

Responsible Data Science

Fairness Review

Prof. George Wood

Center for Data Science
New York University

Fairness module, key ideas

Week 1:

- Goals, benefits, and harms of DS systems
- Stakeholders

Week 2:

- Fairness in classification and risk assessment
- Individual fairness vs group fairness
- Disparate treatment vs disparate impact
- Impossibility result (calibration versus balance of errors)
- Three types of bias in computer systems (pre-existing, technical, emergent)

Week 3:

- Five fairness definitions (FTU, individual fairness, demographic parity, equalized odds, calibration)
- Causal models, causal framework for fairness (causal diagrams, counterfactual fairness)

Week 4:

- Causal framework for fairness continued (causal pathways, counterfactual privilege)
- Philosophical frameworks for fairness
- Fairness as equal opportunity (EOP), formal EOP, substantive EOP

Fairness module, HW1

DS-UA 202, Responsible Data Science, Spring 2022

Homework 1: Algorithmic Fairness
Due on Friday, March 3 at 11:59pm EST

Objectives

This assignment consists of written problems and programming exercises on algorithmic fairness.

After completing this assignment, you will:

- Understand that different notions of fairness correspond to points of view of different stakeholders, and are often mutually incompatible.
- Gain hands-on experience with incorporating fairness-enhancing interventions into machine learning pipelines.
- Learn about the trade-offs between fairness and accuracy.
- Observe the effect of hyperparameter tuning on performance, in terms of both accuracy and fairness.

You must work on this assignment individually. If you have questions about this assignment, please send an email to all instructors.

Submission instructions

Provide written answers to Problems 1, 2, and 3 in a **single PDF file**. (It is recommended that you [Google Docs](#) to prepare this PDF, but you may instead use Word or LaTeX). Provide code in answer to Problem 2 in a **Google Colaboratory notebook**. Both the PDF and the notebook should be turned in as Homework 1 on Brightspace. Please clearly label each part of each question.

Fairness module, HW1

The screenshot shows a Google Colab notebook interface. At the top, the title is "HW1 template DS UA 2022" with a star icon. Below the title is a menu bar with options: File, Edit, View, Insert, Runtime, Tools, Help, and "All changes saved".

On the left side, there is a "Table of contents" panel with a search icon and a close button. It lists the following sections:

- Setup
 - Packages
 - Load data
 - Set protected attribute and target
 - Split data
 - Scale features in the data
- 2 (a)
 - Train a random forest model (baseline)
 - Calculate metrics
- 2 (b)
 - Transform the original data using Disparate Impact Remover at five repair levels and calculate metrics
- 2 (c)
 - Train a Prejudice Remover model at three eta values and calculate metrics

On the right side, the main content area has a "+ Code" and "+ Text" button. The title of the notebook is "RDS (DS-UA 2022) Spring 2022: Homework 1 Template". Below the title, there is a paragraph: "This notebook is a template for problem 2. You should save a copy of code to setup the analysis is provided for you here. You should not expect..." followed by another paragraph: "Some suggested steps are included as comments in the below code solutions or approaches are acceptable)."

Below the text, there are two expandable sections: "Setup" and "Packages". The "Packages" section is expanded, showing a code cell with the following commands:

```
!git clone https://github.com/lurosenb/superquail
!pip install aif360==0.3.0
!pip install BlackBoxAuditing
!pip install tensorflow==1.13.1
```

See Brightspace for:

- Assignment details
- Colab template

Note: you must use Colab!