

# HOMWORK 6: LEARNING THEORY, FAIRNESS METRICS, AND SOCIETAL IMPACT

10-301/10-601 Introduction to Machine Learning (Spring 2024)

<https://www.cs.cmu.edu/~mgormley/courses/10601/>

OUT: Monday, March 18th

DUE: Sunday, March 24th

TAs: Aadit, Annie, Erin, Hailey, Markov

Homework 6 covers topics on Learning Theory, Fairness Metrics, and Societal Impacts. The homework includes multiple choice, True/False, and short answer questions. There will be no consistency points in general, so please make sure to double check your answers to all parts of the questions!

## START HERE: Instructions

- **Collaboration Policy:** Please read the collaboration policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Late Submission Policy:** For this homework, you will only have 2 late days instead of the usual 3. This allows us to provide feedback before the exam. See the late submission policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Submitting your work:** You will use Gradescope to submit answers to all questions and code. Please follow instructions at the end of this PDF to correctly submit all your code to Gradescope.
  - **Written:** For written problems such as short answer, multiple choice, derivations, proofs, or plots, please use the provided template. Submissions can be handwritten onto the template, but should be labeled and clearly legible. If your writing is not legible, you will not be awarded marks. If your scanned submission misaligns the template, there will be a 5% penalty. Alternatively, submissions can be written in LaTeX. Each derivation/proof should be completed in the boxes provided. If you do not follow the template, your assignment may not be graded correctly by our AI assisted grader.

## Instructions for Specific Problem Types

For “Select One” questions, please fill in the appropriate bubble completely:

**Select One:** Who taught this course?

- ☒ Matt Gormley
- ☐ Marie Curie
- ☐ Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

**Select One:** Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☒ Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

**Select all that apply:** Which are instructors for this course?

- ☒ Matt Gormley
- ☒ Henry Chai
- ☒ Hoda Heidari
- ☐ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

**Select all that apply:** Which are the instructors for this course?

- ☒ Matt Gormley
- ☒ Henry Chai
- ☒ Hoda Heidari
- ☒ I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

**Fill in the blank:** What is the course number?

10-601

10-~~6~~301

## Written Questions (68 points)

### 1 L<sup>A</sup>T<sub>E</sub>X Bonus Point and Template Alignment (1 points)

1. (1 point) **Select one:** Did you use L<sup>A</sup>T<sub>E</sub>X for the entire written portion of this homework?

☒ Yes

☐ No

2. (0 points) **Select one:** I have ensured that my final submission is aligned with the original template given to me in the handout file and that I haven't deleted or resized any items or made any other modifications which will result in a misaligned template. I understand that incorrectly responding yes to this question will result in a penalty equivalent to 2% of the points on this assignment.

**Note:** Failing to answer this question will not exempt you from the 2% misalignment penalty.

☒ Yes

### 2 Learning Theory (19 points)

1. Neural the Narwhal is given a classification task to solve, which he decides to use a decision tree learner with 2 binary features  $X_1$  and  $X_2$ . On the other hand, you think that Neural should not have used a decision tree. Instead, you think it would be best to use logistic regression with 16 real-valued features in addition to a bias term. You want to use PAC learning to check whether you are correct. You first train your logistic regression model on  $N$  examples to obtain a training error  $\hat{R}$ .

- (a) (1 point) Which of the following case of PAC learning should you use for your logistic regression model?

☐ Finite and realizable

☐ Finite and agnostic

☐ Infinite and realizable

☒ Infinite and agnostic

- (b) (2 points) What is the upper bound on the true error  $R$  in terms of  $\hat{R}$ ,  $\delta$ , and  $N$ ? You may use big- $\mathcal{O}$  notation if necessary. Write only the final answer. Your work will *not* be graded.

**Note:** Your answer may not contain any other symbols.

Your Answer

for infinite and agnostic case, the true error rate is:

$R(h) \leq \hat{R}(h) + O(\sqrt{\frac{1}{N}(VC(H) + \log \frac{1}{\delta})})$ , since we are using logistic regression with 16 features,  $VC(H) = 17$ . so

$$R(h) \leq \hat{R}(h) + O(\sqrt{\frac{1}{N}(17 + \log \frac{1}{\delta})})$$

(c) (3 points) **Select one:** You want to argue your method has a lower bound on the true error as compared to the Neural's true error bound. Assume that you have obtained enough data points to satisfy the PAC criterion with the same  $\epsilon$  and  $\delta$  as Neural. Which of the following is true?

- ☐ Neural's model will always classify unseen data more accurately because it only needs 2 binary features and therefore is simpler.
- ☐ You must first regularize your model by removing 14 features to make any comparison at all.
- ☐ It is sufficient to show that the VC dimension of your classifier is higher than that of Neural's, therefore having a lower bound for the true error.
- ☒ It is necessary to show that the training error you achieve is lower than the training error Neural achieves.

2. In lecture, we saw that we can use our sample complexity bounds to derive bounds on the true error for a particular algorithm. Consider the sample complexity bound for the infinite, agnostic case:

$$N = O\left(\frac{1}{\epsilon^2} \left[ \text{VC}(\mathcal{H}) + \log \frac{1}{\delta} \right]\right).$$

(a) (2 points) What is the big- $\mathcal{O}$  bound of  $\epsilon$  in terms of  $N$ ,  $\delta$ , and  $\text{VC}(\mathcal{H})$ ?

**Note:**  $A = \mathcal{O}(B)$  (for some value  $B$ )  $\Leftrightarrow$  there exists a constant  $c \in \mathbb{R}$  such that  $A \leq cB$ .

Your Answer

since by definition  $A = \mathcal{O}(B)$  (for some value  $B$ )  $\Leftrightarrow$  there exists a constant  $c \in \mathbb{R}$  such that  $A \leq cB$ , we can choose our constant to be  $c^2$  here. Hence

$$N \leq \frac{c^2}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log \frac{1}{\delta}]$$

$$\epsilon^2 \leq \frac{c^2}{N} [\text{VC}(\mathcal{H}) + \log \frac{1}{\delta}]$$

$$\epsilon \leq \sqrt{\frac{c^2}{N} [\text{VC}(\mathcal{H}) + \log \frac{1}{\delta}]}$$

$$\epsilon \leq c \sqrt{\frac{1}{N} [\text{VC}(\mathcal{H}) + \log \frac{1}{\delta}]}$$

$$\epsilon = O\left(\sqrt{\frac{1}{N} [\text{VC}(\mathcal{H}) + \log \frac{1}{\delta}]}\right)$$

- (b) (2 points) Now, using the definition of  $\epsilon$  (i.e.  $|R(h) - \hat{R}(h)| \leq \epsilon$ ) and your answer to part a, prove that with probability at least  $(1 - \delta)$ :

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N} \left[ \text{VC}(\mathcal{H}) + \log \frac{1}{\delta} \right]}\right).$$

#### Your Answer

since from part (a) we derived that  $\epsilon = O\left(\sqrt{\frac{1}{N} \left[ \text{VC}(\mathcal{H}) + \log \frac{1}{\delta} \right]}\right)$   
 by rearranging the inequality  $|R(h) - \hat{R}(h)| \leq \epsilon$ , we get  $R(h) \leq \hat{R}(h) + \epsilon$   
 hence we have  
 $R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N} \left[ \text{VC}(\mathcal{H}) + \log \frac{1}{\delta} \right]}\right).$

3. (3 points) Consider the hypothesis space of functions that map  $M$  binary attributes to a binary label. A function  $f$  in this space can be characterized as  $f : \{0, 1\}^M \rightarrow \{0, 1\}$ . Neural the Narwhal says that regardless of the value of  $M$ , a hypothesis class containing functions in this space can always shatter  $2^M$  points. Is Neural wrong? If so, provide a counterexample. If Neural is right, briefly explain why in 1-2 *concise* sentences.

#### Your Answer

Yes Neural is Correct, according to definition of shattering, if  $|H(S)| = 2^{|S|}$  then S is shattered by H. The S Narwhal use is  $2^M$  and the total number of dichotomies induced by H need to be  $2^{2^M}$ . since H contains all  $f : \{0, 1\}^M \rightarrow \{0, 1\}$ , we have total of  $2^M$  possible inputs and  $2^{2^M}$  total unique boolean formula. Hence we can indeed shatter  $2^M$  points

4. Consider an instance space  $\mathcal{X}$  which is the set of real numbers.

(a) (3 points) **Select one:** What is the VC dimension of hypothesis class  $H$ , where each hypothesis  $h$  in  $H$  is of the form “if  $a < x < b$  or  $c < x < d$  then  $y = 1$ ; otherwise  $y = 0$ ”? (i.e.,  $H$  is an infinite hypothesis class where  $a, b, c$ , and  $d$  are arbitrary real numbers).

☐ 2

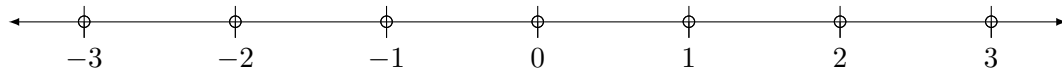
☐ 3

☒ 4

☐ 5

☐ 6

(b) (3 points) Given the set of points in  $\mathcal{X}$  below, construct a labeling of some subset of the points to show that any dimension larger than the VC dimension of  $H$  by *exactly* 1 is incorrect (e.g. if the VC dimension of  $H$  is 3, only fill in the answers for 4 of the points). Fill in the boxes such that for each point in your example, the corresponding label is either 0 or 1. For points you are not using in your example, write N/A (do *not* leave the answer box blank).



Answer for -3 1	Answer for -2 0	Answer for -1 1	
Answer for 0 0	Answer for 1 1	Answer for 2 N/A	Answer for 3 N/A

### 3 Fairness Metrics (21 points)

Neural works for the Bank of ML and is given the following dataset from another bank on whether or not to issue a loan to individuals. Each row in this dataset represents one individual's data, which includes their FICO credit score, their savings rate (percentage of their income that goes into their savings), and credit history in months. The data was collected in two different cities, city A and city B, as denoted in the first column. The "Label" column refers to the true label, where "1" refers to loan issued, and "0" refers to no loan issued. A csv file of this dataset could be found in the handout folder.

Region	FICO Score	Savings Rate (%)	Credit History (months)	Label
A	544.0625	28.0	21	1
A	489.0625	33.9	40	0
A	433.125	62.3	100	0
A	429.0625	56.7	203	1
A	417.8125	56.5	5	0
A	506.5625	32.7	75	1
A	400.625	60.7	216	0
A	836.875	10.7	86	1
A	471.875	36.2	92	1
A	402.8125	62.0	199	0
B	809.4285714	5.6	213	1
B	480.9375	40.2	72	1
B	505.0	31.1	20	0
B	438.4375	51.3	122	0
B	385.9375	76.2	89	0
B	505.625	34.7	39	1
B	514.0625	31.0	41	1
B	385.9375	76.2	89	0
B	446.25	44.5	51	0
B	428.75	55.6	215	1

1. Neural took the average value of the features (for example, the average value for the first data point is 197.69), and developed the following observation. In general, for all three features in this dataset, a high value indicates better credibility. Hence Neural trained the following decision stump on this dataset: if the average feature value is above the median (198.09), then we determine that the individual will receive the loan (prediction = 1). Otherwise, we decide that the individual will not receive the loan. **For parts (a), (b), (c) below, please round your answer to three decimal places.**

(a) (1 point) Using the model that Neural proposed, what is the training error rate on the entire dataset?

Your Answer

0.400

(b) (1 point) What is the training error rate for region A?

Your Answer
0.400

(c) (1 point) What is the training error rate for region B?

Your Answer
0.400

(d) (1 point) How many false positives were there in region A?

Your Answer
3

(e) (1 point) How many false negatives were there in region A?

Your Answer
1

(f) (1 point) How many false positives were there in region B?

Your Answer
1

(g) (1 point) How many false negatives were there in region B?

Your Answer
3



2. (2 points) **True or False:** Using your responses to the previous question, we achieve statistical parity between regions A and B. Justify your answer.

☐ True  
☒ False

Your Answer

for A the selection rate is  $\frac{3+4}{10} = 0.7$ , for B the selection rate is  $\frac{1+2}{10} = 0.3$ , hence we don't have statistical parity between region A and B.

3. (2 points) **True or False:** We achieve equality of accuracy between regions A and B. Justify your answer.

☒ True  
☐ False

Your Answer

the accuracy of region A and B are both 0.6, since these two numbers are equal, there is equality of accuracy between A and B.

4. (2 points) **True or False:** We achieve equality of FPR/FNR between regions A and B. Justify your answer.

☐ True  
☒ False

Your Answer

the FPR/FNR of A is  $0.6/0.2 = 3$   
the FPR/FNR of B is  $0.2/0.6 = 1/3$ .  
since the two ratios are not the same, we don't have equality of  $FPR/FNR$  between regions A and B

5. (2 points) **True or False:** We achieve equality of PPV/NPV between regions A and B. Justify your answer.

☐ True  
☒ False

Your Answer

the PPV/NPV of A is  $0.571/0.667 = 0.856(6/7)$   
the PPV/NPV of B is  $0.667/0.571 = 1.168(7/6)$   
since the two ratio are not the same, we do not have equality of PPV/NPV between region A and B

6. (3 points) Using your responses from the previous questions, comment on the fairness of this model between cities A and B.

Your Answer

we don't have statistical parity between A and B, so there is no equal treatment in terms of label for the 2 region regardless of the correctness of the labeling, yet we have the same accuracy for our predictions. Region A has FPR/FNR ratio  $> 1$ , we produce more false positives than false negative. it is acceptable when false negative are problematic. Region B has FPR/FNR ratio  $< 1$ , it produce more false negatives. this is acceptable when false positive is more harmful. PPV/NPV tell us how reliable the positive and negative predictions are and region A is more reliable on negative prediction and B is more reliable on positive one.

7. (3 points) A Type I error occurs when you erroneously predict a positive label (false positive), and a Type II error is when you erroneously predict a negative label (false negative). Compare and contrast the consequences of making a Type I error and Type II error in this setting. Which would cause more significant consequences?

Your Answer

if we make type I error, we predict someone should be issued loan while they shouldn't. if we make type II error, we predict someone shouldn't get loan while they should. In the first case, the consequence might be lending money to someone who cannot pay it back, in the second case, we might loss potential revenue because we are not lending to someone who is capable to payback the money. From the bank's perspective, type I error has more significant consequences; from borrower's perspective, type II is more significant.

## 4 Societal Impacts (27 points)

The fictional country, Xtopia, is in the midst of an epidemic. The Xtopian healthcare system has been under a great deal of strain in the past year due to a regional epidemic caused by an airborne virus called Xvid. The number of hospital beds is limited and as the result, healthcare professionals have to frequently make very difficult choices about which subset of Xvid patients can be hospitalized. Hospital care greatly increases the chance of recovering from the illness with no subsequent long-term health complications.

To save time and make these decisions more efficient and consistent, a team of ML practitioners have been brought in to automate the decision-making process. They have been given access to a data set consisting of the information about prior Xvid patients who sought hospital care along with the binary decision made about them by the hospital doctors ('+' indicates hospitalization and '-' indicates no hospitalization). The ML team has determined that the decision about each patient is highly correlated with his/her age as well as his/her prior utilization of medical insurance. This observation reflects the fact that Xtopian doctors are on average more likely to allocate scarce medical resources to the young and the vulnerable (i.e., those with prior medical conditions and comorbidities). Here, the insurance utilization serves as a proxy for severity of the patient's health conditions.

Figure 1 provides a snapshot of the Xvid training data and the predictive model that the ML team has come up with. Each instance corresponds to an individual patient, and each patient belongs to one of the two socially salient groups in Xtopia, indicated by blue and red.

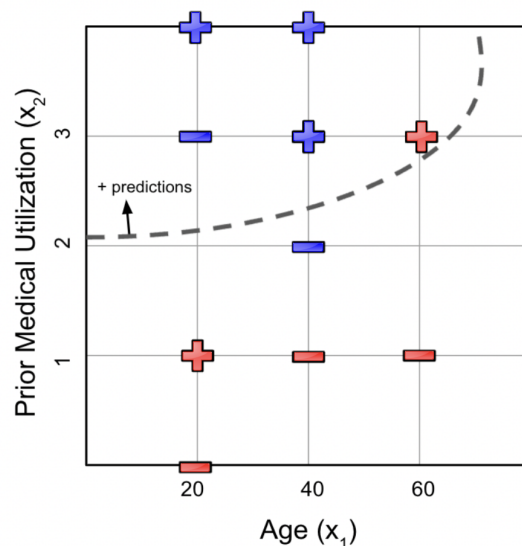


Figure 1: Xvid Training Data

Answer the following questions with respect to the above hypothetical context and data set.

1. (2 points) **Select all that apply:** Does the predictive model above satisfy the following notions of fairness across blue and red groups?

- ☐ False Negative Rate (FNR) parity
- ☐ False Positive Rate (FPR) parity
- ☐ Negative Predictive Value (NPV) parity
- ☐ Positive Predictive Value (PPV) parity
- ☒ Error parity
- ☐ Statistical parity (or Selection rate parity)
- ☐ None of the above

2. (2 points) **Select all that apply:** Which of the above notions of fairness would be satisfied if the ML team could train a model with 0 true error (i.e., a model that always predicts the correct label for every patient)?

- ☒ False Negative Rate (FNR) parity
- ☒ False Positive Rate (FPR) parity
- ☒ Negative Predictive Value (NPV) parity
- ☒ Positive Predictive Value (PPV) parity
- ☒ Error parity
- ☐ Statistical parity (or Selection rate parity)
- ☐ None of the above

3. (2 points) **Select all that apply:** Which of the above notions of fairness would be satisfied in expectation by a random classifier (i.e., a model that makes a randomized prediction for every patient: with probability 0.5 the patient is hospitalized regardless of their attributes)?

- ☐ False Negative Rate (FNR) parity
- ☐ False Positive Rate (FPR) parity
- ☐ Negative Predictive Value (NPV) parity
- ☐ Positive Predictive Value (PPV) parity
- ☐ Error parity
- ☒ Statistical parity (or Selection rate parity)
- ☐ None of the above

4. (2 points) From the perspective of a patient subject to the predictions made by this model, the violation of which of the parity conditions below would be most problematic? Justify your answer.

- ☒ False Negative Rate (FNR) parity
- ☐ False Positive Rate (FPR) parity
- ☐ Negative Predictive Value (NPV) parity
- ☐ Positive Predictive Value (PPV) parity
- ☐ Error parity
- ☐ Statistical parity (or Selection rate parity)
- ☐ None of the above

Your Answer

for patient, false negative parity matters the most because high FNR indicate that people in a group is more likely to miss out beneficial opportunity, in this case, getting medical resources.

5. (2 points) **Causes of unfairness:** Name one potential cause of disparity in false negative rates across the two groups in the above context.

Your Answer

from the plot above, the person that is false negative has age 20 and prior medical utilization of 1. I think potential cause of false negative is that this person has low prior medical utilization, meaning they don't have much prior health condition.

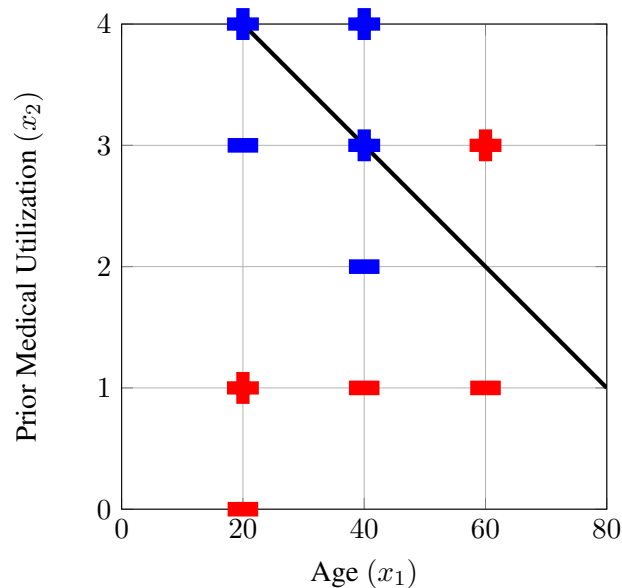
6. **Fairness interventions:** consider the following pre-processing method to improve statistical parity:

While the selection rate is unequal across the two groups:

- (a) Pick the group with lowest selection rate.
- (b) From this group in the training data, pick the data point closest to the decision boundary predicted as negative.
- (c) Change the label of this instance to positive.
- (d) Retrain the model on the modified training data by finding the highest accuracy classifier in the hypothesis class.

Suppose our hypothesis class is the class of all linear separators defined over  $\mathbb{R}^2$ .

- (a) (1 point) The highest accuracy linear separator is shown in the figure below (assume that points on the decision boundary are characterized as '+'):

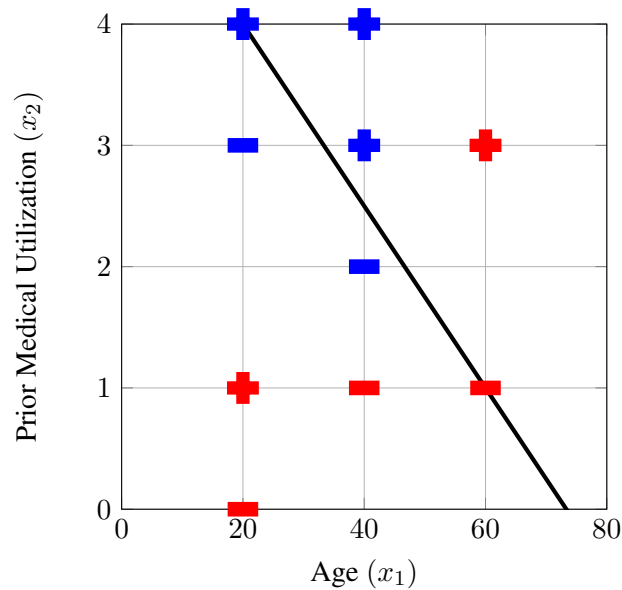


What are the coordinates of the data point whose label would be flipped first by this pre-processing method?

Your Answer

(60, 1)

- (b) (2 points) After flipping the label of the point you identified in the previous question, plot the linear separator with the highest accuracy. For your convenience, we have provided a mechanism for you to input your answer by specifying the coordinates of two points on the decision boundary.



- (c) (1 point) Using the linear decision boundary you plotted in the previous question, what are the coordinates of the data point whose label would be flipped next by this pre-processing method?

Your Answer

(40, 1)

- (d) (1 point) **True or False:** The algorithm terminates at this point.

☐ True

☒ False

7. **The fairness impossibility theorem:** Prove by contradiction that if prevalence rate,  $r_s = P[Y = 1|S = s]$  across the two groups  $s \in \{\text{blue}, \text{red}\}$  is different, then there does not exist a classifier that can satisfy PPV parity, FPR parity, and FNR parity simultaneously.

- (a) (4 points) Verify that the following identity holds for any  $s \in \{\text{blue}, \text{red}\}$ :

$$FPR_s = \frac{r_s}{1 - r_s} \times (1 - FNR_s) \times \frac{(1 - PPV_s)}{PPV_s}$$

Your Answer

$$\begin{aligned}
 r_s &= \frac{TP+FN}{(TP+FN+TN+FP)}, 1-r_s = \frac{TN+FP}{TP+FN+TN+FP} \\
 \frac{r_s}{1-r_s} &= \frac{TP+FN}{TN+FP} \\
 FNR_s &= \frac{FN}{FN+TP} \\
 1-FNR_s &= \frac{TP}{FN+TP} \\
 PPV_s &= \frac{TP}{TP+FP} \\
 1-PPV_s &= \frac{FP}{TP+FP} \\
 \frac{1-PPV_s}{PPV_s} &= \frac{FP}{TP} \\
 \frac{r_s}{1-r_s} \times (1-FNR_s) \times \frac{(1-PPV_s)}{PPV_s} &= \frac{TP+FN}{TN+FP} \times \frac{TP}{FN+TP} \times \frac{FP}{TP} = \frac{FP}{TN+FP} = FPR_s
 \end{aligned}$$

hence the above identity holds for  $s \in \{blue, red\}$

(b) (4 points) Show that the expression from part (a) can be rewritten as

$$1/r_s = 1 + \frac{(1-FNR_s)}{FPR_s} \times \frac{(1-PPV_s)}{PPV_s}$$

Your Answer

$$\begin{aligned}
 &\text{since we have } FPR_s = \frac{r_s}{1-r_s} \times (1-FNR_s) \times \frac{(1-PPV_s)}{PPV_s} \\
 1 &= \frac{r_s}{1-r_s} \times \frac{(1-FNR_s)}{FPR_s} \times \frac{(1-PPV_s)}{PPV_s} \\
 \frac{1-r_s}{r_s} &= \frac{(1-FNR_s)}{FPR_s} \times \frac{(1-PPV_s)}{PPV_s} \\
 \frac{1}{r_s} - 1 &= \frac{(1-FNR_s)}{FPR_s} \times \frac{(1-PPV_s)}{PPV_s} \\
 1/r_s &= 1 + \frac{(1-FNR_s)}{FPR_s} \times \frac{(1-PPV_s)}{PPV_s} \quad Q.E.D
 \end{aligned}$$



- (c) (4 points) Finally, using results from parts (a) and (b), show that if  $FPR_s, FNR_s$ , and  $PPV_s$  are equal for  $s \in \{\text{blue}, \text{red}\}$ , then  $r_s$  must be equal for  $s \in \{\text{blue}, \text{red}\}$ , which is a contradiction.

Your Answer

from part *a* and part *b*, we have:

$$1/r_r = 1 + \frac{(1-FNR_r)}{FPR_r} \times \frac{(1-PPV_r)}{PPV_r}$$

$$1/r_b = 1 + \frac{(1-FNR_b)}{FPR_b} \times \frac{(1-PPV_b)}{PPV_b}$$

if  $FPR_s, FNR_s$ , and  $PPV_s$  are equal for  $s \in \{\text{blue}, \text{red}\}$ , then

$$\frac{(1-FNR_r)}{FPR_r} \times \frac{(1-PPV_r)}{PPV_r} = \frac{(1-FNR_b)}{FPR_b} \times \frac{(1-PPV_b)}{PPV_b}$$

hence  $1/r_r = 1/r_b$ ,  $r_r = r_b$ , which is a contradiction.

## 5 Collaboration Questions

After you have completed all other components of this assignment, report your answers to these questions regarding the collaboration policy. Details of the policy can be found [here](#).

1. Did you receive any help whatsoever from anyone in solving this assignment? If so, include full details.
2. Did you give any help whatsoever to anyone in solving this assignment? If so, include full details.
3. Did you find or come across code that implements any part of this assignment? If so, include full details.

### Your Answer

1. No
2. No
3. No