# Creation and Evaluation of a Self-Learning Course for Data Hazard Labels

## Bachelor Thesis

Joost Krüger

17.07.2024

# Table of Contents

# Introduction

- **Creation of a self learing course using LiaScript**

- **Aimed at university students**

- **Intersection between ethics and data science**

- **Course can be used as supplementary material for lectures**

# Data Hazards

**Created and maintained by the Data Hazard Project**

**Goal: Create a tool to help initiate ethics discussion on data science**

**Data hazard labels:**

- **Similar to real Hazard labels**

- **Currently 11 labels for data science and 5 for synthetic biology**

- **Can be assigned to any projects using data science**

- **Not strict true or false**

- **Multiple labels can apply**

# Data Hazards


Generic Data Hazard


Reinforces Existing Biases


Automates Decisionmaking


Ranks or Classifies People


Danger of Misuse


May cause Harm

**Not present:**

- **High Environmental Cost**

- **Risk to privacy**

- **Lacks Informed Consent**

- **Difficult to understand**

# Current Prototype

# Current Prototype



Data Hazard Labels

**Reinforces Existing Biases**
- Definition
- Examples
- Prevention of Bias
- Videos
- Quiz

**Ranks or classifies people**
- Definition
- Examples
- Prevention of improper ranking and classification
- Videos
- Quiz

**Automates decision making**
- Definition
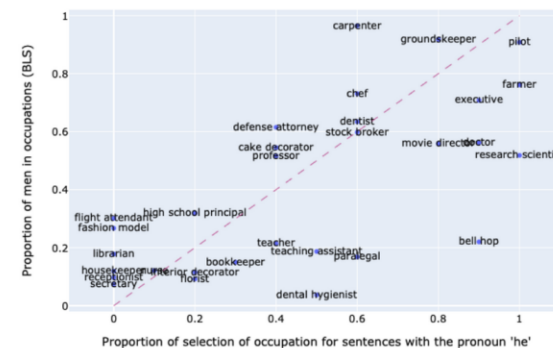- Examples
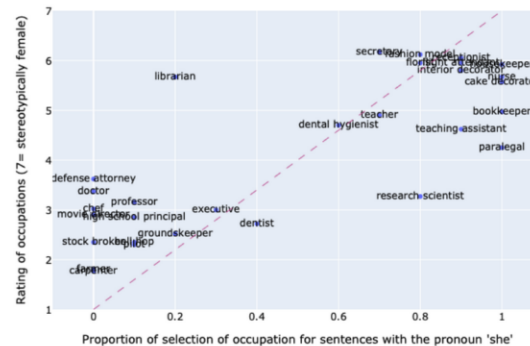- Precautions for automated decision making

## Examples

**Input data**
An algorithm that uses historic employment data that comes to the conclusion that men are more suited to managerial position, as historically men were favoured or even the only allowed candidates for such positions.

**Societal Bias**
Natural Language processing data can reinforce sexist biases due to a bias in training data. This could mean that a model evaluates certain jobs such as secretary or caretaker as intrinsically linked to women.

Such cases were studied and both natural and large language models were found perpetuate stereotypes. Since these models are used more, great care should be taken when working with such cases and active measures taken to prevent the spread of such stereotyping. Such cases prove furthermore that

# Current Prototype

# Current Prototype

# Current Prototype

# Evaluation and Interviews

- Next step is interviews
- Interviewing experts from relevant fields
- Getting a feedback on design and content
- Use feedback to improve learning course

# Questions

Thank you for your attention!

# Quellen

- The data hazard project (https://datahazards.com/index.html)

- Design based research icon (https://commons.wikimedia.org/wiki/File:DBR_german_colour.svg)