1 Fairness in Machine Learning

2 # 1 Team members:

3 Atharva Barwe
4 Siddharth Das
5 Bharath kalathuru
6 Jaswanth Gorthi
7 Sagar Rudagi

8 ## 2 Method attributes and presentation:

9 For our project we are considering a binary classification problem, i.e, to predict the annual income of an

individual, if 10 its greater or smaller than 50k. We can interpret our classifier as a random variable by considering $\hat{Y}$ = $f(X)$.
11

12 For initial study we considered the the adult dataset, we first checked whether the data is unbalanced or not. An 13 imbalance in data is dangerous because it can lead to bias. A Machine Learning model may learn the wrong lessons 14 simply because data was not collected thoroughly enough.
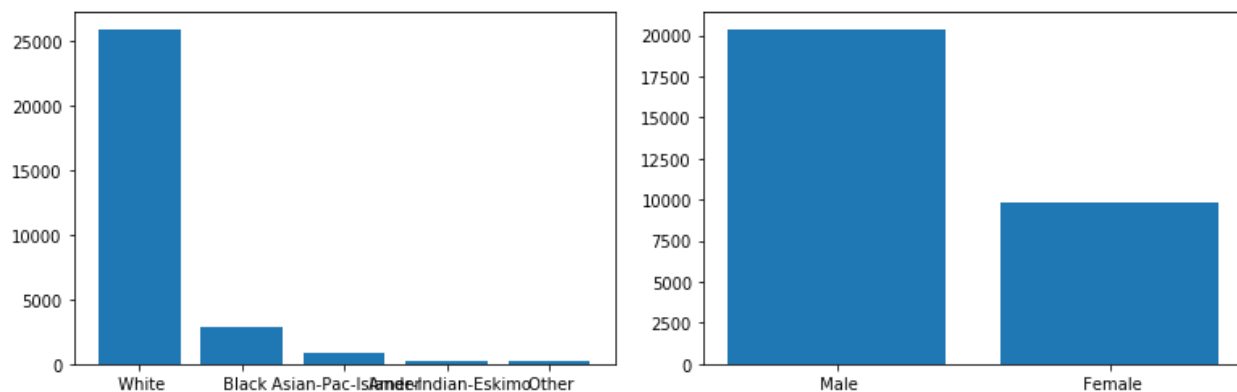15

16 Our dataset contains several protected classes i.e. classes that should not have the outcome but they sometimes do 17 anyway. Race, sex, and native country are the protected classes we have.
18

19 We tested if our dataset for disparate impact. Disparate Impact is a metric to evaluate fairness. It is responsible for 20 comparing the proportion of individuals that have a positive outcome. There are two groups we compare - privileged 21 and unprivileged. Essentially, what we are doing is measuring the success rate of the less favored group against that of 22 the more favored one.
23

24 We can see that for all groups except the Asian Pacific Islander, we have calculated the ratio of probabilities to be lesser 25 than 0.8, which is the industry standard. This means that the unprivileged group receives a positive outcome less than 26 80 percent of their proportion of the privileged group, making it a disparate impact violation.
27

28 This is the formula we use to know if our dataset has disparate impact. If the ratio is lesser than 0.8, it confirms that our 29 dataset has disparate impact -

$$P(class => 50K|X = Black)/P(class => 50K|X = White)$$



30

31 Race - Black , Probability ratio = 0.49266804540332104
32 Race - Amer-Indian-Eskimo , Probability ratio = 0.45078871998012227
33 Race - Asian-Pac-Islander , Probability ratio = 1.0507243618386497

34   Race - Other , Probability ratio = 0.3447207858671523
35

36 These numbers show us that our dataset has disparate impact and we will need to address this issue so that a trained
37 model is not unfairly biased toward or against a certain group.