

## 526 A Appendix

### 527 A.1 Proofs

#### 528 A.1.1 Explanation of Assumption 1

529 Assumption 1: For  $i, j \in \{0, 1\}$ , the classifiers  $\hat{Y} = f(X)$  and  $\hat{A} = h(X)$  satisfy

$$\frac{U_i}{\hat{r}_{i,j}} \leq \hat{\alpha}_{i,j} \leq 1 - \frac{U_i}{\hat{r}_{i,j}}.$$

530 We now expand on the implications of this assumptions. Recall that  $U_i = \mathbb{P}(\hat{A} = i, A \neq i)$ ,  
 531  $\hat{r}_{i,j} = \mathbb{P}(\hat{A} = i, Y = j)$ , and  $\hat{\alpha}_{i,j} = \mathbb{P}(\hat{Y} = 1 \mid \hat{A} = i, Y = j)$ . Thus Assumption 1 states,

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = j) \leq \mathbb{P}(\hat{A} = i, Y = j) - \mathbb{P}(\hat{A} = i, A \neq i). \quad (10)$$

532 Immediately, it is clear that Eq. (10) implies

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{A} = i, Y = j) - \mathbb{P}(\hat{A} = i, A \neq i) \quad (11)$$

$$\implies \mathbb{P}(\hat{A} = i, A \neq i) \leq \frac{1}{2} \mathbb{P}(\hat{A} = i, Y = j) \quad (12)$$

533 Now, the left inequality of Eq. (10) states

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = j) \quad (13)$$

534 and the right inequality of Eq. (10) states

$$\mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = j) \leq \mathbb{P}(\hat{A} = i, Y = j) - \mathbb{P}(\hat{A} = i, A \neq i)$$

535 which implies

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{A} = i, Y = j) - \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = j) \quad (14)$$

$$= \mathbb{P}(\hat{Y} = 0, \hat{A} = i, Y = j) \quad (15)$$

536 If  $j = 1$ , then this implies

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = 1) \quad (16)$$

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 0, \hat{A} = i, Y = 1) \quad (17)$$

537 and if  $j = 0$ ,

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = 0) \quad (18)$$

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = 0, \hat{A} = i, Y = 0) \quad (19)$$

538 Any reasonable classifier  $\hat{Y}$  would have the properties

$$\mathbb{P}(\hat{Y} = 0, \hat{A} = i, Y = 1) \leq \mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = 1) \quad (20)$$

$$\mathbb{P}(\hat{Y} = 1, \hat{A} = i, Y = 0) \leq \mathbb{P}(\hat{Y} = 0, \hat{A} = i, Y = 0) \quad (21)$$

539 Thus, Assumption 1 is met when

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = j, \hat{A} = i, Y \neq j) \quad (22)$$

### 540 A.1.2 Proof of Theorem 1

541 We only prove the result for  $|\Delta_{\text{TPR}}(f)|$  as the proof for  $|\Delta_{\text{FPR}}(f)|$  is completely analogous.

542 *Proof.* The rules of conditional probability and the law of total probability allow us to decompose  
 543  $\alpha_{1,1}$  and  $\alpha_{0,1}$  in the following manner,

$$\alpha_{1,1} = \mathbb{P}(\hat{Y} = 1 \mid A = 1, Y = 1) \quad (23)$$

$$= \frac{\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1)}{\mathbb{P}(A = 1, Y = 1)} \quad (24)$$

$$= \frac{\sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = 1, A = 1, Y = 1, \hat{A} = i)}{\sum_{j \in \{0,1\}} \sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = j, A = 1, Y = 1, \hat{A} = i)} \quad (25)$$

$$= \frac{\sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = i) \cdot \mathbb{P}(\hat{A} = i)}{\sum_{j \in \{0,1\}} \sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = j, A = 1, Y = 1 \mid \hat{A} = i) \cdot \mathbb{P}(\hat{A} = i)} \quad (26)$$

544 and

$$\alpha_{0,1} = \mathbb{P}(\hat{Y} = 1 \mid A = 0, Y = 1) \quad (27)$$

$$= \frac{\mathbb{P}(\hat{Y} = 1, A = 0, Y = 1)}{\mathbb{P}(A = 0, Y = 1)} \quad (28)$$

$$= \frac{\mathbb{P}(\hat{Y} = 1, Y = 1) - \mathbb{P}(\hat{Y} = 1, A = 1, Y = 1)}{\mathbb{P}(Y = 1) - \mathbb{P}(A = 1, Y = 1)} \quad (29)$$

$$= \frac{\mathbb{P}(\hat{Y} = 1, Y = 1) - \left[ \sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = i) \cdot \mathbb{P}(\hat{A} = i) \right]}{\mathbb{P}(Y = 1) - \left[ \sum_{j \in \{0,1\}} \sum_{i \in \{0,1\}} \mathbb{P}(\hat{Y} = j, A = 1, Y = 1 \mid \hat{A} = i) \cdot \mathbb{P}(\hat{A} = i) \right]} \quad (30)$$

545 Therefore,  $\Delta_{\text{TPR}}(f) = \alpha_{1,1} - \alpha_{0,1}$  is a function of the four probabilities given by

$$\mathbb{P}(\hat{Y} = j, A = 1, Y = 1 \mid \hat{A} = i) \quad (31)$$

546 which are unidentifiable in a demographically scarce regime and therefore not computable.

547

548 The Fréchet inequalities tell us that for  $i, j \in \{0, 1\}$

$$\mathbb{P}(\hat{Y} = j, A = 1, Y = 1 \mid \hat{A} = i) \geq \max\{\mathbb{P}(\hat{Y} = j, Y = 1 \mid \hat{A} = i) - \mathbb{P}(A = 0 \mid \hat{A} = i), 0\} \quad (32)$$

$$\mathbb{P}(\hat{Y} = j, A = 1, Y = 1 \mid \hat{A} = i) \leq \min\{\mathbb{P}(\hat{Y} = j, Y = 1 \mid \hat{A} = i), \mathbb{P}(A = 1 \mid \hat{A} = i)\}. \quad (33)$$

549 Observe that,  $\Delta_{\text{TPR}}(f)$  is an increasing function with respect to the two probabilities

$$\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = i) \quad (34)$$

550 and a decreasing one with respect to the two probabilities,

$$\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1 \mid \hat{A} = i). \quad (35)$$

551 As a result,  $\Delta_{\text{TPR}}(f)$  is maximal when  $\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = i)$  achieve their maximum  
 552 values and  $\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1 \mid \hat{A} = i)$  achieve their minimum values. On the other hand,  
 553  $\Delta_{\text{TPR}}(f)$  is minimal when  $\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = i)$  achieve their minimum values and  
 554  $\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1, \mid \hat{A} = i)$  achieve their maximum values. With these facts, we now provide  
 555 the upper bound. Recall from Appendix A.1.1 that Assumption 1 implies

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \frac{1}{2} \mathbb{P}(\hat{A} = i, Y = j) \quad (36)$$

$$\mathbb{P}(\hat{A} = i, A \neq i) \leq \mathbb{P}(\hat{Y} = j, \hat{A} = i, Y \neq j) \leq \mathbb{P}(\hat{Y} = j, \hat{A} = i, Y = j) \quad (37)$$

$$(38)$$

556 With Assumption 1 we first provide the values of  $\min [\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1, | \hat{A} = i)]$ . First,

$$\mathbb{P}(\hat{Y} = 0, Y = 1, | \hat{A} = 1) - \mathbb{P}(A = 0 | \hat{A} = 1) = \frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 1, Y = 1)}{P(\hat{A} = 1)} - \frac{U_1}{\mathbb{P}(\hat{A} = 1)} \quad (39)$$

$$\geq 0 \quad (40)$$

557 because  $\mathbb{P}(\hat{Y} = 0, \hat{A} = 1, Y = 1) - U_1 \geq 0$ . Second,

$$\mathbb{P}(\hat{Y} = 0, Y = 1, | \hat{A} = 0) - \mathbb{P}(A = 0 | \hat{A} = 0) = \frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} - \frac{\mathbb{P}(A = 0, \hat{A} = 0)}{\mathbb{P}(\hat{A} = 0)} \quad (41)$$

$$= \frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} - \frac{\mathbb{P}(\hat{A} = 0) - U_0}{\mathbb{P}(\hat{A} = 0)} \quad (42)$$

$$= \frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} - \frac{\mathbb{P}(\hat{A} = 0) - U_0}{\mathbb{P}(\hat{A} = 0)} \quad (43)$$

$$= \frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} + \frac{U_0}{\mathbb{P}(\hat{A} = 0)} - 1 \quad (44)$$

558 Now note that,

$$\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1) = \mathbb{P}(\hat{A} = 0, Y = 1) - \mathbb{P}(\hat{Y} = 1, \hat{A} = 0, Y = 1) \quad (45)$$

$$\leq \mathbb{P}(\hat{A} = 0, Y = 1) - U_0 \quad (46)$$

559 where the second equality is due to Assumption 1. As a result

$$\frac{\mathbb{P}(\hat{Y} = 0, \hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} + \frac{U_0}{\mathbb{P}(\hat{A} = 0)} - 1 \leq \frac{\mathbb{P}(\hat{A} = 0, Y = 1) - U_0}{P(\hat{A} = 0)} + \frac{U_0}{\mathbb{P}(\hat{A} = 0)} - 1 \quad (47)$$

$$= \frac{\mathbb{P}(\hat{A} = 0, Y = 1)}{P(\hat{A} = 0)} - 1 \leq 0 \quad (48)$$

560 Therefore,

$$\min [\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1 | \hat{A} = 1)] = \mathbb{P}(\hat{Y} = 0, Y = 1, | \hat{A} = 1) - \mathbb{P}(A = 0 | \hat{A} = 1) \quad (49)$$

$$\min [\mathbb{P}(\hat{Y} = 0, A = 1, Y = 1 | \hat{A} = 0)] = 0 \quad (50)$$

561 Now we provide the values of  $\max [\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1, | \hat{A} = i)]$ . First,

$$\mathbb{P}(A = 1 | \hat{A} = 1) = \frac{\mathbb{P}(A = 1, \hat{A} = 1)}{\mathbb{P}(\hat{A} = 1)} \quad (51)$$

$$= \frac{\mathbb{P}(\hat{A} = 1) - U_1}{\mathbb{P}(\hat{A} = 1)} \quad (52)$$

$$\geq \frac{\mathbb{P}(\hat{A} = 1, Y = 1) - U_1}{\mathbb{P}(\hat{A} = 1)} \quad (53)$$

$$\geq \frac{\mathbb{P}(\hat{A} = 1, Y = 1) - \mathbb{P}(\hat{Y} = 0, \hat{A} = 1, Y = 1)}{\mathbb{P}(\hat{A} = 1)} \quad (54)$$

$$= \frac{\mathbb{P}(\hat{Y} = 1, \hat{A} = 1, Y = 1)}{\mathbb{P}(\hat{A} = 1)} = \mathbb{P}(\hat{Y} = 1, Y = 1 | \hat{A} = 1) \quad (55)$$

562 Second,

$$\mathbb{P}(A = 1 | \hat{A} = 0) = \frac{\mathbb{P}(A = 1, \hat{A} = 0)}{\mathbb{P}(\hat{A} = 0)} \quad (56)$$

$$\leq \frac{\mathbb{P}(\hat{Y} = 1, \hat{A} = 0, Y = 1)}{\mathbb{P}(\hat{A} = 0)} = \mathbb{P}(\hat{Y} = 1, Y = 1 | \hat{A} = 0) \quad (57)$$

563 Therefore,

$$\max [\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = 1)] = \mathbb{P}(\hat{Y} = 1, Y = 1 \mid \hat{A} = 1) \quad (58)$$

$$\max [\mathbb{P}(\hat{Y} = 1, A = 1, Y = 1 \mid \hat{A} = 0)] = \mathbb{P}(A = 1 \mid \hat{A} = 0) \quad (59)$$

564 Plugging these 4 values into  $\Delta_{\text{TPR}}$  will yield the upper bound,

$$B_1 + C_{0,1} = \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \hat{\alpha}_{1,1} - \frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \hat{\alpha}_{0,1} + U_0 \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right) \quad (60)$$

565 One can similarly use the assumptions to derive the lower bound,

$$B_1 - C_{1,1} = \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \hat{\alpha}_{1,1} - \frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \hat{\alpha}_{0,1} - U_1 \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right) \quad (61)$$

566 and thus  $|\Delta_{\text{TPR}}| \leq \max\{|B_1 + C_{0,1}|, |B_1 - C_{1,1}|\}$ . One can use same arguments to derive the upper  
 567 bound for  $|\Delta_{\text{FPR}}|$ .  $\square$

### 568 A.1.3 Proof of Theorem 2

569 We prove the result for  $|\Delta_{\text{TPR}}(f)|$ . We first start by proving the existence part of the theorem.

570

571 Let  $\hat{A} = h(X)$  be a sensitive attribute classifier with errors  $U_0$  and  $U_1$  that produces rates  $\hat{r}_{i,j} =$   
 572  $\mathbb{P}(\hat{A} = i, Y = j)$ . Let  $\mathcal{F}$  be the set of classifiers for  $Y$  such that  $\forall f \in \mathcal{F}$ ,  $f$  and  $h$  satisfy  
 573 Assumption 1. Consider any  $f \in \mathcal{F}$  with group conditional probabilities,  $\hat{\alpha}_{i,j} = \mathbb{P}(\hat{Y} = 1 \mid \hat{A} =$   
 574  $i, Y = j)$ . Since we are only proving the result for  $|\Delta_{\text{TPR}}(f)|$ , set  $j = 1$ . Consider the  $xy$  plane, with  
 575 the  $x$ -axis being  $\hat{\alpha}_{0,1}$  and the  $y$ -axis being  $\hat{\alpha}_{1,1}$ . We know,

$$\frac{U_i}{\hat{r}_{i,1}} \leq \hat{\alpha}_{i,1} \leq 1 - \frac{U_i}{\hat{r}_{i,1}} \quad (62)$$

576 which implies

$$\frac{U_i}{\hat{r}_{i,1}} \leq \frac{1}{2}. \quad (63)$$

577 The two equations above define a rectangular region in the  $xy$  plane with a center  $(\frac{1}{2}, \frac{1}{2})$ , meaning  
 578 any classifier  $f \in \mathcal{F}$ , has  $\hat{\alpha}_{i,j}$  that are in this region.

579

580 Now, denote  $\bar{\mathcal{F}}$  to be a the set of classifiers for  $Y$ , with group conditional probabilities  $\hat{\alpha}_{i,1}$ , that  
 581 satisfy the condition,

$$\frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \hat{\alpha}_{0,1} - \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \hat{\alpha}_{1,1} = \frac{\Delta U}{2} \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right). \quad (64)$$

582 This condition defines a line in the  $xy$  plane meaning any classifier in  $\bar{\mathcal{F}}$  has  $\hat{\alpha}_{i,1}$  that are on this line.

583 Now observe that the classifier  $\bar{f} \in \bar{\mathcal{F}}$  with  $\hat{\alpha}_{i,1} = \frac{1}{2}$ , satisfy the above condition because,

$$\frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \left( \frac{1}{2} \right) - \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \left( \frac{1}{2} \right) = \frac{1}{2} \left( \frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} - \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \right) \quad (65)$$

$$= \frac{1}{2} \left( \frac{\hat{r}_{0,1} - \Delta U + \Delta U}{\hat{r}_{0,1} - \Delta U} - \frac{\hat{r}_{1,1} + \Delta U - \Delta U}{\hat{r}_{1,1} + \Delta U} \right) \quad (66)$$

$$= \frac{1}{2} \left( 1 + \frac{\Delta U}{\hat{r}_{0,1} - \Delta U} - 1 + \frac{\Delta U}{\hat{r}_{1,1} + \Delta U} \right) \quad (67)$$

$$= \frac{\Delta U}{2} \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right) \quad (68)$$

584 This implies that the line defined by Eq. (64) intersects the rectangular region that Assumption 1  
 585 defines. As a result,  $\mathcal{F} \cap \bar{\mathcal{F}}$  is not empty, meaning there exists a classifier  $\bar{f} \in \mathcal{F}$  with group  
 586 conditional probabilities  $\hat{\alpha}_{i,1}$  that also satisfies the condition,

$$\frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \hat{\alpha}_{0,1} - \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \hat{\alpha}_{1,1} = \frac{\Delta U}{2} \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right). \quad (69)$$

587 Now we prove that such a classifier has minimal bounds. Theorem 1 tells us that for  $f \in \mathcal{F}$

$$|\Delta_{\text{TPR}}(f)| \leq B_{\text{TPR}}(f) \triangleq \max\{|B_1 + C_{0,1}|, |B_1 - C_{1,1}|\}$$

588 Note that  $B_1$  is linear in  $\hat{\alpha}_{1,1}$  and  $\hat{\alpha}_{0,1}$  and that  $C_{0,1}$  and  $C_{1,1}$  are constants such that  $B_1 + C_{0,1} \geq$   
 589  $B_1 - C_{1,1}$  simply because  $B_1 + C_{0,1}$  is the upper bound for  $\Delta_{\text{TPR}}$  and  $B_1 - C_{0,1}$  is the lower bound.  
 590 Since these bounds are affine functions shifted by a constant, then  $\min \max\{|B_1 + C_{0,1}|, |B_1 - C_{1,1}|\}$   
 591 necessarily occurs when

$$B_1 + C_{0,1} = -B_1 - C_{1,1} \quad (70)$$

592 meaning

$$2B_1 = -(C_{1,1} + C_{0,1}) \quad (71)$$

593 have minimal upper bounds on  $|\Delta_{\text{TPR}}|$ . After rearranging terms, this condition is precisely

$$\frac{\hat{r}_{0,1}}{\hat{r}_{0,1} - \Delta U} \hat{\alpha}_{0,1} - \frac{\hat{r}_{1,1}}{\hat{r}_{1,1} + \Delta U} \hat{\alpha}_{1,1} = \frac{\Delta U}{2} \left( \frac{1}{\hat{r}_{1,1} + \Delta U} + \frac{1}{\hat{r}_{0,1} - \Delta U} \right). \quad (72)$$