

# Data-Driven Analysis of Gender Fairness in the Software Engineering Academic Landscape

Giordano d'Aloisio, Andrea D'Angelo, Francesca Marzi,  
Diana Di Marco, Giovanni Stilo, Antiniscia Di Marco

Università degli Studi dell'Aquila / Italy

UNIVERSITÀ  
DEGLI STUDI  
DELL'AQUILA



DISIM  
Dipartimento di Ingegneria  
e Scienze dell'Informazione  
e Matematica



SOBIGDATA.it

ITALIAN RESEARCH INFRASTRUCTURE

In this work, we study the problem of **Gender Bias in academic promotions** in the Informatics and Software Engineering Italian communities.

We mine public data about role promotions and **academic productivity** to compute **Disparate Impact**, a formal definition of bias.

Researcher

Associate  
Professor

Full  
Professor



Literature Review

Analysis Description

Experimental Results

# Gender Bias in Classic Academic Systems: A review

Search string:

```
allintitle : gender bias OR academic recruitment OR  
gender discrimination OR  
Women ' s faculty recruitment OR faculty equity OR  
career advancements OR  
Italian universities OR selection processes
```



21 papers

Classified based on:

Context

Process

Privacy

Analytical Method Used

Year

# Gender Bias review (II)

## Process



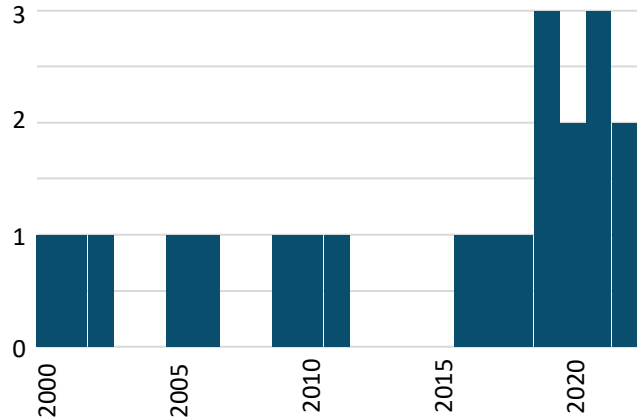
Only 2 of the papers focused on Gender Bias in Productivity.

## Privacy of the Data



Only 38% of the papers use Public data from trusted sources. Private data were usually collected through interviews or surveys.

# Gender Bias review (III)



Research on the subject drastically increased in recent years.

**None** of the papers focus on the problem in the Informatics community.

Reliance on formal metrics of bias is **severely lacking**.



Motivation!



# Gender Bias review (IV)

Our study aims to **formally** analyze the issue of gender bias in academic promotions in the **Informatics and SE Italian communities**.



We mined **public data from trusted sources** and included **productivity** metrics. We used **formal bias definitions** to analyze our results.



Literature Review

Analysis Description

Experimental Results



# Data Gathering

- First, we scraped publicly available data from official sources.

Career and affiliations data was obtained from the MIUR (Ministry of University and Research) and National Scientific Qualification websites.



Productivity and publication metrics were obtained via the Scopus API.



Scopus

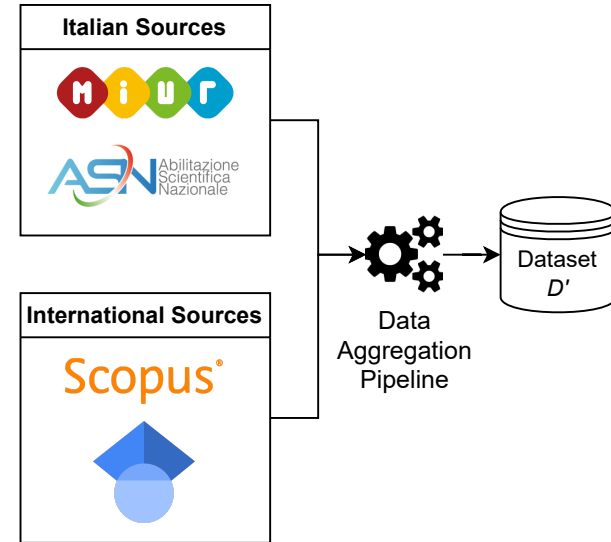
Public data from official sources

# Data Pipeline

- Second, we merged the data from different sources.

We employed **Regular Expression** logic to split Full Names into Name and Surname and remove special characters.

We merged the two datasets and split the productivity metrics in order to obtain a **time series** of publications and citations for each record.



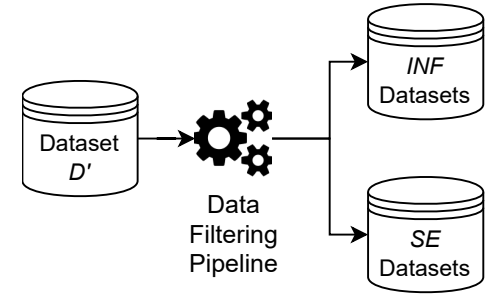
# Data Pipeline (II)

- Third, we filtered the resulting data.

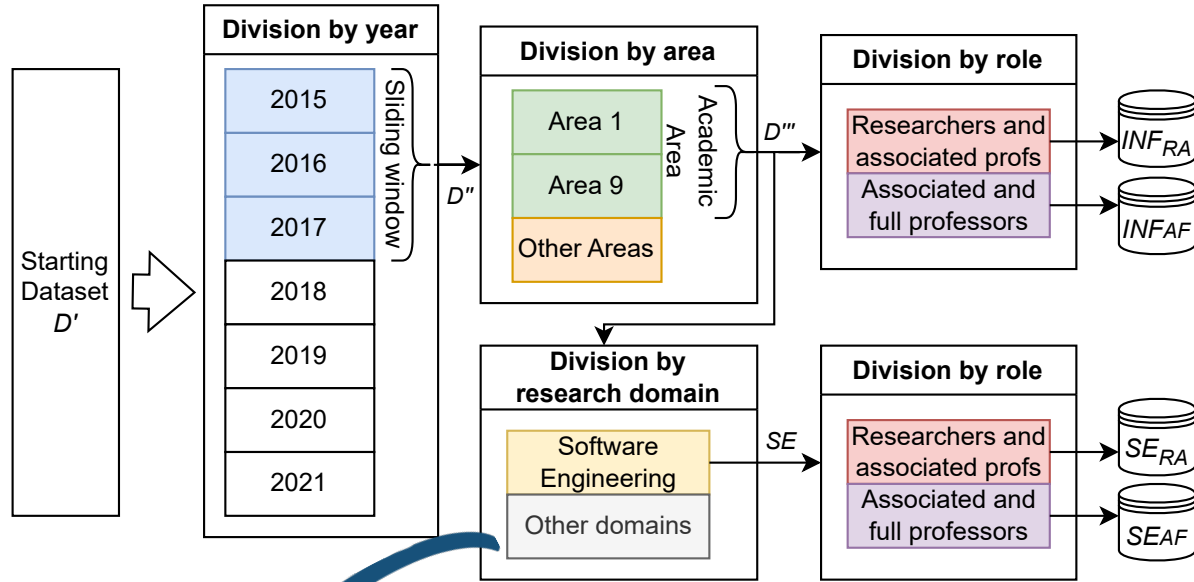
The dataset  $D'$  was split according to a **sliding time window** of fixed size (3 years).

We only selected **specific fields of research**, related to Computer Science and Engineering.

We split the dataset into **Informatics and Software Engineering** and by Academic Role.



# Data Pipeline (III)



# Data Pipeline (IV)

➤ The end results are 4 datasets, each divided into time windows:



Informatics Researchers and Associated ( $INF_{RA}$ )



Software Engineering Researchers and Associated ( $SE_{RA}$ )



Informatics Associated and Full ( $INF_{AF}$ )



Software Engineering Associated and Full ( $SE_{AF}$ )



Literature Review

Analysis Description

Experimental Results

- To compute the Bias in each dataset, we refer to the formal definition of Disparate Impact:

## Disparate Impact (DI):

Disparate Impact compares the probability of having a *Positive Outcome* while being in the *privileged* or *unprivileged* group. Formally:

$$DI = \frac{P(Y = y_p | X = x_{unpriv})}{P(Y = y_p | X = x_{priv})}$$

- The closer this value is to 1, the «fairer» the dataset.

# Bias Metric (II)

- To compute the Bias in each dataset, we refer to the formal definition of Disparate Impact:

## Disparate Impact (DI):

Disparate Impact compares the probability of having a *Positive Outcome* while being in the *privileged* or *unprivileged* group. Formally:

$$DI = \frac{P(Y = y_p | X = x_{unpriv})}{P(Y = y_p | X = x_{priv})}$$

Associate/Full professors

$x_{unpriv}$  are women,  
 $x_{priv}$  are men



- To compute the Bias in each dataset, we refer to the formal definition of Disparate Impact:

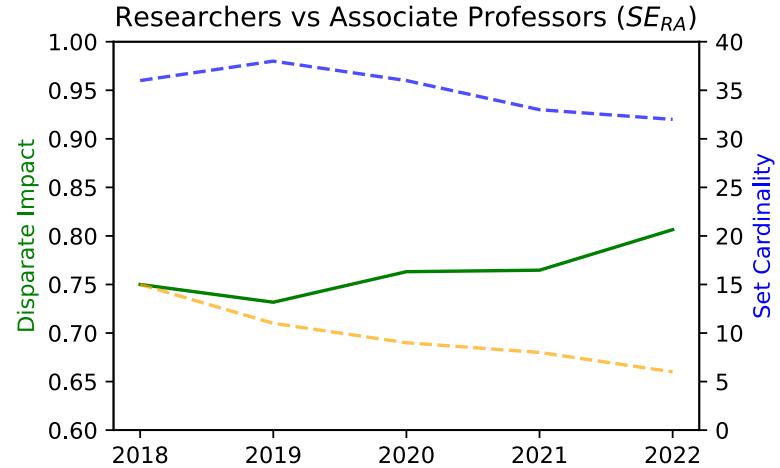
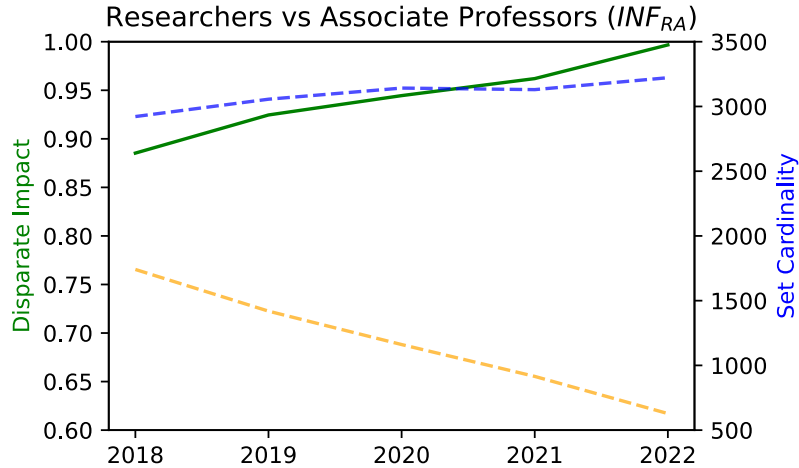
## Disparate Impact (DI):

Disparate Impact compares the probability of having a *Positive Outcome* while being in the *privileged* or *unprivileged* group. Formally:

$$DI = \frac{P(Y = y_p | X = x_{unpriv})}{P(Y = y_p | X = x_{priv})}$$

**Disparate Impact** does not need a classifier to be computed, as it can be calculated on the dataset itself. We compute the probabilities by appropriately slicing the dataset.

# Experimental Results

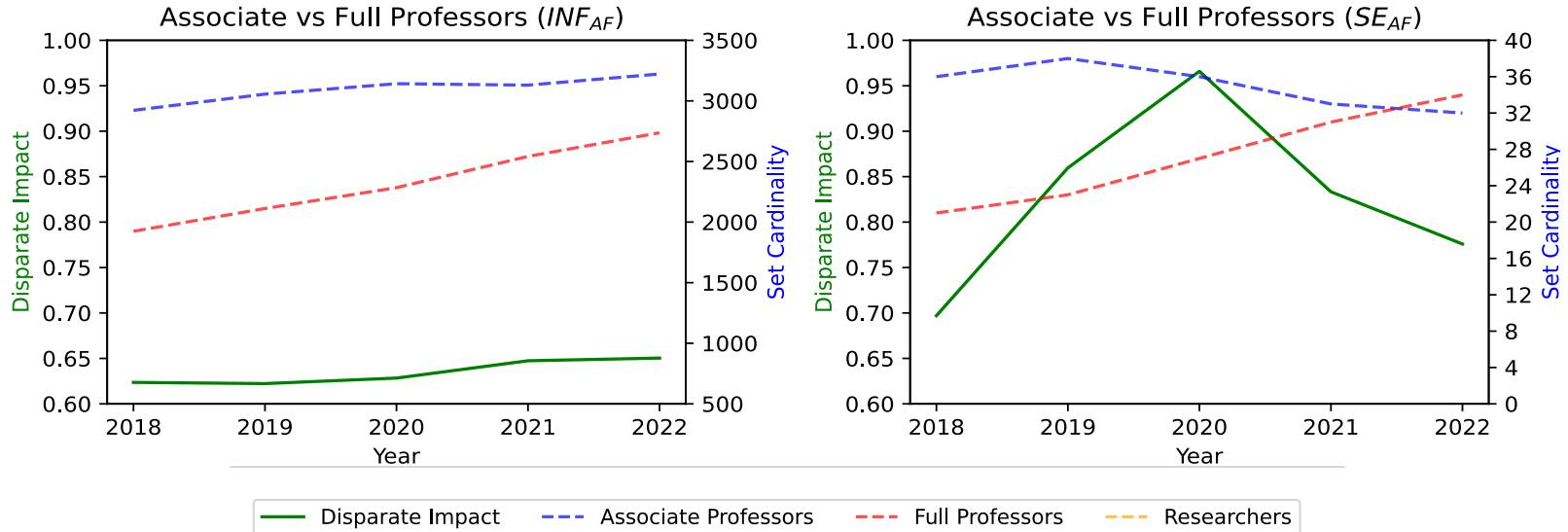


— Disparate Impact    - - Associate Professors    - - Full Professors    - - Researchers

The **SE Community** appears to have **more gender bias** in the career promotion **from Researchers to Associate Professors**.

In both contexts, bias is steadily decreasing.

# Experimental Results (II)



On the other hand, the **SE Community is much fairer** w.r.t. promotions from Associate to Full Professors.

The peak in fairness was registered in 2020.

# Main Takeaways

- We performed a Literature Review on the subject which highlighted several critical points;
- We built a joint dataset from several different official sources and processed it through a pipeline;
- We used a formal metric to show that the SE community is lagging behind in fairness for promotions from Researchers to Associate Professors, but is fair from promotions from Associate to Full Professors

# Open Problems

- Expanding the study to other Areas and Countries (need public data from official sources);
- Train a ML classifier to predict the Academic Position of a Researcher, study feature importance and possible bias related to gender;



**Thank you for your attention**

UNIVERSITÀ  
DEGLI STUDI  
DELL'AQUILA



DISIM  
Dipartimento di Ingegneria  
e Scienze dell'Informazione  
e Matematica