

Deadline 8 Apr 2024 23:59 (CPH time).

This is a group submission (2-3 person group).

Using the Folktales dataset (see jupyter notebooks on github) use the two techniques presented in lecture 7 to create fairer representations of the dataset. For this assignment you will work with 2 protected attributes: gender and race (in the dataset they are called **SEX** and **RAC1P**).

The Task: Train three binary classifiers to predict income (label = True if income > \$25k, otherwise label = False)

1. Train one classification model on the raw dataset and calculate its general accuracy, and respectively the accuracies for men and women and for different races. The model does not need to be fancy, logistic regression or Random Forest are completely fine choices. Remember to evaluate the model using cross validation.
2. Using the “fairer” (reprojected) versions of the dataset, build two classification models (see more below) and calculate: a) their overall accuracy, b) their accuracies split for men and women, and c) their accuracies split for different races. **Subtasks:**
 - Build one classification model trained on data reprojected using the de-correlation method from the paper “A Geometric Solution to Fair Representations”. Record your results and create a plot of how accuracies vary as functions of $\lambda \in [0,1]$.
 - Build one classification model using reprojected data from FairPCA.
3. Write a report, **max 3 pages long**, that describes your approach (e.g. which machine learning model did you use, how did you split the data, etc.) and summarizes your findings and conclusions. You must address the following:
 - How does debiasing affect the overall classification accuracy? For the de-correlation method, what effect do different values of λ have on the overall accuracy? What value of λ is the best to use? Does such a value even exist?

Assignment 2 – debiasing data and fairer models

- How are classification accuracies of different genders, and race groups affected when debiasing data? I.e. does debiasing data work? I.e. are gaps in accuracies less biased?

Please submit your report (in PDF format) & jupyter notebook on learnIT.