# After-class exercise for lecture 15

## 36-313, Fall 2022

## Due by 6 pm on Wednesday, 26 October 2022

Consider the following tables which derived from the COMPAS data set. The outcome $Y$ here is being arrested for a violent offense within two years of one's initial arrest. The prediction $\hat{Y}$ was derived by taking the COMPAS score for risk of violence, between 1 and 10, and setting $\hat{Y} = 1$ if the score was $\geq 6$. The "protected attribute" $X$ here is race/ethnicity[1].

```
##            recid
## score_factor    0    1
##    HighRisk  325  208
##    LowRisk  1189  196
```

*Table 1: Black arrestees*

```
##            recid
## score_factor    0    1
##    HighRisk  117   51
##    LowRisk  1168  123
```

*Table 2: White arrestees*

```
##            recid
## score_factor    0   1
##    HighRisk  33   7
##    LowRisk  287  28
```

*Table 3: Hispanic arrestees*

1. (2 points) What percentage of each group is classified as a risk of violence? (In symbols, $\mathbb{P}\left(\hat{Y} = 1 | X = x\right)$.) Does this classifier have demographic parity?

2. (2) For each of the three groups, what's the classification accuracy? (In symbols, $\mathbb{P}\left(Y = \hat{Y} | X = x\right)$.) Does this classifier have parity of classification accuracy?

3. (3) For each of the three groups, what are the false positive and false negative rates? (In symbols, $\mathbb{P}\left(\hat{Y} = 1 | Y = 0, X = x\right)$ and $\mathbb{P}\left(\hat{Y} = 0 | Y = 1, X = x\right)$.) Does this classifier have parity of false positive rates? Of false negative rates?

4. (3) For each of the three groups, what are the positive predictive value and negative predictive value? (In symbols, $\mathbb{P}\left(Y = 1 | \hat{Y} = 1, X = x\right)$ and $\mathbb{P}\left(Y = 0 | \hat{Y} = 0, X = x\right)$.) Does this classifier have parity of positive predictive values? Of negative predictive values?

---

[1]According to the Census Bureau &c., "Hispanic" is about cultural background and "Hispanics can be of any race". In practice, lots of local governments, like this county in Florida, treat "Hispanic" as though it were a race distinct from black, white, Asian, etc.