

AI Ethics and Society Final Exam

The test is open-book, open-note, and open-internet: you may consult any materials you want if you do NOT interact live with another human being. This means you may not post about the exam on Ed discussion, text others about the exam, email others about the exam, talk on the phone, or otherwise gain live support from another person by any other means.

Professor Woke decided to create a new course around the theme of Ethical AI. As she browsed the internet for stories about the abuse and misuse of AI, she became saddened by all the problems out there those developers didn't even seem to be aware of. She worried that the tech backlash against technologists because of this misuse would override the potential benefits that AI could have for society. She also worried that developers would not be mindful enough to want to fix the problem even if they heard about this abuse. She wonders if she should give up and take those offered corporate management jobs, which pay much better with fewer worries.

It is your task as a student to convince Professor Woke that, as a technologist, you understand some of these problems with AI misuse and have ideas about how to address some of the corresponding bias and unfairness issues.

Task 1 (25 pts): You must find a **public artifact** (newspaper/magazine article, blog post, YouTube video, etc.) that has been released within the **last six months** (Feb 1, 2022 – Aug 1, 2022). The public artifact must identify some aspect of AI misuse as applied to an application/scenario/domain; the misuse must impact a regulated domain and/or a legally recognized protected class; the public artifact must have associated with it some form of data evidence (a research publication, a released dataset, results from a survey, etc.).

Note: The evidence does not need to have been released within the last six months, only the public artifact.

Provide information related to your artifact: Title of artifact, release date, link to artifact, application/scenario/domain of misuse; regulated domain/protected class impacted; link to evidence. As an example (caveat: this public artifact was released outside of the six months):

- Public Artifact:
 - Title - “Can you make AI fairer than a judge? Play our courtroom algorithm game.”
 - Released - October 17, 2019
 - Link - <https://www.technologyreview.com/s/613508/ai-fairer-than-judge-criminal-risk-assessment-algorithm/>
- Application/Scenario/Domain of Misuse: Criminal Risk Assessment (Predictive Algorithm)
- Regulated Domain/Protected Classes Impacted: Public Accommodation/Race
- Evidence: Dataset - <https://www.propublica.org/datastore/dataset/compas-recidivism-risk-score-data-and-analysis>

Task 2 (20 pts): Provide a 1-2 paragraph summary of the bias identified by the public artifact. In your description, your summary should be specific and reference definitions, concepts, and ideas discussed during this course (i.e., in lectures, assignments, case studies, written critiques, etc.).

Task 3 (30 pts): Provide specific details (in bullet or table format) on all available quantifiable metrics that can be derived from compiling information from the artifact and associated evidence. **There should be enough metrics and details** provided for us to validate your ability to synthesize course concepts based on the overarching topics:

- Privileged/unprivileged groups
- Any misleading graphs?
- Sources of Data Bias
- Sources of Sampling Bias
- Sampling Methods Used to Collect Data
- Correlations found in the data
- Outcome measures: Averages, Standard Deviations, Quartiles, Frequency Distributions, Margins of Error
- Bias & Fairness (or other) metrics used to identify differences in outcomes

Task 4 (25 pts): Identify an issue related to one of the quantifiable metrics listed above (Task 3) that, if addressed, might help mitigate bias and/or unfairness. **Design a method** to help address the issue identified. The method should relate to a **concept discussed in the lectures**. You should explain the method **using a pseudo-code with a 2-4 paragraph summary or a python script with comments**. Remember to identify the issue, **the data inputs and outputs** (based on the evidence), and the **anticipated change** in outcomes.

Note: You do not need a working code or quantitative results.