



# **Closing the Gender Wage Gap: Adversarial Fairness in Job Recommendation**

Clara Rus, Jeffrey Luppès, Harrie Oosterhuis, and Gido Schoenmacker

# GOAL



Have a fair job recommendation system

→ does not discriminate against gender, race, religion etc.

$$F(X, S = 0) = F(X, S = 1) = Y$$

# THINGS CAN GO WRONG...

## Facebook's ad delivery system still has gender bias, new study finds

*Women are excluded from seeing some job listings 'beyond what can be legally justified'*

**Man is to Computer Programmer as Woman is to Homemaker**

## Amazon Scraps Secret AI Recruiting Engine that Showed Biases Against Women

AI Research scientists at Amazon uncovered biases against women on their recruiting machine learning engine

## Twitter taught Microsoft's AI chatbot to be a racist





# SOLVING THE PROBLEM



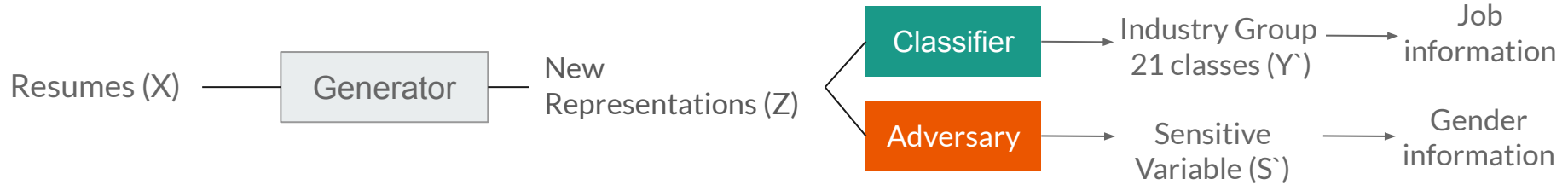
**Goal:** Obtain vector representation that lead to fair predictions

**Solutions:**

- (1) Word Substitution Method
- (2) Adversarial Method



# ADVERSARIAL METHOD



$$L = \alpha L_{cla}(Z, Y') + \beta L_{adv}(Z, S')$$

# EVALUATION



## Performance

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

## Fairness: Statistical Parity

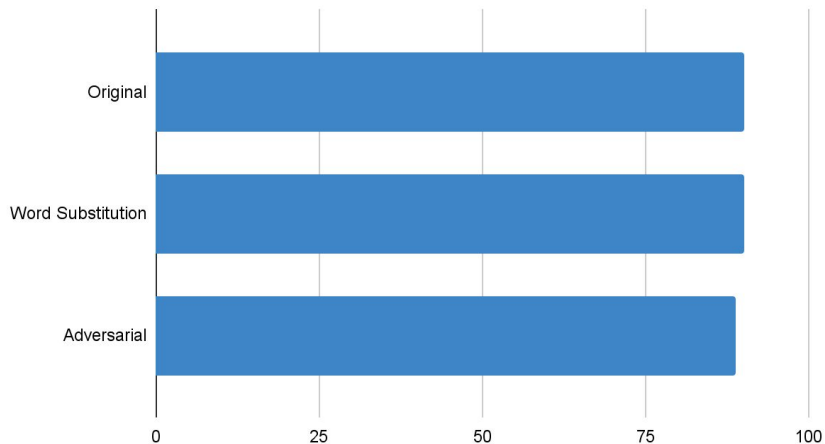
$$P(\text{cla}(Z) = 1 | S = 1) - P(\text{cla}(Z) = 1 | S = 0) < \epsilon$$



# PREDICTION OF INDUSTRY GROUP

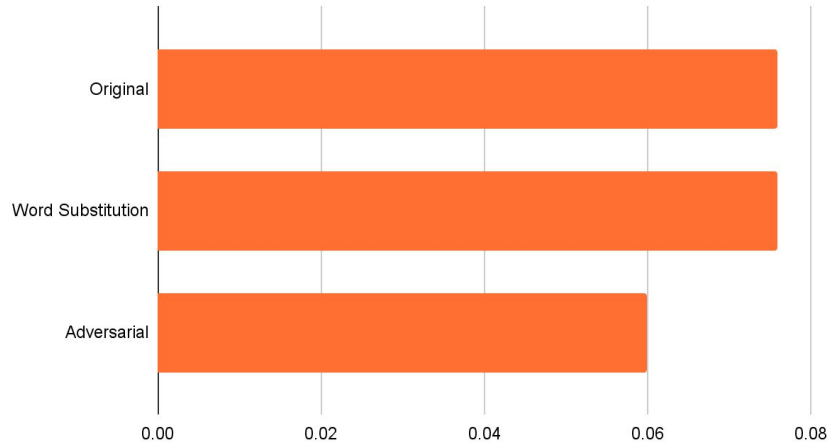
## PERFORMANCE FAIRNESS TRADE-OFF

Performance: Accuracy



Higher is better

Unfairness: Statistical Parity

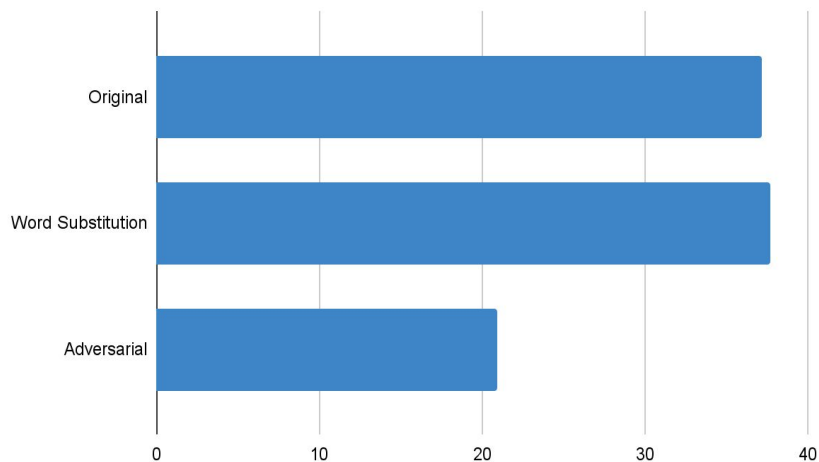


Lower is better

# PREDICTION OF INDUSTRY GROUP

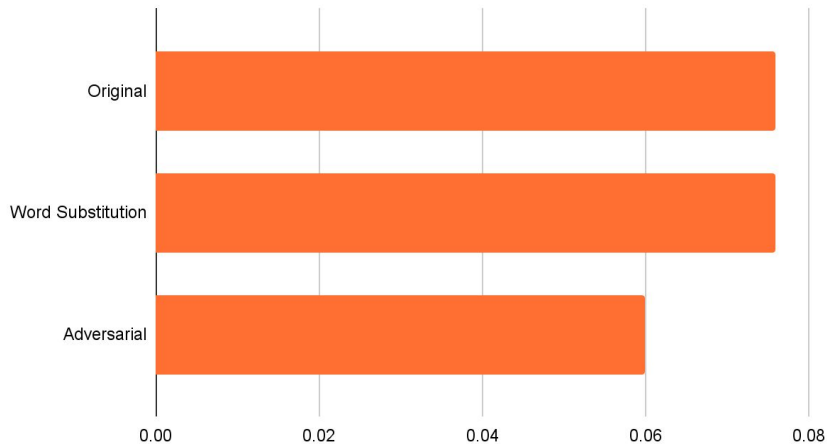
## PERFORMANCE FAIRNESS TRADE-OFF

Performance: True Positive Rate



Higher is better

Unfairness: Statistical Parity

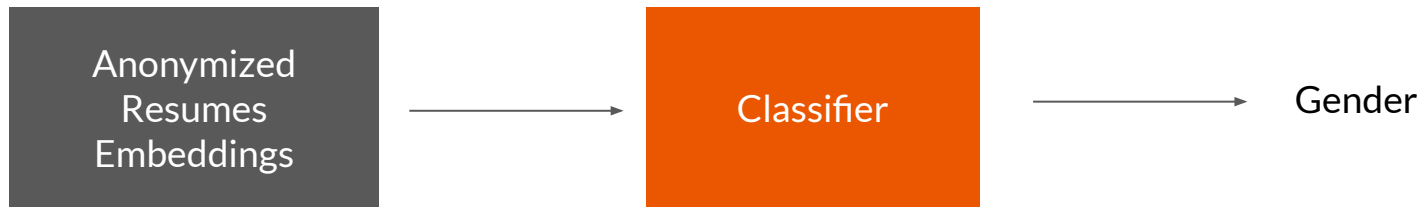


Lower is better

# PREDICTION OF SENSITIVE VARIABLE

**Goal:** Have a model make predictions independent of the gender

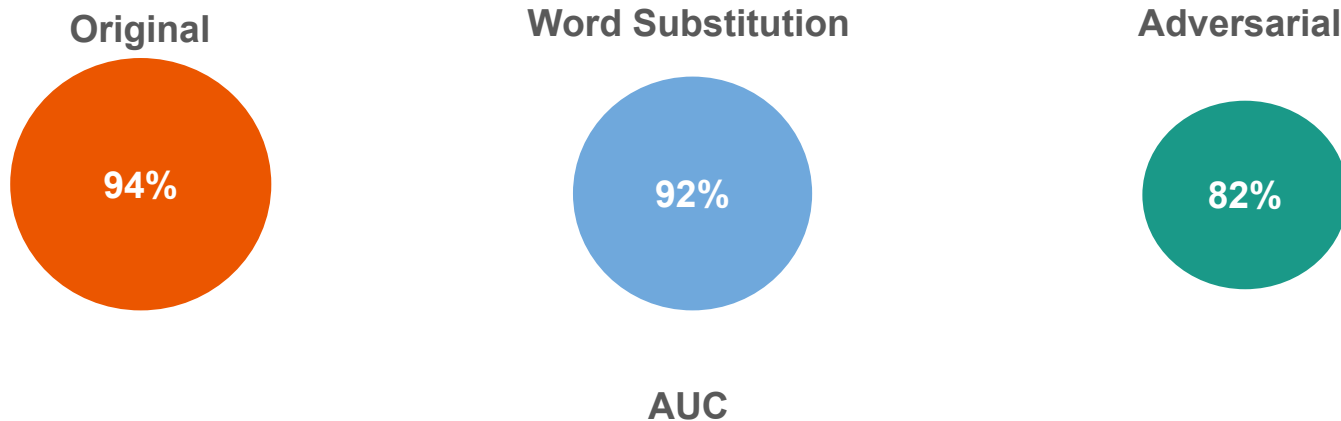
If a model can learn to predict the gender → not independent of the gender



## PREDICTION OF SENSITIVE VARIABLE

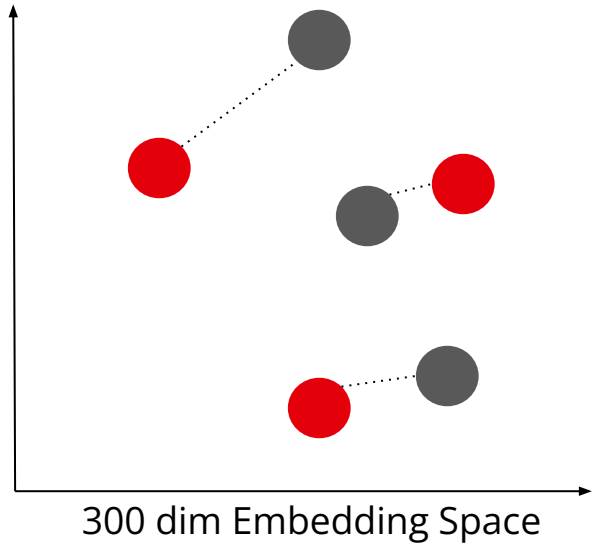
**Goal:** Have a model make predictions independent of the gender

If a model can learn to predict the gender → not independent of the gender

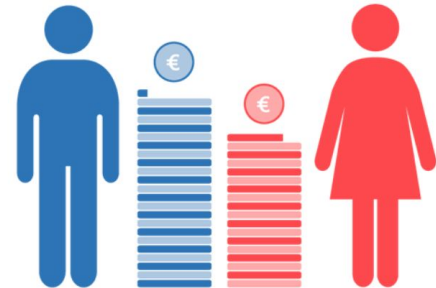


# SALARY ASSOCIATION TEST

For each resume find the most similar job description.



- Resumes
- Vacancies
- ..... Euclidean Distance



# SALARY ASSOCIATION TEST

For each resume find the most similar job description.

Original



€1680

Word Substitution



€1900

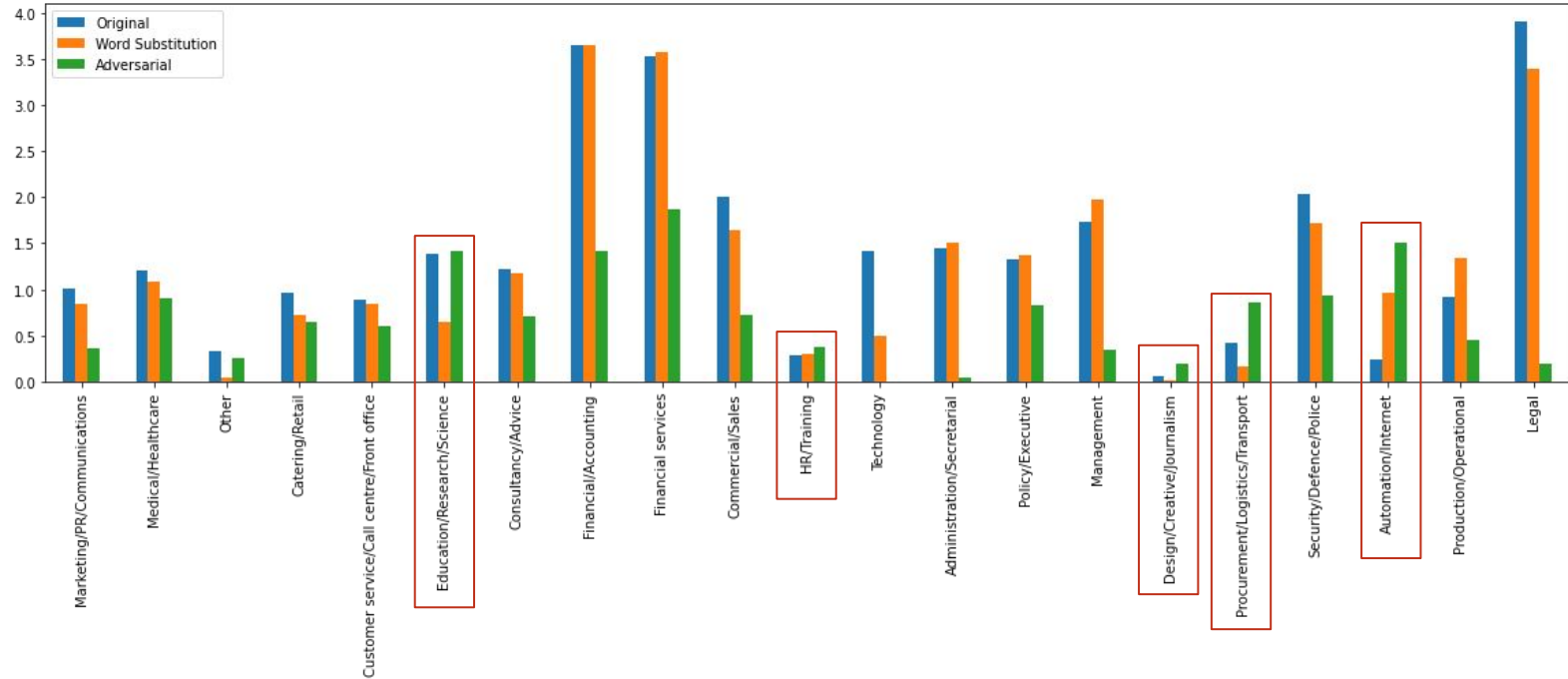
Adversarial



€180

Average wage gap per year

# SALARY ASSOCIATION TEST



## FUTURE WORK



- Evaluate the quality of the matches
- More complex matching system for the Salary Association Test
- Include non-binary genders
- Use other types of embeddings



# MAIN CONTRIBUTIONS



- (1) Applied adversarial debiasing methods on large scale real text data
- (2) Trade-off between fairness and performance on a multi-label classification task
- (3) The Salary Association Test: showing the impact on a real business case scenario
- (4) Provided open-source software

# CONCLUSIONS



- Ignoring the bias in the data can lead to unwanted discriminatory behaviour
  - > producing a wage gap
- The adversarial method:
  - improved fairness with the cost of lowering the performance
  - reduced the wage gap by 89%



**THANK YOU!**