

---

# Fairness in Federated Learning via Core-Stability

---

Bhaskar Ray Chaudhury Yinyi Li Mintong Kang Bo Li Ruta Mehta

University of Illinois at Urbana Champaign

## Abstract

1 Federated learning provides an effective paradigm to jointly optimize a model  
2 benefited from rich distributed data while protecting data privacy. Nonetheless, the  
3 heterogeneity nature of distributed data, especially in the non-IID setting, makes it  
4 challenging to define and ensure fairness among local agents. For instance, it is  
5 intuitively “unfair” for agents with data of high quality to sacrifice their performance  
6 due to other agents with low quality data. Currently popular *egalitarian* and  
7 *weighted equity-based* fairness measures suffer from the aforementioned pitfall. In  
8 this work, we aim to formally represent this problem and address these fairness  
9 issues using concepts from co-operative game theory and social choice theory. We  
10 model the task of learning a shared predictor in the federated setting as a *fair public*  
11 *decision making* problem, and then define the notion of *core-stable fairness*: Given  
12  $N$  agents, there is no subset of agents  $S$  that can benefit significantly by forming a  
13 coalition among themselves based on their utilities  $U_N$  and  $U_S$  (i.e.,  $\frac{|S|}{N}U_S \geq U_N$ ).  
14 Core-stable predictors are robust to low quality local data from some agents, and  
15 additionally they satisfy *Proportionality* (each agent gets at least  $1/n$  fraction of  
16 the best utility that she can get from any predictor) and Pareto-optimality (there  
17 exists no model that can increase the utility of an agent without decreasing the  
18 utility of another), two well sought-after fairness and efficiency notions within  
19 social choice. We then propose an efficient federated learning protocol CoreFed to  
20 optimize a core stable predictor. CoreFed determines a core-stable predictor when  
21 the loss functions of the agents are convex. CoreFed also determines approximate  
22 core-stable predictors when the loss functions are not convex, like smooth neural  
23 networks. We further show the existence of core-stable predictors in more general  
24 settings using Kakutani’s fixed point theorem. Finally, we empirically validate our  
25 analysis on two real-world datasets, and we show that CoreFed achieves higher  
26 core-stable fairness than FedAvg while maintaining similar accuracy.

## 27 1 Introduction

28 The success of many deployed machine learning (ML) systems crucially hinges on the availability of  
29 high-quality data. However, a single entity might not own all the data it needs to train the ML model  
30 it wants; instead, valuable data instances or features might be scattered in different organizations or  
31 entities. Distributed learning schemes such as federated learning (FL) [16] provide a training scheme  
32 that focuses on training a single ML model using all the data available in a cooperative way without  
33 moving the training data across the organizational or personal boundaries to protect data privacy.

34 On the other hand, given the heterogeneity nature of local data in FL, ensuring fairness among agents  
35 has attracted intensive interests. A significant part of the existing literature mainly focus on ensuring  
36 equal accuracy for different agents [32, 23, 20] or the fairness for the final aggregated model regarding  
37 the protected attributes without taking the imbalanced contribution of agents into account [14, 35].  
38 In this work, we ask: *Is it possible to jointly optimize a centralized model with fairness guarantees*  
39 *regarding the heterogeneity of local agents? How to define such fairness such that no agents would*  
40 *intend to form an alternative coalition with a subset of agents? What could be the federated learning*  
41 *protocol that is able to ensure such fairness?*

42 To address the above research questions, we bring to bear notions from game theory and social choice  
 43 theory. We first observe that federated learning can be cast into *public decision making*, where all  
 44 agents derive their respective utilities from a common global decision, namely the globally learned  
 45 model. Now the goal is to make this global decision fairly. One of the fundamental fairness measures  
 46 in public decision making is that of core-stability [25]. Intuitively, we say that a set of agents can  
 47 form a “blocking coalition” if each one of them can gain utility significantly (proportional to the size  
 48 of their coalition) by training a unified model amongst themselves than the globally trained model. A  
 49 globally trained model is core stable if there are no blocking coalitions.

50 We briefly justify the advantages of core-stability (a.k.a. core-stable fairness) over some of the  
 51 existing notions of fairness in federated learning. Two commonly used fairness notions in federated  
 52 learning are the *egalitarian fairness* [32, 23, 8, 35, 36, 26, 23] and *proportional fairness* [7, 6]. The  
 53 egalitarian fairness aims to maximize the utility of the least happy agent<sup>1</sup>. In a proportional fairness,  
 54 we want the ratios of the losses of any pair of the agents to be (super/ sub) proportional to the size of  
 55 their respective datasets (this incentivizes the agents to share more of their data with the server). To  
 56 avoid naming conflicts with our notion of *proportionality*, from here on, we refer to the proportional  
 57 fairness introduced in [7] as *weighted equity based fairness* as this fairness compares the losses  
 58 of every pair of agents. Usually, fairness notions that compare the utilities/ losses of agents with  
 59 each other are called equity based fairness in social choice theory. We remark that both the notions  
 60 are vulnerable if some agents have poor quality datasets. In particular, if one of the agents have  
 61 high levels of noise in their data, call them noisy-agent, then their loss values will tend to be higher  
 62 for most learnt predictors. The egalitarian fairness and the weighted equity based fairness may be  
 63 unfair to the other agents as both may make decisions aiming to reduce the large loss incurred by  
 64 the noisy-agent, thereby biasing the learning towards the data of the noisy-agent. A more desirable  
 65 fairness property in this scenario maybe to compare the *loss percentage* of agents, i.e., the ratio of  
 66 the loss incurred to the maximum loss that can be incurred, or equivalently *utility percentage* of  
 67 agents, i.e., the ratio of the utility incurred to the maximum utility that can be incurred. Core-stability  
 68 achieves this, together with other desirable properties (elaborated in Section 3).

69 **Our contribution.** We formally define the *core-stable fairness* in federated learning by appropri-  
 70 ately modeling agent’s utility functions to capture their learning loss error. In particular, given a group  
 71 of  $N$  local agents, an aggregation protocol  $P$ , and an aggregated model  $f$ , we say that the model  $f$   
 72 achieves core-stability if there are no coalition  $S$  of agents that could benefit significantly by training  
 73 a model with only their data (see Definition 1). Intuitively, this means that under a core-stable FL  
 74 model, no agent has the incentive to deviate from current group and thereby obtain proportionally  
 75 better aggregate utility from the final trained model. Additionally, we note that such a model  $f$   
 76 will ensure sought-after guarantees of Proportionality (each agent gets  $1/n$  times their best possible  
 77 utility)[31] and Pareto-optimality (there is no predictor that can increase the utility of any agent  
 78 without decreasing the utility of another agent) that equal-treatment based models [7, 32] may fail to.

79 Core-stability is a well-sought-after but a rare-to-exist notion. In case of *public goods* that resembles  
 80 FL, existence of core-stable outcome was known only when agent’s utility functions are linear [9].  
 81 While the utility functions that capture learning errors are inherently non-linear and highly complex  
 82 making existing results inapplicable. We summarize our *contributions* as below.

- 83 • We formally formulate core-stability from co-operative game theory to fairness in federated learning.  
 84 We show that core-stability exists (in Section 4.1) as long as agent’s utility functions are continuous  
 85 with respect to the model parameters, and their non-negative conical combinations have a convex  
 86 set of (local) optima. We prove this result using a fixed point formulation. In particular, we define a  
 87 correspondence  $\phi : P \rightarrow P$  on the set of all feasible predictors  $P$ , and ensure that any predictor  
 88  $\theta^* \in P$  such that  $\theta^* \in \phi(\theta^*)$  is core-stable. Thereafter we show that  $\phi$  satisfies nice continuity like  
 89 properties and therefore must admit a fixed point by Kakutani’s fixed-point theorem [15].
- 90 • Next, we design an effective federated learning protocol CoreFed, which optimizes the final model  
 91 by maximizing the protocol of agent’s utilities. We prove that this protocol efficiently finds the  
 92 core-stable model whenever the underlying utility functions are concave (see Section 4.2). Our  
 93 protocol only needs gradient information from agents in each round.
- 94 • We prove that above method directly applies to learning through linear regression or logistic  
 95 regression, since their resulting utility functions are convex (see Section 4.3). For *Smooth Neural*

---

<sup>1</sup>Equivalently maximize the minimum loss.

96 *Nets (DNN)*, although the utility functions are (highly) non-convex, we manage to show that an  
 97 approximate core-stable model can be learned within a local neighborhood (see Section 4.4).  
 98 • To capture cases where agents may have varying importance, we extend core-stability to *weighted*  
 99 core-stability (in Section 4.5). We show that a *weighted* core-stable model is *weighted* proportional  
 100 and Pareto-optimal, and that CoreFed protocol can be generalized to Weighted CoreFed to get the  
 101 desired weighted guarantees.  
 102 • We conduct experiments on three datasets, and show that CoreFed achieves the core-stable fairness,  
 103 while maintaining similar utility with the standard FedAvg protocol (see Section 5).

## 104 2 Related Work

105 **Fairness in social choice.** Fairness is one of the fundamental goals in many multi-agent settings.  
 106 Over the years, motivated by applications, several notions of fairness have been proposed and  
 107 investigated. Two fairness notions that are studied in many applications are that of *proportionality* [31]  
 108 and *envy-freeness* [10]. Proportionality requires every agent to receive their proportional share of the  
 109 best outcome, i.e., at least  $1/n$  times their best possible utility. The notion of envy-freeness is defined  
 110 in the context of resource allocation, where one aims to divide a set of items among agents fairly.  
 111 In an envy-free allocation, no agent prefers the bundle of the other agent to her own. However, in  
 112 *public decision making*, where all agents derive utility from a common global decision, this notion  
 113 is not applicable! In public decision making, one of the most sought out fairness notion is that of  
 114 *core-stability* [25]. Core-stability generalizes the notion of proportionality alongside other desirable  
 115 properties like Pareto-optimality. The concepts of Core-stability find applications in many other  
 116 settings in social choice and game theory and is well known to exist in some special settings [1]  
 117 Another popular fairness notion is *equitability* which states that every agent should be equally happy,  
 118 i.e., the utilities/ losses of all the agents should be the same. However, as explained in 1, using  
 119 relaxations of equitability, may lead to undesirable outcomes if some agents have poor data quality.  
 120 Over the last seven decades, several relaxations of envy-freeness [21, 27], proportionality [5, 24] and  
 121 equitability [11] have been studied in computational social choice.

122 **Fairness in federated learning.** There have been several results on fairness in federated learning,  
 123 each focusing on a particular aspect of the entire paradigm. For instance some work [13, 33] aim to  
 124 establish fairness at the agent selection phase, where the server requests for updates from a selected  
 125 subset of the agents. There are studies that aim to study fairness while training the global model  
 126 such that it does not discriminate against protected groups [34] or that the model does not overfit the  
 127 data of some agents at the expense of others [20, 23]. Earlier mentioned *egalitarian fairness* will fall  
 128 under this category. Then, there are studies that consider fairness by evaluating the contribution of  
 129 the agents towards training the joint model– for instance *weighted equity fairness* [7] does this based  
 130 on the size of the data shared by the agents. Other studies assign significance to the agents based  
 131 on Shapley values [30]. For a full taxonomy of fairness in federated learning, we urge the reader to  
 132 check [29]. At large, most of the fairness notions are incomparable. As remarked in [7], “no one set  
 133 of definitions is going to resolve the complex questions it raises”.

## 134 3 Core-Stability in Federated Learning

135 **Problem Setup.** In any *predictive modelling* task, one would like to learn a function mapping from  
 136  $\mathcal{X} \subseteq \mathbb{R}^d$  to  $\mathcal{Y} \subseteq \mathbb{R}$ . This includes both *regression* and binary *classification* where extension to multi-  
 137 class classification is also feasible. We denote the space of such mappings as  $\mathcal{F} = \{f_\theta \mid \theta \in P \subseteq \mathbb{R}^d\}$ ,  
 138 where each  $f_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$  is a mapping function parameterized by the model weights vector  $\theta$ . Our goal  
 139 is to determine  $f_\theta \in \mathcal{F}$ , such that for data sample  $(x, y)$  drawn from the distribution  $\mathcal{P}$ ,  $f_\theta(x) \approx y$ ,  
 140 i.e.,  $f_\theta(x)$  learns  $y$  well. Since we identify a mapping function with each  $\theta \in P$ , we refer to  $\theta$  as a  
 141 *predictor* for the model <sup>2</sup>.

142 **Utility Functions of the Agents.** The quality of a predictor  $\theta$  is usually measured by the  
 143 expected loss over the data distribution  $\mathcal{P}$ , i.e.,  $\mathbb{E}_{(x,y) \sim \mathcal{P}} \ell(f_\theta(x), y)$ , where  $\ell(\cdot, \cdot)$  is a pre-  
 144 defined loss function. Ideally, the training process would minimize this expected loss, i.e., attain  
 145  $\theta^* = \arg \min_{\theta \in P} \mathbb{E}_{(x,y) \sim \mathcal{P}} \ell(f_\theta(x), y)$ . Since we are trying to determine a single predictor for

---

<sup>2</sup>each  $\theta$  is a predictor

146 several heterogeneous agents/ groups, we may not be able to give every group its best predictor.  
 147 However, we want to choose the predictor that achieves *fairness* across all the groups. To define any  
 148 notion of fairness from the classical economics literature, we need to define the **utility function** of a  
 149 group for a predictor  $\theta$ . Intuitively, the utility is a measure of how good the predictor is for the group  
 150 and its data. We define,

$$u(\theta) = M - \mathbb{E}_{(x,y) \sim \mathcal{P}} \ell(f_\theta(x), y) \quad (1)$$

151 where  $M$  is a constant more than  $(1 + \varepsilon)$  times the loss incurred from the worst predictor for agent  
 152  $i$ , i.e.,  $M \geq (1 + \varepsilon) \sup_{\theta \in P, (x,y) \in \mathcal{X} \times \mathcal{Y}} \ell(f_\theta(x), y)$ . The scaling by  $(1 + \varepsilon)$  is to avoid unnecessary  
 153 degeneracies resulting from zero utilities, and we choose  $\varepsilon \ll 10^{-5}$ . Observe that the range of the  
 154 utility function is from  $0 < M\varepsilon$  to  $M(1 + \varepsilon)$ .

155 **Federated Learning and Fairness.** In the federated learning setting, we are given  $n$  groups. Each  
 156 group  $i$  has their loss function  $\ell_i(\cdot)$  and correspondingly a utility function  $u_i(\cdot)$  for each choice of  
 157 a predictor  $\theta \in P$ . We now define the fairness criterion. Recall that given  $n$  groups, our goal is  
 158 to choose a  $\theta \in P$ , such that we are fair to all the involved agents. The fairness notion here is  
 159 *core-stability*.

160 **Definition 1 (Core-Stability).** A predictor  $\theta \in P$ , is called core stable if there exists no other  $\theta' \in P$ ,  
 161 and no subset  $S$  of agents, such that  $\frac{|S|}{n} u_i(\theta') \geq u_i(\theta)$  for all  $i \in S$ , with at least one strict inequality.

162 Intuitively, core-stability implies that there is no subset of agents that can benefit “significantly” by  
 163 forming a coalition among themselves, i.e., if we were to choose any other  $\theta' \in P$  only considering  
 164 the utility functions of the agents in  $S \subseteq n$ , then there is some agent who’s (multiplicative) gain in  
 165 utility will be strictly less than a factor  $n/|S|$ , i.e., there is no substantial benefit for this agent if she  
 166 chooses to belong to the set  $S$ . Furthermore, *core-stability* gives some classical fairness and welfare  
 167 guarantees. In particular, note that every agent  $i$  gets at least  $1/n$  times her best utility, i.e., the utility  
 168 derived from the best possible mapping for agent  $i$ . Mathematically  $u_i(\theta) \geq 1/n \cdot u_i(\theta')$  for all  
 169  $\theta' \in P$  (setting  $S = \{i\}$  in Definition 1). This fairness is called *Proportionality* [31]. Formally,

170 **Definition 2 (Proportionality).** A predictor  $\theta \in P$  is proportional if and only if for all  $\theta' \in P$ , we  
 171 have  $u_i(\theta) \geq \frac{u_i(\theta')}{n}$  for all  $i \in [n]$ .

172 Similarly, observe that there exists no predictor  $\theta' \in P$  where  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} > n$  (setting  $S = [n]$   
 173 in Definition 1). This implies that there is no predictor that can increase the utility of some agent  
 174 without decreasing the utility of another. We call this property *Pareto-optimality*. Formally,

175 **Definition 3 (Pareto-optimality).** A predictor  $\theta \in P$  is Pareto-optimal if and only if there exists no  
 176 other  $\theta' \in P$ , such that  $u_i(\theta') \geq u_i(\theta)$  for all  $i \in [n]$  with at least one strict inequality.

177 Core is a central concept within cooperative game theory, defined to capture “no deviating sub-group”  
 178 property and is considered very strong. However, it is well known to exist only in special cases [1].  
 179 We now elaborate the advantages of core-stability over some of the existing fairness concepts in  
 180 federated learning.

### 181 3.1 Advantages of Core-Stability

182 As briefly mentioned in the introduction, core-stability is robust to low local data qualities of some  
 183 agents, unlike the FedAvg or federated learning based on egalitarian or weighted equity based fairness.  
 184 We elaborate this with a small example: consider three agents that contribute equal amount of data,  
 185 and say agent 3 has poor data quality, i.e., there exists no proper predictor for agent 3’s data, or  
 186 equivalently for all predictors  $\theta \in P$ , the loss function of this agent is very high (and utility is very  
 187 low). Concretely, consider two predictors  $\theta_1$  and  $\theta_2$ . Under  $\theta_1$ , agents 1 and 2 incur a loss of  $a$   
 188 and agent 3 incurs a loss of  $M \cdot a$  (think of  $M$  as a very large integer). Now, under  $\theta_2$ , agents 1  
 189 and 2 have a loss of  $10a$  and agent 3 has a loss of  $0.9Ma$ . Observe that  $\theta_2$  is preferable over  $\theta_1$   
 190 under egalitarian fairness (as  $0.9Ma \gg 10a$ ) and also by FedAvg as it has a lower total average loss  
 191 ( $0.1Ma \gg 9a$ ). Similarly, a learning algorithm based on weighted equity fairness would prefer  $\theta_2$ ,  
 192 as ratio of the losses between agents 1 (or 2) and 3 is significantly high in both  $\theta_1$  and  $\theta_2$  and is lesser  
 193 in  $\theta_2$ . However, intuitively,  $\theta_1$  seems fairer, as agent 3 is not substantially worse off in  $\theta_1$  than it is in  
 194  $\theta_2$  (by a factor 1.1), while agents 1 and 2 are significantly better off in  $\theta_1$  (by a factor 10). Note that  
 195 in this example  $\theta_2$  is not a core-stable predictor, as agents 1 and 2 have an incentive to break off and

196 improve substantially (intuitively  $\theta_2$  is very unfair to them). We say that core-stable predictors are  
 197 robust to low data quality of specific agents, as we never compare the losses of two agents with each  
 198 other; rather our comparison is more along the lines of *loss percentages*, i.e., the ratio of the loss to  
 199 the maximum possible loss incurred by the agent.

200 The robustness to poor local data quality of agents of core-stable predictors is a parallel to the  
 201 property of *scale-invariance* that core-stable allocations exhibit in social choice theory. In particular,  
 202 scaling the utility of any single agent does not alter the core-stable allocation. Similarly, Egalitarian,  
 203 utilitarian<sup>3</sup> and equity based fairness suffer from being responsive to scaling [4].

## 204 4 Core-Stability in Federated Learning

205 In this section, we prove the existence of core-stability under certain assumptions on the loss functions  
 206 of the individual agents/ groups (Section 4.1). Then, in Section 4.2, we give a distributive training  
 207 protocol CoreFed to determine a core-stable predictor when the loss functions are convex<sup>4</sup>. Finally,  
 208 by applying the theory developed in Sections 4.1 and 4.2, we show that CoreFed determines a core  
 209 stable predictor in Linear Regression, and in Classification with Logistic Regression (Section 4.3).  
 210 Finally, in Section 4.4 we show that CoreFed determines an approximate core stable predictor in  
 211 Deep Neural Networks.

### 212 4.1 Existence of Core-Stability in Federated Learning.

213 We show that core stable predictors exist in the federated setting if the utility functions of the agents  
 214 satisfy the following conditions:

- 215 1. The utility function of each agent is continuous.
- 216 2. The set of maximizers of any conical combination of the utility functions is convex i.e., for  
 217 all  $\langle \alpha_1, \alpha_2, \dots, \alpha_n \rangle \in \mathbb{R}_{\geq 0}^n$ , the set  $C = \{ \theta \mid \sum_{i \in [n]} \alpha_i u_i(\theta) \text{ is maximum} \}$  is convex.

218 To the best of our knowledge, prior to this work, the existence of core-stability in public fair division  
 219 was shown only for linear utility functions by [9]. Utility functions that satisfy the above two  
 220 conditions cover several other utility functions and is therefore a strict generalization of linear utility  
 221 functions. We show the existence of core-stability for instances satisfying conditions 1 and 2 above  
 222 using Kakutani’s fixed point theorem. For completeness, we state the Kakutani’s fixed point theorem.

223 **Definition 4.** [Kakutani’s Fixed Point Theorem] A *correspondence* or equivalently a *set valued*  
 224 *function*  $\phi: D \rightarrow 2^D$  admits a fixed point, i.e., there exists a point  $d \in D$ , such that  $d \in \phi(d)$ , if

- 225 1.  $D$  is non-empty, compact, and convex.
- 226 2. For all  $d \in D$ ,  $\phi(d)$  is non-empty, convex and compact.
- 227 3.  $\phi(\cdot)$  has a closed graph, i.e., for all sequences  $(d_i)_{i \in \mathbb{N}}$  converging to  $d^*$  and  $(e_i)_{i \in \mathbb{N}}$  converg-  
 228 ing to  $e^*$ , such that  $d_i \in D$  and  $e_i \in \phi(d_i)$ , we have  $e^* \in \phi(d^*)$ .

229 We define a correspondence or a set valued function  $\phi: P \rightarrow 2^P$  where  $P$  is the set of all feasible  
 230 predictors. In particular, for all  $\theta \in P$ , we set

$$\phi(\theta) = \left\{ d \mid \sum_{i \in [n]} \frac{u_i(d)}{u_i(\theta)} \text{ is maximum} \right\}$$

231 We first observe that any fixed point of  $\phi$  corresponds to a core stable classifier.

232 **Lemma 4.1.** *Let  $\theta \in P$  be such that  $\theta \in \phi(\theta)$ . Then,  $\theta$  is a core-stable predictor.*

233 The proof of Lemma 4.1 can be found in the Appendix. Now, it suffices to show that  $\phi$  admits a fixed  
 234 point. In particular, note that the domain of  $\phi$ ,  $P$  is non-empty, compact, and convex. Similarly, for  
 235 every  $\theta \in P$ ,  $\phi(\theta)$  is non-empty, compact, and convex. By Kakutani’s fixed point theorem, it only  
 236 remains to show that  $\phi(\cdot)$  has a closed graph, to ensure that  $\phi(\cdot)$  admits a fixed point.

<sup>3</sup>This is the parallel to FedAvg in social choice theory.

<sup>4</sup>The assumptions made in Section 4.1 to show only existence of core-stability are weaker than the convexity assumptions in Section 4.2.

237 **Lemma 4.2.** *The correspondence  $\phi(\cdot)$  has a closed graph.*

238 The detailed proof can be found in the Appendix. Intuitively, since the utility functions are continuous  
 239 and non-zero, the optima of  $\sum_{i \in [n]} \frac{u_i(d)}{u_i(\theta)}$  over  $d \in P$ , also changes continuously with  $\theta$ . We are  
 240 ready to prove the main result of this section.

241 **Theorem 1.** *In any federated learning setting, where the agent’s utilities are continuous and the set  
 242 of maximizers of any conical combination of the agents utilities is convex, a core-stable predictor  
 243 exists.*

244 *Proof Sketch.* Any fixed point of  $\phi(\cdot)$  corresponds to core stable predictor (Lemma 4.1). It suffices to  
 245 show that  $\phi(\cdot)$  admits a fixed point under assumptions in Theorem 1. To this end, note that domain  $P$   
 246 of  $\phi(\cdot)$  is non-empty, compact, and convex. For all  $\theta \in P$ ,  $\phi(\theta)$  is convex by assumption in Theorem 1.  
 247 Finally  $\phi(\cdot)$  has a closed graph by Lemma 4.2, and thus  $\phi(\cdot)$  admits a fixed point.

248 **Implications.** Theorem 1 describes the conditions under which core-stable predictors exist. We  
 249 briefly state how to adapt the proofs to show the existence of *locally core-stable* predictors for more  
 250 general utility functions. Lemmas 4.1 4.2, and Theorem 1 are valid even if we change the definition  
 251 of  $\phi(\theta)$  to the set of maximizers of  $\sum_{i \in [n]} u_i(d)/u_i(\theta)$  over  $d \in \mathcal{B}(\theta, r)$  (instead of  $d \in P$ ), i.e., we  
 252 define  $\phi(\theta)$  to be the set of maximizers in the local neighbourhood of  $\theta$  (within distance  $r$  to  $\theta$ ). In  
 253 such a case, we only need conditions 1 and 2 to be true within a radius of  $r$  of every point, i.e., within  
 254  $\mathcal{B}(\theta, r)$  for all  $\theta \in P$ . These guarantees typically tend to be true for small values of  $r$  in Deep Neural  
 255 Networks. Thus, the predictor corresponding to the fixed point of  $\phi$  will satisfy core-stability when  
 256 restricted to predictors within distance  $r$  to it, i.e., it is a locally core-stable predictor.

## 257 4.2 Computation of a Core-Stable Predictor When Utility Functions are Concave.

258 In this section, we show that under certain assumptions on the utility functions, we can describe an  
 259 efficient distributed protocol that computes the core-stable predictor. In particular, we look into the  
 260 scenario, where the utility function of each agent is concave. Note that this would automatically  
 261 satisfy the conditions in Theorem 1, as any conical combination of concave functions is also concave  
 262 and will admit a convex set of maximizers.

263 We first show that a core stable predictor can be expressed as an optima of a convex program. In partic-  
 264 ular, any predictor that maximizes the product of utilities of the agents, i.e.,  $\operatorname{argmax}_{\theta \in P} \prod_{i \in [n]} u_i(\theta)$   
 265 (or equivalently the sum of logarithms of the utilities of the agents), is core stable.

$$\begin{aligned} & \text{maximize} && \mathcal{L}(\theta) = \sum_{i \in [n]} \log(u_i(\theta)) \\ & \text{subject to} && \theta \in P \end{aligned} \tag{2}$$

266 Observe that if the utility of each agent is concave, then the above program is convex. Since the  
 267 logarithm is a concave increasing function, each  $\log(u_i(\cdot))$  is a concave in  $\theta$  and the sum of concave  
 268 functions is concave. Thus, 2 is a concave maximization subject to convex constraints.

269 **Theorem 2.** *If  $u_i(\cdot)$  is concave for all  $i \in [n]$ , then any predictor  $\theta^*$  that maximizes the convex  
 270 program 2 is core-stable.*

271 *Proof Sketch.* We defer a formal proof to Appendix C. The main technical ingredient of our proof  
 272 is to show that if  $\theta^*$  is a solution to the convex program 2, then, for any other predictor  $\theta' \in P$ ,  
 273 we have  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} \leq n$ . Now if  $\theta^*$  is not core-stable, then there exists an  $S \subseteq [n]$  and  $\theta' \in P$ ,  
 274 such that  $u_i(\theta') \geq n/|S|u_i(\theta^*)$  for all  $i \in S$  with at least one strict inequality, then we have  
 275  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} \geq \sum_{i \in S} \frac{u_i(\theta')}{u_i(\theta^*)} > n/|S| \cdot |S| = n$ , which is a contradiction.

276 **Implications.** The proof of Theorem 2 shows that the strong utilitarian property of  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} \leq n$   
 277 for any  $\theta' \in P$  implies core-stability of  $\theta^*$  under *any* (non-negative) utility functions. Clearly, such a  
 278  $\theta^*$  must be Pareto-optimal, and furthermore, the inequality implies that at  $\theta'$  computed by any other  
 279 classical method, if some agents gain, then some other agents must be loosing by a lot. Secondly,  
 280 under concave utilities optima of convex program equation 2 satisfies this property, and hence can be  
 281 computed in efficiently. Below we discuss a distributed protocol for the same.

282 We propose a distributed SGD framework to determine a core stable predictor. We call our Algorithm  
 283 as CoreFed (Fully outlined in Algorithm 1 in the appendix).

284 **CoreFed.** Observe that for convex losses, we can directly solve this maximization to the optimal to  
 285 achieve core-stability. For non-convex losses such as those for DNNs, we apply gradient descent to  
 286 maximize the objective. Suppose we are training on  $n$  finite samples  $\{(x_i, y_i)\}_{i \in [n]}$  drawn from the  
 287 data distribution  $\mathcal{P}$ , which constitute empirical distribution  $\hat{\mathcal{P}}_n$ . We observe that, the gradient can be  
 288 expressed as an conical combination of the gradients of each group:

$$\nabla_{\theta} \mathcal{L}(\theta) = \sum_{i \in [n]} \frac{\nabla_{\theta} u_i(\theta)}{u_i(\theta)} = \sum_{i \in [n]} \frac{-\nabla_{\theta} \mathbb{E}_{(x,y) \sim \hat{\mathcal{P}}_n^{(i)}} \ell(f_{\theta}(x), y)}{M_i - \mathbb{E}_{(x,y) \sim \hat{\mathcal{P}}_n^{(i)}} \ell(f_{\theta}(x), y)}. \quad (3)$$

289 Therefore, for each group, we reweight its gradients or weight updates by  $(M_i -$   
 290  $\mathbb{E}_{(x,y) \sim \hat{\mathcal{P}}_n^{(i)}} \ell(f_{\theta}(x), y))^{-1}$  and then sum up to get the final weight update in each iteration, which  
 291 leads to a distributed training protocol shown in Algorithm 1.

292 This protocol is similar to standard FedAvg. However, in CoreFed the model weight updates are  
 293 weighted then aggregated at each iteration, while in FedAvg, model weights are directly averaged  
 294 and aggregated at each iteration. In the limit that each local update uses single step with entire  
 295 dataset,  $\Delta \theta_s = -\eta \frac{1}{|\mathcal{D}_s|} \sum_{i=1}^{|\mathcal{D}_s|} \nabla_{\theta^t} \ell(f_{\theta^t}(x_s^{(i)}), y_s^{(i)})$ , where  $\mathcal{D}_s = \{(x_s^{(i)}, y_s^{(i)}) : 1 \leq i \leq |\mathcal{D}_s|\}$  is  
 296 the training dataset on device  $s$ . Therefore, the global update is a unbiased gradient descent step of the  
 297 objective  $\sum_s \log(M_s - \frac{1}{|\mathcal{D}_s|} \sum_{i=1}^{|\mathcal{D}_s|} \ell(f_{\theta^t}(x_s^{(i)}), y_s^{(i)})) = \mathcal{L}(\theta_t)$  where  $\mathcal{L}(\cdot)$  is defined in 2.

### 298 4.3 Core-Stability in Linear Regression and Classification with Logistic Regression.

299 We now discuss some of the predictive models, where the concavity requirements of the utility  
 300 function is satisfied. Note that a necessary and sufficient condition for  $u_i(\cdot)$  to be concave is that  
 301 the loss function  $\ell(\cdot)$  should be convex in  $c$ . Here we elaborate few scenarios where this is true.  
 302 Suppose we are training on  $n$  finite samples  $\{(x_i, y_i)\}_{i \in [n]}$  drawn from the data distribution  $\mathcal{P}$ , which  
 303 constitute empirical distribution  $\hat{\mathcal{P}}_n$ .

304 **Linear Regression.** These are the scenarios where we map our input variables to a real number  
 305 (not discrete class labels). In this case, we have  $f_{\theta}(x) = \theta^{\top} x$ . Observe that the *regression loss*, then  
 306  $\mathbb{E}_{(x,y) \sim \hat{\mathcal{P}}_n} \ell(f_{\theta}(x), y) = \mathbb{E}_{(x,y) \sim \hat{\mathcal{P}}_n} (\theta^{\top} x - y)^2$  would be convex in  $\theta$ .

307 **Lemma 4.3.** *CoreFed determines a core-stable predictor in a federated learning setting training*  
 308 *linear regression.*

309 **Classification with Logistic Regression.** In classification tasks we map the input variables to dis-  
 310 crete class labels. A commonly used loss function in classification is logistic regression. Given a clas-  
 311 sifier  $\theta$  and a scalar  $c \in \mathbb{R}$ , an agent  $i$ 's loss is given by  $\ell_i(\theta, c) = \frac{\|\theta\|_2}{2} + \alpha \cdot \sum_{i \in [n]} \log(e^{-y_i(\theta^{\top} x_i + c)} +$   
 312  $1)$  [28, 2, 12]. It is well known that  $\ell_i(\theta, c)$  is convex (see, e.g., [12]). Thus,  $u_i(\theta, c) = M_i - \ell_i(\theta, c)$   
 313 where  $M_i = \arg \max_{\theta \in P, c \in \mathbb{R}} (\ell_i(\theta, c))$ , is concave.

314 **Lemma 4.4.** *CoreFed determines a core-stable predictor in a federated learning setting training*  
 315 *classification with logistic regression.*

### 316 4.4 Approximate Core-Stability in Deep Neural Networks

317 Theorem 2 requires that  $u_i(\theta)$  is concave in terms of  $\theta$  and global optimality for the objective  
 318  $\mathcal{L}(\theta) = \prod_{i \in [n]} u_i(\theta)$  to achieve core-stability. However, these two conditions are challenging to be  
 319 satisfied for DNNs, where the training loss is non-convex and common training methods, which are  
 320 based on first-order gradients, are not guaranteed to absolutely converge. In this more general scenario,  
 321 the following theorem shows the relaxed local core-stability that we can attain for approximately  
 322 first-order converged predictor (i.e., predictor with small local gradient  $\|\nabla_{\theta} \mathcal{L}(\theta)\|_2 \leq \epsilon$ ).

323 **Definition 5.** A predictor  $\theta \in P$ , is called  $(d, k)$ -pseudo core stable, where  $d > 0, k > 1$  if there  
 324 exists no other  $\theta' \in P$  such that  $\|\theta' - \theta\|_2 < d$ , and no subset  $S$  of agents, such that  $\frac{|S|}{kn} u_i(\theta') \geq u_i(\theta)$   
 325 for all  $i \in S$ , with at least one strict inequality.

326 **Theorem 3.** For all  $i \in [n]$ , if  $u_i(\theta)$  is  $\beta$ -smooth in terms of  $\theta$  within  $\{\theta' : \|\theta - \theta'\|_2 \leq d\}$ , and  
 327  $\|\nabla_{\theta} \mathcal{L}(\theta)\|_2 \leq \epsilon$ , then  $\theta$  is a  $(d, k)$ -pseudo core stable predictor, where

$$d = \frac{-\epsilon + \sqrt{\epsilon^2 + 2\beta(k-1)n \sum_{i \in [n]} u_i(\theta)^{-1}}}{\beta \sum_{i \in [n]} u_i(\theta)^{-1}}. \quad (4)$$

328 **Implications.** We defer the proof to Appendix E. Theorem 3 states that, for smooth neural net-  
 329 works, there exists no predictor  $\theta'$  in the neighborhood with  $\ell_2$  radius  $d$ , that any subset of agents  
 330 prefer “significantly”. Although our guarantees are local guarantees, we remark that global fairness  
 331 guarantees are unlikely for DNNs. Most of the fairness guarantees in computational social choice  
 332 and game theory crucially require the agents to have convex preferences, i.e., the level sets of the  
 333 utility functions need to be convex. There are impossibility results for fairness when the agent’s  
 334 preferences are non-convex. However, while non-convex consumer preferences are not interesting  
 335 from an economic standpoint, our current work finds an application for these preferences in fairness  
 336 in federated learning with DNNs.

## 337 4.5 Weighted Core-Stability

338 In this section, we show how to generalize all of our results (Theorems 1, 2, and 3) when we want to  
 339 train the joint predictor to fit the data of certain agents more than some others. In particular, for each  
 340 agent  $i$ , if we assign weight  $w_i$ , indicating the desired bias of the final trained model towards agent  $i$ <sup>5</sup>,  
 341 then with subtle modifications, we can show the existence of a *weighted core stable* predictor, when  
 342 the utility functions of the agents satisfy the conditions in Theorem 1. Formally,

343 **Definition 6** (Weighted Core-Stability). Given the weight vector  $w = \langle w_1, w_2, \dots, w_n \rangle$ , a predictor  
 344  $\theta \in P$ , is weighted core-stable if and only if there exists no other predictor  $\theta' \in P$  and a subset of  
 345 agents  $S \subseteq [n]$  such that  $\frac{\sum_{j \in S} w_j}{\sum_{j \in [n]} w_j} \cdot u_i(\theta') \geq u_i(\theta)$  for all  $i \in S$  with at least one strict inequality.

346 Note that, when all the agents have the same weight, e.g.,  $w_i = 1, \forall i \in [n]$ , then weighted core-  
 347 stability matches core-stability. At a high-level the concept is the same, no group of agents can  
 348 significantly benefit by forming a coalition within themselves. However “significantly” means a  
 349 multiplicative increase by  $\frac{\sum_{j \in [n]} w_j}{\sum_{j \in S} w_j}$  (instead of  $|S|/n$  for the unweighted case), i.e., it is dependent  
 350 on the total weight of the set  $S$ . We make the aforementioned guarantee more intuitive by considering  
 351 special cases of  $S$ . When  $S = \{i\}$ , our guarantees say that agent  $i$  gets a utility of  $w_i / (\sum_{j \in [n]} w_j)$   
 352 fraction of her maximum utility, i.e., the utility of agents with higher weights are prioritized. We call  
 353 this *weighted proportionality*. Also note that by setting  $S = [n]$ , we get Pareto-optimality (similar to  
 354 the unweighted case).

355 Furthermore, by changing the convex program 2 to maximizing  $\sum_{j \in [n]} w_j \log(u_j(\theta))$  instead of  
 356  $\sum_{j \in [n]} \log(u_j(\theta))$ , we can get the weighted version of Theorem 2. This also suggests a simple  
 357 generalization of CoreFed to Weighted-CoreFed and all our extensions in Sections 4.3 and 4.4 will  
 358 also generalize to the weighted setting.

## 359 5 Empirical Evaluation

360 We evaluate our fair ML method CoreFed and baseline FedAvg [22] on three datasets (Adult, MNIST  
 361 and CIFAR-10) on linear model and deep neural networks. We show that the model trained with  
 362 CoreFed is able to achieve core-stable fairness, while maintaining similar utility with the standard  
 363 FedAvg protocol, which cannot guarantee to achieve core-stable fairness.

<sup>5</sup>Following the light of [7], one possible candidate can be  $w_i = \mathcal{D}_i$ , i.e., set  $w_i$  to the size of the data shared  
 by agent  $i$  with the model.



364 **5.1 Experiment Setup**

365 **Dataset.** We evaluate our algorithm CoreFed on Adult [3], MNIST [18] and CIFAR-10 [17] datasets.  
 366 To perform federated learning on heterogeneous data, we construct the non-IID data by sampling the  
 367 proportion of each label from Dirichlet distribution for every agent, following the literature [19].

368 **Models.** We train a logistic regression classifier on Adult data. We use a CNN, which has two 5x5  
 369 convolution layers followed by 2x2 max pooling and two fully connected layer with ReLU activation  
 370 for MNIST and CIFAR-10. For Adult dataset, the utility is selected as  $M - \ell_{log}$  where  $\ell_{log}$  is the  
 371 logistic loss. For CIFAR-10 and MNIST, we use cross entropy loss  $\ell_{ce}$  as the training loss with utility  
 372  $U$  becomes  $M - \ell_{ce}$ .  $M$  is set to be 3.0, 1.0 and 3.0 for Adult, MNIST, and CIFAR-10, respectively,  
 373 based on statistical analysis during training. All experiments are conducted on a 1080 Ti GPU.

374 **5.2 Evaluation Results**

375 We demonstrate that our CoreFed distributed training protocol in Algorithm 1 achieves the core-stable  
 376 fairness through comparison with FedAvg on different datasets and settings. Concretely, we perform  
 377 training with FedAvg and our proposed CoreFed, and then validate whether the utility inequality  
 378  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} < n$  (see *Implications* after Theorem 2) holds under different settings. Here we treat the  
 379 model trained by FedAvg parameterized by  $\theta'$ , while the model trained by our CoreFed parameterized  
 380 by  $\theta^*$ . That is to say, since the model trained by CoreFed achieves core-stable fairness, we expect  
 381 the model parameterized by  $\theta$  would have pareto-optimality. Indeed, results in Table 1 suggest that  
 382 CoreFed achieves core-stable fairness compared with FedAvg while maintaining similar utility. We  
 383 also report the average and multiplicative utility of the trained global model in "U(Average)" and  
 384 "U(Multi)" columns. We can see that CoreFed achieves higher overall utilities, especially for the  
 385 multiplicative case since FedAvg favors the average case in general.

Table 1: Comparison of the utility ( $M - \ell_{ce}$ ) of each agent trained with CoreFed and FedAvg. We see that  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} < n$  holds, where  $\theta'$  denotes the weights of shared model trained by FedAvg and  $\theta^*$  by CoreFed.

Dataset	Method	Agent 0	Agent 1	Agent 2	U(Average)	U(Multi)	$\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)}$
Adult	FedAvg	2.59	0.77	1.46	1.61	2.91	2.80 (<3)
	CoreFed	2.62	0.90	1.53	1.68	3.61	
MNIST	FedAvg	0.34	0.29	0.92	0.52	0.091	2.66 (<3)
	CoreFed	0.36	0.41	0.91	0.56	0.13	
CIFAR-10	FedAvg	0.63	1.40	0.51	0.84	0.45	2.62 (<3)
	CoreFed	0.73	1.35	0.71	0.93	0.70	

386 **6 Conclusion**

387 In this work, we aim to bridge the algorithmic fairness with game theory and social choice theory by  
 388 formally defining the core-stable fairness in federated learning, especially for the non-IID setting. We  
 389 prove the core stability of our proposed fair FL algorithm with empirical validation. We believe this  
 390 work would open up new research directions on connecting game theoretic analysis and statistical  
 391 machine learning under different learning paradigms, objective, and utilities.

## References

- 392 [1] [https://en.wikipedia.org/wiki/Core\\_\(game\\_theory\)](https://en.wikipedia.org/wiki/Core_(game_theory)).
- 393 [2] Alan Agresti. *Categorical data analysis*. John Wiley & Sons, 2003.
- 394 [3] Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- 395 [4] Haris Aziz. Justifications of welfare guarantees under normalized utilities. *SIGecom Exch.*,  
396 17(2):71–75, 2019.
- 397 [5] Eric Budish. The combinatorial assignment problem: Approximate competitive equilibrium  
398 from equal incomes. *Journal of Political Economy*, 119(6):1061–1103, 2011.
- 399 [6] Kate Donahue and Jon M. Kleinberg. Models of fairness in federated learning. *CoRR*,  
400 abs/2112.00818, 2021.
- 401 [7] Kate Donahue and Jon M. Kleinberg. Optimality and stability in federated learning: A game-  
402 theoretic approach. In *NeurIPS*, pages 1287–1298, 2021.
- 403 [8] Wei Du, Depeng Xu, Xintao Wu, and Hanghang Tong. Fairness-aware agnostic federated  
404 learning. In *SDM*, pages 181–189. SIAM, 2021.
- 405 [9] Brandon Fain, Kamesh Munagala, and Nisarg Shah. Fair allocation of indivisible public goods.  
406 In *EC*, pages 575–592. ACM, 2018.
- 407 [10] Duncan Karl Foley. *Resource allocation and the public sector*. Yale University, 1966.
- 408 [11] Laurent Gourvès, Jérôme Monnot, and Lydia Tlili. Near fairness in matroids. In *ECAI*,  
409 volume 263 of *Frontiers in Artificial Intelligence and Applications*, pages 393–398. IOS Press,  
410 2014.
- 411 [12] David W Hosmer Jr, Stanley Lemeshow, and Rodney X Sturdivant. *Applied logistic regression*,  
412 volume 398. John Wiley & Sons, 2013.
- 413 [13] Tiansheng Huang, Weiwei Lin, Wentai Wu, Ligang He, Keqin Li, and Albert Y. Zomaya. An  
414 efficiency-boosting client selection scheme for federated learning with fairness guarantee. *IEEE*  
415 *Trans. Parallel Distributed Syst.*, 32(7):1552–1564, 2021.
- 416 [14] Wei Huang, Tianrui Li, Dexian Wang, Shengdong Du, Junbo Zhang, and Tianqiang Huang.  
417 Fairness and accuracy in horizontal federated learning. *Information Sciences*, 589:170–185,  
418 2022.
- 419 [15] Shizuo Kakutani. A generalization of Brouwer’s fixed point theorem. *Duke mathematical*  
420 *journal*, 8(3):457–459, 1941.
- 421 [16] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,  
422 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*  
423 *preprint arXiv:1610.05492*, 2016.
- 424 [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
425 2009.
- 426 [18] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
427 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 428 [19] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data  
429 silos: An experimental study. *arXiv preprint arXiv:2102.02079*, 2021.
- 430 [20] Tian Li, Maziari Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in  
431 federated learning. In *ICLR*. OpenReview.net, 2020.
- 432 [21] Richard J. Lipton, Evangelos Markakis, Elchanan Mossel, and Amin Saberi. On approximately  
433 fair allocations of indivisible goods. In *5th*, pages 125–131, 2004.
- 434

- 435 [22] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y  
436 Arcas. Communication-efficient learning of deep networks from decentralized data. 2017.
- 437 [23] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In  
438 *International Conference on Machine Learning*, pages 4615–4625. PMLR, 2019.
- 439 [24] Herve Moulin. Fair division in the internet age. *Annual Review of Economics*, 11, 2019.
- 440 [25] Thomas J Muench. The core and the lindahl equilibrium of an economy with a public good: An  
441 example. *Journal of Economic Theory*, 4(2):241–255, 1972.
- 442 [26] Afroditi Papadaki, Natalia Martínez, Martín Bertrán, Guillermo Sapiro, and Miguel R. D.  
443 Rodrigues. Federating for learning group fair models. *CoRR*, abs/2110.01999, 2021.
- 444 [27] Benjamin Plaut and Tim Roughgarden. Almost envy-freeness with general valuations. *SIAM J.*  
445 *Discret. Math.*, 34(2):1039–1068, 2020.
- 446 [28] scikit-learn developers. 1.1. linear models — scikit-learn 1.1.0 documentation. [https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression),  
447 2022.  
448
- 449 [29] Yuxin Shi, Han Yu, and Cyril Leung. A survey of fairness-aware federated learning. *CoRR*,  
450 abs/2111.01872, 2021.
- 451 [30] Tianshu Song, Yongxin Tong, and Shuyue Wei. Profit allocation for federated learning. In *2019*  
452 *IEEE International Conference on Big Data (Big Data)*, pages 2577–2586, 2019.
- 453 [31] Hugo Steinhaus. The problem of fair division. *Econometrica*, 16:101–104, 1948.
- 454 [32] Xinyi Xu and Lingjuan Lyu. Towards building a robust and fair federated learning system.  
455 *CoRR*, abs/2011.10464, 2020.
- 456 [33] Miao Yang, Ximin Wang, Hongbin Zhu, Haifeng Wang, and Hua Qian. Federated learning with  
457 class imbalance reduction. In *2021 29th European Signal Processing Conference (EUSIPCO)*,  
458 pages 2174–2178. IEEE, 2021.
- 459 [34] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi.  
460 Fairness beyond disparate treatment & disparate impact: Learning classification without dis-  
461 parate mistreatment. In *WWW*, pages 1171–1180. ACM, 2017.
- 462 [35] Yuchen Zeng, Hongxu Chen, and Kangwook Lee. Improving fairness via federated learning.  
463 *CoRR*, abs/2110.15545, 2021.
- 464 [36] Fengda Zhang, Kun Kuang, Yuxuan Liu, Chao Wu, Fei Wu, Jiaxun Lu, Yunfeng Shao, and Jun  
465 Xiao. Unified group fairness on federated learning. *CoRR*, abs/2111.04986, 2021.

## 466 Checklist

- 467 1. For all authors...
- 468 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
469 contributions and scope? [Yes]
- 470 (b) Did you describe the limitations of your work? [Yes] We clearly specify the assumptions  
471 and application scopes, if exists, of all our results.
- 472 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We discuss  
473 broader impacts in Appendix A.
- 474 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
475 them? [Yes]
- 476 2. If you are including theoretical results...
- 477 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 478 (b) Did you include complete proofs of all theoretical results? [Yes] We include complete  
479 proofs for all results in the appendix.

- 480 3. If you ran experiments...
- 481 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
482 mental results (either in the supplemental material or as a URL)? [Yes]
- 483 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
484 were chosen)? [Yes] We include them in Section 5.1
- 485 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
486 ments multiple times)? [Yes] Indeed, we rigorously compute the confidence intervals  
487 due to the finite sampling (see Section 5.2), and guarantee that all fairness certificates  
488 reported in the paper holds with probability  $\geq 90\%$ .
- 489 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
490 of GPUs, internal cluster, or cloud provider)? [Yes] We provide them in Section 5.1
- 491 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 492 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 493 (b) Did you mention the license of the assets? [Yes]
- 494 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 495 (d) Did you discuss whether and how consent was obtained from people whose data you're  
496 using/curating? [No] We only use public data.
- 497 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
498 information or offensive content? [No] We only use public data.
- 499 5. If you used crowdsourcing or conducted research with human subjects...
- 500 (a) Did you include the full text of instructions given to participants and screenshots, if  
501 applicable? [N/A]
- 502 (b) Did you describe any potential participant risks, with links to Institutional Review  
503 Board (IRB) approvals, if applicable? [N/A]
- 504 (c) Did you include the estimated hourly wage paid to participants and the total amount  
505 spent on participant compensation? [N/A]

506 **A Broader Impact**

507 This paper aims to provide a fair federated learning algorithm to guarantee that the utilities of the  
 508 trained agents are core-stable fair. We do not expect the work to have any ethics issues or negative  
 509 social impact if it is correctly used. On the other hand, if our evaluation and theory is misused, there  
 510 could be potential negative social impact. For instance, our fairness metrics cannot indicate other  
 511 accuracy or loss utilities and people need to evaluate federated learning algorithms with different  
 512 utility metrics, rather than only using our metrics. We expect that our work will provide a way to  
 513 measure and achieve fairness for different federated learning paradigms.

514 **B Missing proofs from Section 4.1**

515 **B.1 Proof of Lemma 4.1**

516 *Proof.* We prove by contradiction. Assume that there exists a  $\theta' \in P$ , and a  $S \subseteq [n]$ , such that  
 517  $(|S|/n) \cdot u_i(\theta') \geq u_i(\theta)$  for all  $i \in S$  with at least one strict inequality. Then, we have  $\frac{u_i(\theta')}{u_i(\theta)} \geq n/|S|$   
 518 for all  $i \in S$  with at least one strict inequality, implying  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} \geq \sum_{i \in S} \frac{u_i(\theta')}{u_i(\theta)} > n$ . However,  
 519 since  $\theta \in \phi(\theta)$ , we have  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} \leq \sum_{i \in [n]} \frac{u_i(\theta)}{u_i(\theta)} = n$ , which is a contradiction.  $\square$

520 **B.2 Proof of Lemma 4.2**

521 *Proof.* We need to show that for every sequence  $(\theta)_i$  converging to  $\theta$ , and  $(\gamma)_i$  converging to  $\gamma$ , such  
 522 that  $\gamma_i \in \phi(\theta_i)$  for all  $i$ , we have  $\gamma \in \phi(\theta)$ . We prove this by contradiction. Let us assume otherwise,  
 523  $\gamma \notin \phi(\theta)$ . Let  $\gamma' \in \phi(\theta)$  and let  $\delta = \frac{\sum_{i \in [n]} u_i(\gamma')/u_i(\theta)}{\sum_{i \in [n]} u_i(\gamma)/u_i(\theta)} > 1$ . We now make a technical claim about  
 524 the utility functions of the agents.

525 **Claim B.1.** For all  $i \in [n]$ , and  $x, y$  such that  $\|x - y\|_2 \leq \beta$ , we have

- 526 1.  $|u_i(x) - u_i(y)| \leq h(\beta)$  where  $h: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  is continuous increasing function with  
 527  $h(0) = 0$ , and
- 528 2. for each  $i \in [n]$ , we have  $u_i(y) \cdot h'(\beta)^{-1} \leq u_i(x) \leq u_i(y) \cdot h'(\beta)$ , where  $h'(\beta) = (1 + \frac{h(\beta)}{M\varepsilon})$   
 529 and  $M = \min_{i \in [n]} M_i$ .

530 *Proof.* Claim (1) follows immediately from the continuity of the utility functions.

531 For claim (2), we have

$$\begin{aligned}
 u_i(x) &\leq u_i(y) + h(\beta) \\
 &\leq u_i(y) \cdot \left(1 + \frac{h(\beta)}{u_i(y)}\right) \\
 &\leq u_i(y) \cdot \left(1 + \frac{h(\beta)}{M_i\varepsilon}\right) && (u_i(y) \geq M_i\varepsilon) \\
 &\leq u_i(y) \cdot \left(1 + \frac{h(\beta)}{M\varepsilon}\right) && (M_i \geq M) \\
 &\leq u_i(y) \cdot h'(\beta).
 \end{aligned}$$

532 In a similar way, we can prove that  $u_i(y) \leq u_i(x) \cdot h'(\beta)$ , which would then imply that  $u_i(x) \geq$   
 533  $u_i(y) \cdot (h'(\beta))^{-1}$ .  $\square$

534 We choose a  $\delta'$  such that  $h'(\delta')^3 = (1 + \frac{h(\delta')}{M\varepsilon})^3 \ll \delta$ . Such a  $\delta'$  exists as  $h()$  is a continuous  
 535 increasing function with  $h(0) = 0$ , and  $\delta > 1$ . Since the sequences  $(\theta)_i$  and  $(\gamma)_i$  converges to  $\theta$   
 536 and  $\gamma$  respectively, there exists a  $n' \in \mathbb{N}$  such that for all  $\ell \geq n'$ , we have  $\|\gamma_\ell - \gamma\|_2 < \delta'$  and

537  $\|\theta_\ell - \theta\|_2 < \delta'$ . Now observe that

$$\begin{aligned}
\sum_{i \in [n]} \frac{u_i(\gamma')}{u_i(\theta_\ell)} &\geq h'(\delta')^{-1} \cdot \sum_{i \in [n]} \frac{u_i(\gamma')}{u_i(\theta)} && \text{(by Claim B.1)} \\
&= h'(\delta')^{-1} \cdot \delta \cdot \sum_{i \in [n]} \frac{u_i(\gamma)}{u_i(\theta)} && \text{(by definition of } \delta) \\
&\geq h'(\delta)^{-2} \cdot \delta \cdot \sum_{i \in [n]} \frac{u_i(\gamma_\ell)}{u_i(\theta)} && \text{(by Claim B.1)} \\
&\geq h'(\delta)^{-3} \cdot \delta \cdot \sum_{i \in [n]} \frac{u_i(\gamma_\ell)}{u_i(\theta_\ell)} && \text{(by Claim B.1)} \\
&> \sum_{i \in [n]} \frac{u_i(\gamma_\ell)}{u_i(\theta_\ell)} && \text{(as } \delta \gg h'(\delta')^3).
\end{aligned}$$

538 This shows that  $\gamma_\ell \notin \phi(\theta_\ell)$ , which is a contradiction.  $\square$

## 539 C Missing Proofs from Section 4.2

### 540 C.1 Proof of Theorem 2

541 *Proof.* We first show that for any other predictor  $\theta' \in P$ , we have  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} \leq n$ . Consider any  
542 other predictor  $\theta' \in P$ . Since  $P$  is convex, we have  $(\nabla_\theta \mathcal{L}(\theta^*))^T (\theta' - \theta^*) < 0$ . Now, observe that

$$\begin{aligned}
\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} - n &= \sum_{i \in [n]} \frac{u_i(\theta') - u_i(\theta^*)}{u_i(\theta^*)} \\
&\leq \sum_{i \in [n]} \frac{(\nabla u_i(\theta^*))^T (\theta' - \theta^*)}{u_i(\theta^*)} && \text{(from concavity of } u_i(\cdot)) \\
&= \sum_{i \in [n]} \sum_{j \in [d]} \left( \frac{\partial u_i(\theta^*)}{\partial \theta_j} \cdot (\theta'_j - \theta_j^*) \cdot \frac{1}{u_i(\theta^*)} \right) \\
&= \sum_{i \in [n]} \frac{1}{u_i(\theta^*)} \cdot \sum_{j \in [d]} \left( \frac{\partial u_i(\theta^*)}{\partial \theta_j} \cdot (\theta'_j - \theta_j^*) \right) \\
&= \sum_{j \in [d]} (\theta'_j - \theta_j^*) \cdot \sum_{i \in [n]} \left( \frac{1}{u_i(\theta^*)} \cdot \frac{\partial u_i(\theta^*)}{\partial \theta_j} \right) \\
&= (\nabla_\theta \mathcal{L}(\theta^*))^T (\theta' - \theta^*) < 0.
\end{aligned}$$

543 Now if  $\theta^*$  is not core-stable, then there exists an  $S \subseteq [n]$  and  $\theta' \in P$ , such that  $u_i(\theta') \geq n/|S| u_i(\theta^*)$   
544 for all  $i \in S$  with at least one strict inequality, then we have  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta^*)} \geq \sum_{i \in S} \frac{u_i(\theta')}{u_i(\theta^*)} >$   
545  $n/|S| \cdot |S| = n$ , which is a contradiction.  $\square$

## 546 D Algorithm CoreFed

547 Here we present the full description of CoreFed in Algorithm 1.

## 548 E Missing Proofs from Section 4.4

### 549 E.1 Proof of Theorem 3

550 *Proof.* For any  $\theta'$  such that  $\|\theta - \theta'\|_2 \leq d$ , according to the definition of  $\beta$ -smooth, we have

$$u_i(\theta') \leq u_i(\theta) + \nabla_\theta u_i(\theta)^T (\theta' - \theta) + \frac{\beta}{2} \|\theta' - \theta\|_2^2.$$

---

**Algorithm 1:** CoreFed Distributed Training Protocol.

---

**Input:** Number of clients  $K$ , number of rounds  $T$ , epochs  $E$ , learning rate  $\eta$ **Output:** Model weights  $\theta^T$ 

1 **for**  $t = 0, 1, \dots, T - 1$  **do**  
2     Server selects a subset of  $K$  devices  $S_t$ ;  
3     Server sends weights  $\theta^t$  to all selected devices;  
4     Each select device  $s \in S_t$  updates  $\theta^t$  for  $E$  epochs of SGD with learning rate  $\eta$  to obtain new weights  $\bar{\theta}_s^t$ ;  
5     Each select device  $s \in S_t$  computes

$$\Delta\theta_s^t = \bar{\theta}_s^t - \theta^t,$$
$$\mathcal{L}_s^t = \frac{1}{|\mathcal{D}_s|} \sum_{i=1}^{|\mathcal{D}_s|} \ell(f_{\theta^t}(x_s^{(i)}), y_s^{(i)})$$

   where  $\mathcal{D}_s = \{(x_s^{(i)}, y_s^{(i)}) : 1 \leq i \leq |\mathcal{D}_s|\}$  is the training dataset on device  $s$ ;  
6     Each selected device  $s \in S_t$  sends  $\Delta\theta_s^t$  and  $\mathcal{L}_s^t$  back to the server;  
7     Server updates  $\theta^{t+1}$  following

$$\theta^{t+1} \leftarrow \theta^t + \frac{1}{|S_t|} \sum_{s \in S_t} \frac{\Delta\theta_s^t}{M_s - \mathcal{L}_s^t}.$$

8 **end**

---

551 Then we observe that

$$\begin{aligned} & \sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} - kn \\ &= \sum_{i \in [n]} \frac{u_i(\theta') - u_i(\theta)}{u_i(\theta)} - (k-1)n \\ &\leq \sum_{i \in [n]} \frac{\nabla_{\theta} u_i(\theta)^{\top} (\theta' - \theta) + \frac{\beta}{2} \|\theta' - \theta\|_2^2}{u_i(\theta)} - (k-1)n \\ &= (\nabla_{\theta} \mathcal{L}(\theta))^{\top} (\theta' - \theta) + \sum_{i \in [n]} \frac{\beta}{2u_i(\theta)} \|\theta' - \theta\|_2^2 - (k-1)n \\ &\leq \epsilon \|\theta' - \theta\|_2 + \sum_{i \in [n]} \frac{\beta}{2u_i(\theta)} \|\theta' - \theta\|_2^2 - (k-1)n. \end{aligned}$$

552 By plugging in the RHS of Equation (4), we observe that when  $\|\theta' - \theta\|_2 < d$ ,  $\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} - kn < 0$ .553 On the other hand, suppose for any  $S \subseteq [n]$ , if for all  $i \in S$  we have  $\frac{|S|}{kn} u_i(\theta') \geq u_i(\theta)$ , then

$$\sum_{i \in [n]} \frac{u_i(\theta')}{u_i(\theta)} = \sum_{i \in S} \frac{u_i(\theta')}{u_i(\theta)} + \sum_{i \in [n] \setminus S} \frac{u_i(\theta')}{u_i(\theta)} \geq \sum_{i \in S} \frac{u_i(\theta')}{u_i(\theta)} \geq \sum_{i \in S} \frac{kn}{|S|} \geq kn \quad (5)$$

554 which contradicts the above result and concludes the proof.  $\square$