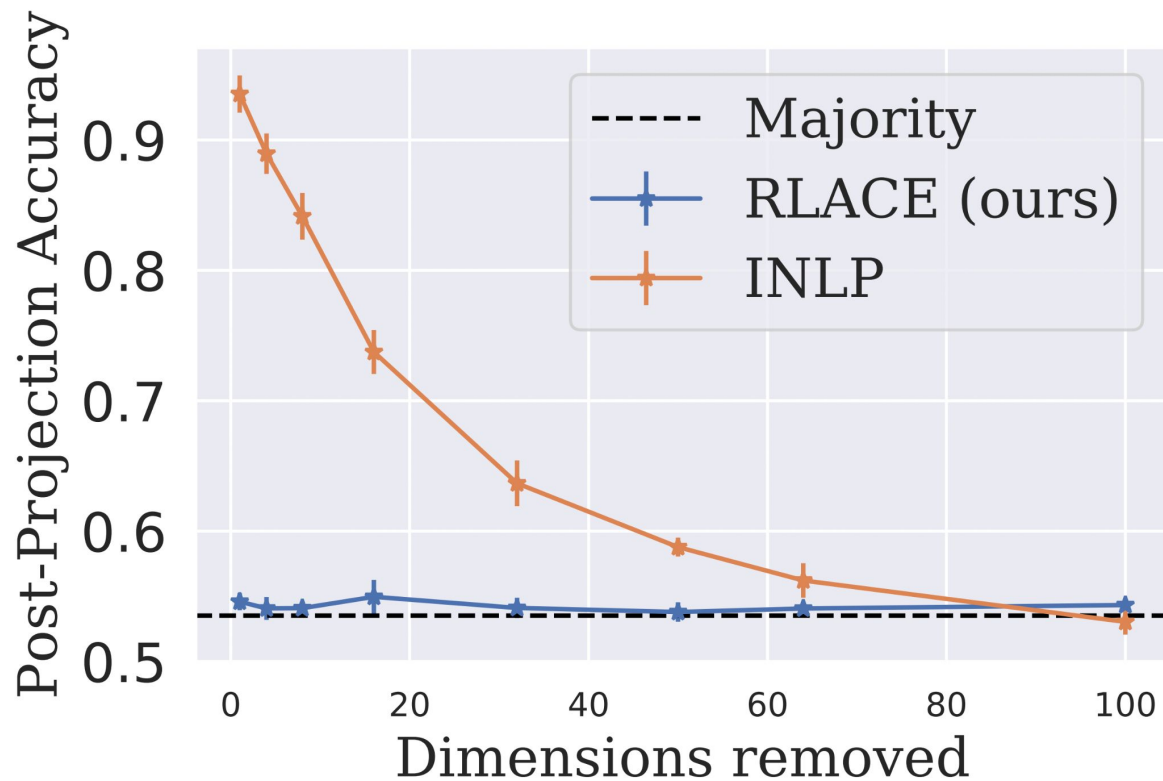
Rank 4 - INLP



# Comparison with INLP - BERT (bios data)



# Comparison with INLP - BERT (bios data)



Additional results in the paper.

# Application on Images



Original



# Application on Images



# Application on Images



Original

Smile

Glasses

# Conclusions

- We have formally defined the problem of removing linear “concept subspaces”

# Conclusions

- We have formally defined the problem of removing linear “concept subspaces”
- We present analytical solutions in some cases, and provide a relaxation which works well for others.

# Conclusions

- We have formally defined the problem of removing linear “concept subspaces”
- We present analytical solutions in some cases, and provide a relaxation which works well for others.
- High level conclusion: it's sometimes valuable to constrain our model (in contrast to the “more parameters is better” trend)

# Thanks!

