# Author Contribution Ckecklist Form for "A Burden Shared is a Burden Halved: A Fairness Adjusted Approach to Classification"

This form documents the artifacts associated with the article (i.e., the data and code supporting the computational findings) and describes how to reproduce the findings.

## Part 1: Data

☐ This paper does not involve analysis of external data (i.e., no data are used or the only data are generated by the authors via simulation in their code).

☒ I certify that the author(s) of the manuscript have legitimate access to and permission to use the data used in this manuscript.

## Abstract

This work contains two data sets.

(1) **COMPAS recidivism data set:**

- Originally collected for the ProPublica story "Machine Bias".
- The features contained in this data set are criminal history, jail and prison time, demographics, and a COMPAS risk scores for defendants from Broward County.
- Goal is to predict whether a given defendant will recidivate, given the current features at hand.

- The original data set can be downloaded from ProPublica's github. https://github.com/propublica/compas-analysis/tree/master

(2) **Adult census income data set:**

- This data set was collected from the 1994 US Census.
- Data set containing demographic information such as age, education, marital status, occupation, etc. . . for individuals in the US from 1994.
- Goal is to predict whether a given person makes above (class 2) or below (class 1) 50K a year.
- The original data set can be downloaded from the UCI machine learning repository. https://archive.ics.uci.edu/dataset/2/adult

## Availability

☒ Data **are** publicly available.
☐ Data **cannot be made** publicly available.

If the data are publicly available, see the *Publicly available data* section. Otherwise, see the *Non-publicly available data* section, below.

## Publicly available data

☒ Data are available online at:

– The COMPAS recidivism data set can be downloaded at https://github.com/propublica/compas-analysis/tree/master.

- The Adult census income data set can be downloaded at https://archive.ics.uci.edu/dataset/2/adult.

☐ Data are available as part of the paper's supplementary material.

☐ Data are publicly available by request, following the process described here:

☐ Data are or will be made available through some other mechanism, described here:

**Non-publicly available data**

## Description

**File format(s)**

☒ CSV or other plain text.
☐ Software-specific binary format (.Rda, Python pickle, etc.):
☐ Standardized binary format (e.g., netCDF, HDF5, etc.):
☐ Other (please specify):

**Data dictionary**

☐ Provided by authors in the following file(s):

☐ Data file(s) is(are) self-describing (e.g., netCDF files)

☒ Available at the following URL:

- COMPAS recidivism data dictionary: https://github.com/propublica/compas-analysis/tree/master
- Adult census income data dictionary: https://archive.ics.uci.edu/dataset/2/adult

**Additional Information (optional)**

# Part 2: Code

## Abstract

Each figure in the paper is accompanied by its own folder containing a function file and a main file. The function file must be executed before running the main file. The main file produces numerical results and generates an accompanying plot. The function file for every figure are nearly identical, with minor modifications implemented to expedite computation in more complex tasks. Although parallel processing is utilized in some files, it is optional; these files can be run without parallel processing by replacing **'mclapply'** with **'lapply'**.

An accompanying R package named **fasi** is available on CRAN, which can be used to implement our algorithm on new data sets. However, please note that this R package was not directly used in the studies presented in the main paper.

## Description

**Code format(s)**

☒ Script files
    ☒ R
    ☐ Python
    ☐ Matlab
    ☐ Other:
☒ Package
    ☒ R

&#9723; Python
&#9723; MATLAB toolbox
&#9723; Other:
&#9723; Reproducible report
   &#9723; R Markdown
   &#9723; Jupyter notebook
   &#9723; Other:
&#9723; Shell script
&#9723; Other (please specify):

## Supporting software requirements

**Version of primary software used**   R version 4.3.1

**Libraries and dependencies used by the code**   There are no packages that are foundational to the implementation of the FASI algorithm. However, since users can estimate confidence scores from any machine learning model, we demonstrated our algorithm using various algorithms. The ones used in our report are,

**ML algorithms for estimating confidence scores**

1. **GAM** - from the **'mgcv'** package on CRAN, version *'1.9.1'*.

2. **adaboost** - from the **'JOUSBoost'** package on CRAN, version *'2.1.0'*.

3. **Logistic Regression** - from the **'MASS'** package on CRAN, version *'7.3.60'*.

4. **Nonparametric Naive Bayes** - from the **'naivebayes'** package on CRAN, version *'1.0.0'*.

5. **XGBoost** - from the **'xgboost'** package on CRAN, version *'1.7.5.1'*.

**Packages for general simulation workflow**

1. **tidyverse** - version *'2.0.0'*, available on CRAN.

2. **ggplot2** - version *'3.4.4'*, available on CRAN.

## Supporting system/hardware requirements (optional)

## Parallelization used

&#9723; No parallel code used
&#8864; Multi-core parallelization on a single machine/node
   &minus; Number of cores used: 10
&#9723; Multi-machine/multi-node parallelization
   &minus; Number of nodes and cores used:

## License

&#8864; MIT License (default)
&#9723; BSD
&#9723; GPL v3.0
&#9723; Creative Commons
&#9723; Other: (please specify)

**Additional information (optional)**

# Part 3: Reproducibility workflow

## Scope

The provided workflow reproduces:

- ☐ Any numbers provided in text in the paper
- ☒ The computational method(s) presented in the paper (i.e., code is provided that implements the method(s))
- ☒ All tables and figures in the paper
- ☐ Selected tables and figures in the paper, as explained and justified below:

## Workflow

### Location

The workflow is available:

- ☐ As part of the paper's supplementary material.
- ☒ In this Git repository: https://github.com/bradleyrava/fasi_experiments
- ☐ Other (please specify):

### Format(s)

- ☒ Single master code file
- ☐ Wrapper (shell) script(s)
- ☐ Self-contained R Markdown file, Jupyter notebook, or other literate programming approach
- ☐ Text file (e.g., a readme-style file) that documents workflow
- ☐ Makefile
- ☐ Other (more detail in *Instructions* below)

### Instructions

### Expected run-time

Approximate time needed to reproduce the analyses on a standard desktop machine:

- ☐ < 1 minute
- ☐ 1-10 minutes
- ☐ 10-60 minutes
- ☒ 1-8 hours
- ☐ > 8 hours
- ☐ Not feasible to run on a desktop machine, as described here:

**Additional information (optional)**

# Notes (optional)