

1 A Algorithm Details

2 We summarize our method in Algorithm 1 as follows:

Algorithm 1: FCFL : obtain a **fair Pareto min-max** model

Input: fairness budget $\{\epsilon_i\}$, gradient threshold ϵ_d , the initial parameters δ_l^0, δ_g^0 , the learning rate η , the parameter decay rate β , the initial hypothesis h_{θ^0} , and the data of all clients $\{D^1, D^2, \dots, D^N\}$.

Stage 1: obtain a fair min-max model h^1

for $t = 1, 2, \dots, T^1$ **do**

Local clients:

 The center server sends the global model h^t to all local client;

 Local clients evaluate the performance of h^t and obtain the local gradients $\nabla_{\theta_i} l_i$;

Center server:

 Calculate the surrogate maximum objective \hat{l} and disparity \hat{g} using the Eq.(10) in the main text

if $\hat{g} \leq 0$ **then**

 Determine the gradient direction d^t by solving the LP problem in the Eq.(12) in the main text;

end

else

 Determine d^t by solving the problem in the Eq.(13) in the main text;

end

 Update the model $\theta^{t+1} = \theta^t + \eta d^t$;

if $\|d^t\| \leq \epsilon_d$ **then**

 Decay the parameter δ_l^t and δ_g^t : $\delta_l^{t+1} = \beta \cdot \delta_l^t$, $\delta_g^{t+1} = \beta \cdot \delta_g^t$

end

else

$\delta_l^{t+1} = \delta_l^t$, $\delta_g^{t+1} = \delta_g^t$

end

 Until converge to **fair min-max solution** h^1 ;

end

Stage 2: continue to optimize h^1 to achieve Pareto optimality

for $t = 1, 2, \dots, T^2$ **do**

Local clients:

 The center server sends the global model h^t to all local client;

 Local clients evaluate the performance of h^t and obtain the local gradients $\nabla_{\theta_i} l_i$;

Center server:

 Determine d^t by solving the LP problem in the Eq.(15) in the main text;

 Update the model $\theta^{t+1} = \theta^t + \eta d^t$;

 Until converge to h^* ;

end

Output: the **fair Pareto min-max** solution h^* .

3 B Proof of Theorem 1

4 Firstly, \mathcal{H}^{FP} and \mathcal{H}^{FU} are always non-empty from their definitions. Suppose $h \in \mathcal{H}^{FU} \subset \mathcal{H}^F$. From
5 the definitions of \mathcal{H}^{FP} and \mathcal{H}^{FU} , we have (1) $\exists h' \in \mathcal{H}^{FP}$, s.t., h' dominates h that $h' \preceq h \in \mathcal{H}^{FU}$;
6 (2) $\max(l(h)) \leq \max(l(h'))$. From (1), $\max(l(h)) \geq \max(l(h'))$ holds. Combining (2), we have
7 $\max(l(h)) = \max(l(h'))$, so $h' \in \mathcal{H}^{FU}$ holds. Therefore, $h' \in \mathcal{H}^{FU} \cap \mathcal{H}^{FP} \neq \emptyset$ holds.

8 C Convergence Derivation of FCFL

9 We will prove the convergence of our method in this section. FCFL reaches a min-max Pareto fair
10 model by a two-stage process and we aim to identify a gradient direction d to optimize the model in

11 each iteration. For N clients with N corresponding objectives $[l_1, l_2, \dots, l_N]$ and a given direction d , we
 12 have the following proposition, i.e.,

13 **Proposition 1.** *Given the gradient direction d with all $d^T \nabla_{\theta} l_i < 0$, there exists η^0 , such that:*

$$\begin{aligned} l_i(h_{\theta^{t+1}}) &< l_i(h_{\theta^t}), \forall i \in \{1, 2, \dots, N\} \\ \theta^{t+1} &= \theta^t + \eta \cdot d, \forall \eta \in [0, \eta^0]. \end{aligned} \quad (\text{S.1})$$

14

15 *Proof.* Consider Taylor's expansion with Peano's form of the differentiable objective function
 16 $l_i(h_{\theta^{t+1}})$:

$$l_i(h_{\theta^{t+1}}) = l_i(h_{\theta^t + \eta \cdot d}) = l_i(h_{\theta^t}) + \eta d^T \nabla_{\theta^t} l_i + o(\eta), \quad (\text{S.2})$$

17 where $o(\eta)$ denotes a function that approaches 0 faster than η . Specifically, $\forall \epsilon > 0$, there exists η^0 ,
 18 such that $o(\eta) < \epsilon \eta$ for all $\eta \in [0, \eta^0]$. Now we set $\epsilon = -d^T \nabla_{\theta^t} l_i > 0$, there exists η^0 , such that
 19 $l_i(h_{\theta^{t+1}}) - l_i(h_{\theta^t}) = \eta d^T \nabla_{\theta^t} l_i + o(\eta) = -\eta \epsilon + o(\eta) < 0$ for all $\eta \in [0, \eta^0]$. \square

20 We will prove the convergence of our method in two steps. First, when optimizing h to achieve
 21 min-max performance with MCF as in Eq.(11) in the main text, we optimize either \hat{l} or \hat{g} and keep \hat{l}
 22 without ascent in each iteration. The monotonic decreasing \hat{l} means h will converge to h^1 satisfying
 23 fairness and min-max constraints. Then, we continue to optimize h^1 to achieve Pareto optimality
 24 by controlling the gradient direction d without causing ascent for all objectives in Eq.(15) in the
 25 main text. The objective $\frac{1}{N} \sum_{i=1}^N l_i$ will monotonically decrease until convergence as we constrain
 26 all objectives to descend or remain unchanged defined in Eq.(15) in the main text. Assuming all
 27 objectives and their derivatives are bounded, the formal and detailed derivations are as follows.

28 **Convergence of Constrained Min-Max Optimization** To prove the convergence of the Constrained
 29 Min-Max Optimization procedure, we firstly give Lemma 1 as follows:

30 **Lemma 1.** *In each iteration, the surrogate maximum function $\hat{l}_{max}(h_{\theta^t}, \delta_l^t)$ decreases:*

$$\hat{l}_{max}(h_{\theta^{t+1}}, \delta_l^t) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t) \quad (\text{S.3a})$$

$$\hat{l}_{max}(h_{\theta^t}, \delta_l^{t+1}) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t) \quad (\text{S.3b})$$

$$\hat{l}_{max}(h_{\theta^{t+1}}, \delta_l^{t+1}) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t). \quad (\text{S.3c})$$

31

32 *Proof.* We prove Eq.(S.3a) in two cases: 1) $\hat{g}_{max}(h_{\theta^t}, \delta_g^t \leq 0)$: we obtain the gradient direction
 33 d by solving Eq.(12) in the main text. As we choose $d = -\nabla_{\theta^t} \hat{l}$, $d^T \nabla_{\theta^t} \hat{l} < 0$ holds. According
 34 to Proposition 1, as $\min d^T \nabla_{\theta^t} \hat{l} \leq d^T \nabla_{\theta^t} \hat{l} < 0$, we prove $\hat{l}_{max}(h_{\theta^{t+1}}, \delta_l^t) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t)$ as in
 35 Eq.(S.3a); 2) as $\hat{g}_{max}(h_{\theta^t}, \delta_g^t > 0)$, we obtain the gradient direction d by solving Eq.(13) in the
 36 main text. If we choose $-d$ which lies on the angular bisector of the angle formed by $\nabla_{\theta^t} \hat{l}$ and
 37 $\nabla_{\theta^t} \hat{g}$, we have $d^T \nabla_{\theta^t} \hat{l} \leq 0$ and $d^T \nabla_{\theta^t} \hat{g} \leq 0$. As the optimal solution d^* of Eq.(13) in the main text
 38 satisfies $d^{*T} \nabla_{\theta^t} \hat{l} \leq 0$ and $d^{*T} \nabla_{\theta^t} \hat{g} \leq 0$, we prove $\hat{l}_{max}(h_{\theta^{t+1}}, \delta_l^t) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t)$ as
 39 in Eq.(S.3a) in this case. While the SMF $\hat{l}_{max}(\theta^t, \delta_l^t)$ monotonically increases with respect to δ_l^t ,
 40 Eq.(S.3b) holds as $\theta^{t+1} = \beta \cdot \theta^t < \theta^t$. From Eq.(S.3a) and Eq.(S.3b), we have $\hat{l}_{max}(h_{\theta^{t+1}}, \delta_l^{t+1}) \leq$
 41 $\hat{l}_{max}(h_{\theta^t}, \delta_l^{t+1}) \leq \hat{l}_{max}(h_{\theta^t}, \delta_l^t)$ as in Eq.(S.3c) shows. \square

42 Combining the conclusion that $\hat{l}_{max}(h_{\theta^t}, \delta_l^t)$ decreases monotonically in each iteration in Lemma 1
 43 with $\hat{l}_{max}(h_{\theta^t}, \delta_l^t) \geq 0$, we prove the convergence of constrained min-max optimization described in
 44 Section 4.2.

45 **Convergence of Constrained Pareto Optimization** Constrained Pareto optimization procedure
 46 ensures the property as follows:

47 **Lemma 2.** *In each iteration, the objective $l(h_{\theta^t})$ decreases or remains unchanged:*

$$l_i(h_{\theta^{t+1}}) \leq l_i(h_{\theta^t}), \forall i \in \{1, \dots, N\}. \quad (\text{S.4})$$

48 *Proof.* We obtain the gradient direction d by solving Eq.(15) in the main text during constrained
 49 Pareto optimization. If there exists a solution d which satisfies all constraints in Eq.(15) in the main
 50 text, we have $d^T \nabla_{\theta^t} l_i \leq 0, \forall i$ which means all objectives will not increase in this iteration as Lemma
 51 2 shows. If there is no solution to Eq.(15) in the main text, we achieve Pareto optimality and the
 52 algorithm converges. \square

53 D Time Complexity Analysis

54 FCFL scales linearly with the dimension (n) of the model parameters. In our constrained min-max
 55 optimization procedure, the computation of $\hat{l}_m(h, \delta_l), \hat{g}'_m(h, \delta_g)$ in Eq.(10) has runtime of $O(N)$.
 56 With the current best LP solver [1], the LP problem with k variables and $\Omega(k)$ constraints has runtime
 57 of $O^*(k^{2.38})$ ¹. The LP problem in Eq.(12) and Eq.(12) has 2 variables and at most 4 constraints (3
 58 constraints for $d \in \bar{G}$ and 1 constraint for $d^T \nabla_{\theta^t} \hat{l} \leq 0$) so the runtime is $O^*(2^{2.38})$. The LP problem
 59 in Eq.(15) has N variables and at most $2N + 2$ constraints so the runtime is $O^*(N^{2.38})$. In deep
 60 model in FL, usually $n \gg N$.

61 E Implementation Details and Additional Experimental Results

62 E.1 Implementation Details

63 Following the experiment setting on Adult dataset in works [2, 3], we use the original Adult dataset
 64 and 66% samples are training samples and 33% are test samples. We also split the 33% as test
 65 sets and train models on the remaining 67% on eICU dataset. For the experiments in fairness
 66 constrained setting, we randomly split the eICU dataset with this ratio to run all experiments five
 67 times and report the average performance. We delete the sensitive attribute when training the model
 68 in fairness-constrained setting. All models are based on Logistic Regression and are trained and
 69 evaluated on randomly split datasets. While the original function for measuring ΔDP or ΔEO in
 70 Eq.(2) in the main text is not differentiable as there is indicator function $\hat{Y} = 1_{h(X) \geq 0.5}$, we use
 71 a surrogate differentiable function to approximate the indicator function $\hat{Y} = \frac{1}{1 + (\frac{1-h(X)}{h(X)})^{10}}$ during
 72 training. We implement our method with Pytorch and determine all hyper-parameters (including
 73 learning rate, the decay rate of δ_l, δ_g , etc.) by evaluating different combinations on the training
 74 set. We run all experiments 5 times and report the average results with stds. For all baselines
 75 except FedAve+FairReg, we use the source implementation for comparison with the optimal hyper-
 76 parameters. For the baselines on fairness-constrained experiments given uniform fairness budgets, 1)
 77 if the baselines can satisfy MCF, we try to optimize the model disparities to achieve MCF with the
 78 optimal model utilities; 2) if the baselines cannot satisfy MCF after exhaustive trying, we minimize the
 79 model disparities with reasonable model utilities. For gradient computation, we use SGD optimizer
 80 for each local client. For more algorithm details, the source code of our method is available at
 81 <https://github.com/cuis15/FCFL>.

82 E.2 Training Devices

83 We train all our models on our local Linux server with 8 GeForce RTX 2080 Ti GPUs.

84 E.3 Data Asset

85 Adult [4] is public data that anyone can download and use it freely. eICU [5] is a dataset for which
 86 permission is required. We followed the procedure on the website <https://eicu-crd.mit.edu>
 87 and got the approval for this dataset.

¹We use O^* to hide $k^{o(1)}$ and $\log^{O(1)}(1/\delta)$, δ being the relative accuracy. Detailed information is in [1]

88 **E.4 Additional Results on Fairness Constraint Experiment**

89 We show the statistic results on eICU in Table 1.

Table 1: the Statistic Results on LoS Prediction with MCF.

Methods	utility mean(%)	utility min (%)	disparity max (%)
MMCF	65.8 ± .0	61.4 ± .0	11.9 ± .0
FA	69.5 ± .05	60.1 ± .09	11.5 ± 1.2
FedAve+FairReg	67.6 ± .9	59.3 ± 1.2	11.0 ± 4.5
ours($\epsilon = 0.1$)	69.1 ± 1.3	61.4 ± .3	9.3 ± 1.4
ours($\epsilon = 0.05$)	69.1 ± 1.5	60.9 ± .5	6.8 ± 1.1

90 Table 1 shows the statistical results of LoS prediction with equal fairness budgets. Our method
 91 achieves the min-max performance as we set the fairness budget $\epsilon = 0.1$. All baselines cannot reduce
 92 the disparities below 0.1. When we constrain all disparities $\Delta DP_i \leq 0.05$, our method reduces the
 93 disparities significantly compared to baselines with a comparable min-max accuracy 60.9%.

94 **Income Prediction with sensitive attribute being *gender* and equal fairness budgets**

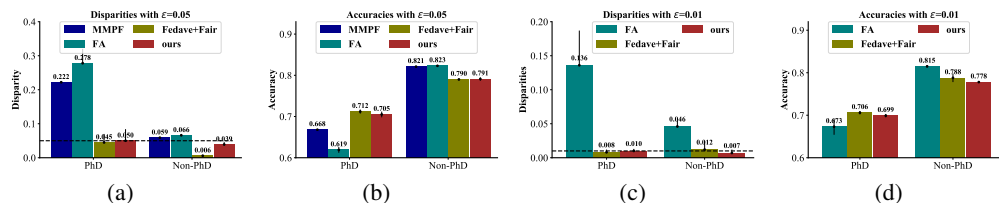
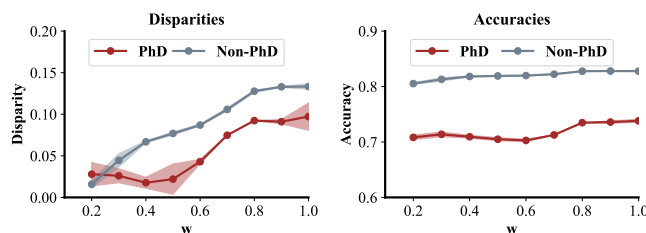


Figure 1: The disparities and accuracies on both clients as $\epsilon = 0.05$ and as $\epsilon = 0.01$ of on Adult dataset when *gender* is the sensitive attribute.

95 In this experiment, we select *gender* as sensitive attribute. From Figure 1, with the budget $\epsilon_i =$
 96 $0.05, \forall i \in \{1, \dots, N\}$, we achieve comparable min-max performance with MCF while MMPF and FA
 97 violate the constraint on PhD client. As the fairness budget $\epsilon_i = 0.01, \forall i \in \{1, \dots, N\}$, all baselines
 98 violate the constraint on PhD client and we maintain the utilities on both clients with MCF.

99 **Income Prediction with sensitive attribute being *gender* and client-specific fairness budgets**

100 The results of Income prediction with client-specific fairness constraints in Figure 2. The disparities
 101 of both clients decrease significantly as in Figure 2(a) as w decreases. With a decreasing fairness
 budget, the utilities of both clients slightly decreases as in Figure 2(b).



(a) The disparities vary as w decreases (b) The accuracies vary as w decreases

Figure 2: Client-specific constraint experiment on Adult dataset with sensitive attribute being *gender*.

102

103 **Income prediction using EO metric** Besides DP [6] metric, we verify the effectiveness of our
 104 method using another metric *Equal opportunity* (EO) [7] which measures difference of the false
 105 negative rates across different groups as in Eq.(2). Here we show the results on Income prediction
 106 with the sensitive attribute being *race* given uniform fairness budgets $\epsilon_i = 0.05$ on both clients in our
 107 main text.

108 The original unconstrained model causes disparities on both clients: $\Delta EO_{PhD} = 0.105$ and
 109 $\Delta EO_{non-PhD} = 0.183$. Our model significantly reduces the disparities with fairness budget

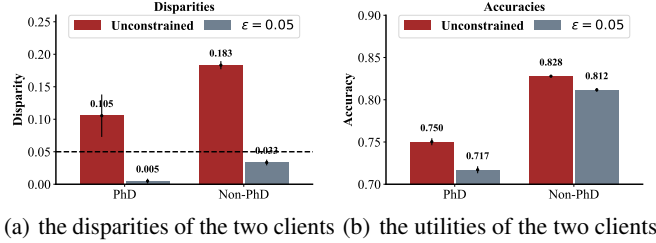


Figure 3: The results of unconstrained optimization and fairness-constrained optimization with the disparities measured by EO and the sensitive attribute is *race*.

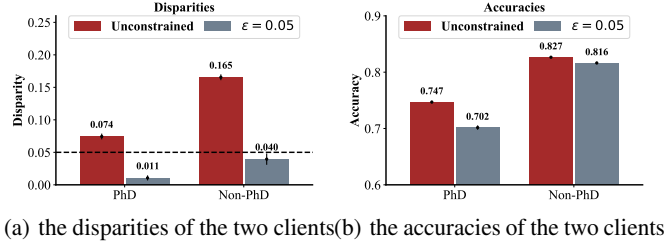


Figure 4: The results of unconstrained optimization and fairness-constrained optimization with the disparities measured by EO and the sensitive attribute is *gender*.

110 $\epsilon = 0.05$ as shown in Figure 3(a). For the model utilities on both clients shown in Figure 3(b), we
 111 achieve reasonable performances and there is a 0.03 drop on PhD client and 0.02 on non-PhD client.
 112 Similar results can also be found in Figure 4 as the sensitive attribute is *gender* in Appendix.

113 F Compatibility with Unconstrained Pareto Min-Max Optimization

114 F.1 Achieving Pareto Min-Max Optimality without Fairness Constraints

115 Though our method is proposed for fairness-constrained multi-objective optimization, FCFL is also
 116 compatible with unconstrained min-max optimization problems which only care about the utilities of
 117 all clients as in [2, 3]:

1. optimize h to achieve min-max performance:

$$\min_{d \in \bar{G}} d^T \nabla_{\theta^t} \hat{l},$$

2. optimize h for Pareto optimality and min-max performance:

$$\min_{d \in \bar{G}} \frac{1}{N} \sum_{i=1}^N d^T \nabla_{\theta^t} l_i, \quad \text{s.t. } d^T \nabla_{\theta^t} l_i \leq 0 \quad \forall i \in \{1, \dots, N\},$$

118 where \bar{G} is the convex hull of $[\nabla_{\theta^t} l_1, \dots, \nabla_{\theta^t} l_N]$ and we obtain the gradient direction d using the
 119 gradient information of all clients without accessing the local data.

120 F.2 Experiments on Improving Consistency without Fairness Constraints

121 We evaluate the performance of our method on the problem of improving consistency without fairness
 122 constraints described with two existing methods q-FFL [2] and AFL [3].

123 Our method seeks for min-max performance by optimizing the SMF \hat{l} which is the upper bound of
 124 all objectives. Experimental results on Income Prediction in Table 2 show that our method achieves
 125 relatively higher performance on the worst-performing client (74.9). Besides, FCFL achieves Pareto

Table 2: The accuracies on Income Prediction

Methods	average (%)	PhD (%)	non-PhD (%)
q-FFL	82.3 \pm .1[2]	74.4 \pm .9[2]	82.4 \pm .1[2]
AFL	82.5 \pm .5[2]	73.0 \pm 2.2[2]	82.6 \pm .5[2]
ours	82.7 \pm .1	74.9 \pm .4	82.8 \pm .1

126 optimality to avoid unnecessary harm to other clients. From Table 2, we maintain the performance on
 127 non-PhD client (82.8).

128 We conduct experiments on eICU dataset in unconstrained min-max setting. We randomly split the
 129 dataset 5 times and run the two prediction tasks. We show the statistical average results in Table 3.
 130 For LoS prediction task, we achieve higher uniformity by a higher utility on the worst-performing
 131 client compared to baselines. We also predict the in-hospital mortality as a prediction task and all
 132 methods achieve similar results.

Table 3: The accuracies on eICU without fairness constraints

Methods	Mortality Prediction		LoS Prediction	
	minimum (%)	average(%)	minimum (%)	average(%)
q-FFL	91.7 \pm .1	88.3 \pm .7	57.6 \pm 2.2	70.0 \pm .3
AFL	91.7 \pm .1	88.2 \pm .7	58.1 \pm 2.0	70.0 \pm .4
ous	91.7 \pm .1	88.3 \pm .6	60.5 \pm 2.0	67.4 \pm .7

133 G Discussion about Fairness Budget

134 G.1 How is the fairness budget of each client relate to one another

135 For each fairness budget ϵ_k , it should be assigned by each client based upon its actual fairness
 136 requirements. However, the fairness budgets ϵ_k assigned by different clients are not unrelated since
 137 all fairness constraints defined by ϵ_k determine the feasible region of the model together. Due to the
 138 potential trade-off between the fairness and the utility, the fairness constraint determined by ϵ_k of the
 139 k-th client can have an impact on the utility of all clients in the federated network.

140 G.2 When the fairness budget should be the same

141 As we stated above, the assignment of the fairness budget ϵ_k depends on the actual scenarios. In
 142 high-stake scenarios (e.g., hospitals, banks, etc.), people are highly concerned about fairness and
 143 different clients should have consistent fairness budgets. In some other scenarios such as advertising
 144 recommendations, people may be more tolerant of the difference of ϵ_k among different clients. In this
 145 case, the assignment of the fairness budget ϵ_k depends on whether it can lead to a satisfying trade-off
 146 between fairness and utility for all clients. For example, we determine the ϵ_k based upon the original
 147 disparity as the experiments presented in Sec 5.4. In addition, to obtain a budget combination that
 148 satisfies all clients, one may try different budget combinations and evaluate the performances of all
 149 clients, then all clients vote for a reasonable budget combination.

150 H Limitations and Future work

151 In this study, we aim to tackle the algorithmic disparity and performance inconsistency issues in
 152 federated learning. As it is hard to realize the min-max optimality and Pareto optimality by one-
 153 stage optimization, we propose a two-stage optimization framework that first achieving consistent
 154 performance then enforcing Pareto optimality. Since our framework encourages a more uniform
 155 model performance distribution, on the one hand, some clients with poor performances may be
 156 significantly improved; on the other hand, other clients whose model utility could be further improved
 157 may come to a halt.

158 Since we focus on addressing the local disparity, it cannot guarantee reasonable global fairness. In
159 reality, the global disparity is also a matter of concern in a federated network. We will study how to
160 give consideration to both local fairness and global fairness in our future work.

161 **References**

- 162 [1] Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix multiplication
163 time. In *Proceedings of the 51st annual ACM SIGACT symposium on theory of computing*, pages 938–942,
164 2019.
- 165 [2] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning.
166 In *International Conference on Learning Representations*, 2019.
- 167 [3] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International*
168 *Conference on Machine Learning*, pages 4615–4625, 2019.
- 169 [4] Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *Kdd*, volume 96,
170 pages 202–207, 1996.
- 171 [5] Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. The eicu
172 collaborative research database, a freely available multi-center database for critical care research. *Scientific*
173 *data*, 5(1):1–13, 2018.
- 174 [6] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through
175 awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226,
176 2012.
- 177 [7] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *Proceedings*
178 *of the 30th International Conference on Neural Information Processing Systems*, pages 3323–3331, 2016.