

---

000 FEDERATED EQUILIBRIUM SOLUTIONS FOR  
001 GENERALIZED METHOD OF MOMENTS APPLIED TO  
002 INSTRUMENTAL VARIABLE ANALYSIS  
003  
004  
005

006 **Anonymous authors**

007 Paper under double-blind review  
008  
009

010  
011 ABSTRACT  
012

013 Instrumental variables (IV) analysis is an important applied tool for areas such as  
014 healthcare and consumer economics. For IV analysis in high-dimensional settings,  
015 the Generalized Method of Moments (GMM) using deep neural networks offer an  
016 efficient approach. With non-i.i.d. data sourced from scattered decentralized  
017 clients, federated learning is a popular paradigm for training the models while  
018 promising data privacy. However, to our knowledge, no federated algorithm for  
019 either GMM or IV analysis exists to date. In this work, we introduce federated  
020 IV analysis (FEDIV) via federated GMM (FEDGMM). We formulate FEDGMM  
021 as a federated zero-sum game defined by a non-convex non-concave minimax  
022 optimization problem. We characterize the solutions to the federated game using  
023 Stackelberg equilibrium and show that it satisfies client-local equilibria up to a  
024 heterogeneity bias. Thereby, we show that the consistency of federated GMM  
025 estimator across clients closely depends on the heterogeneity bias. Our experiments  
026 demonstrate that the federated framework for IV analysis efficiently recover the  
027 consistent GMM estimators for low and high-dimensional data.  
028

029  
030 1 INTRODUCTION  
031

032 Federated Learning (FL) (McMahan et al., 2017) is now an established paradigm for training Machine  
033 Learning (ML) models over decentralized clients, keeping the data local and private. The applications  
034 include important domains such as healthcare (Oh & Nadkarni, 2023), finance & banking (Long  
035 et al., 2020), smart cities & mobility (Gecer & Garbinato, 2024), and many others (Ye et al., 2023).  
036 The scale of FL has also grown large – see the Nature Medicine report by Dayan et al. (2021) on a  
037 global-scale FL to predict the effectiveness of oxygen administration to COVID-19 patients in the  
038 emergency rooms while maintaining data locality. However, the current popular FL methods have a  
039 crucial limitation due to their standard supervised nature of learning. For example, Liang et al. (2023)  
040 suggests that the hypoxia-inducible factors (HIF) (a protein that controls the rate of transcription of  
041 genetic information from DNA to messenger RNA by binding to a specific DNA sequence) play a  
042 vital role in oxygen consumption at the cellular level. Arguably, the Dayan et al. (2021)’s approach  
043 may over- or under-estimate the effects of oxygen treatment as it does not accommodate the influence  
044 of HIF levels on oxygen consumption.

045 A classical approach to address such limitations is Instrumental variables (IV) analysis, which assumes  
046 conditional independence between a confounding variable and the outcome while considering its  
047 causal effect on the treatment variable. IV analysis can very practically apply to training Dayan  
048 et al. (2021)’s ML model wherein the patients’ HIF levels work as an IV that influences the effective  
049 organ-level oxygen consumption (a treatment variable) but does not directly affect the mortality of the  
050 COVID-19 patients (the outcome). IV analysis has been comprehensively explored in econometrics  
051 (Angrist & Krueger, 2001; Angrist & Pischke, 2009) with several decades of history, such as works of  
052 Wright (1928) and Reiersøl (1945). Its efficiency is now accepted for learning even high-dimensional  
053 complex causal relationships, such as those in image datasets (Hartford et al., 2017; Bennett et al.,  
2019). Naturally, the growing demand for FL entails designing methods for federated IV analysis,  
which, to our knowledge, is yet unexplored.

054 In the centralized deep learning setting, Hartford et al. (2017) introduced an IV analysis framework,  
055 namely DEEPIV, which uses two stages of neural networks (NN) training – learn the conditional  
056 treatment distribution as a NN-parametrized Gaussian mixture for the treatment prediction and  
057 then train the outcome model. The two-stage process has precursors in applying the least square  
058 regressions in the two phases (Angrist & Pischke, 2009)[4.1.1]. In the same setting, another approach  
059 for IV analysis applies the generalized method of moments (GMM) (Wooldridge, 2001). GMM is  
060 a celebrated estimation approach in social sciences and economics. It was introduced by Hansen  
061 (1982), for which he won a Nobel Prize in Economics (Steif et al., 2014).

062 Lewis & Syrgkanis (2018) employed neural networks for GMM estimation. Their method, called  
063 the adversarial generalized method of moments (AGMM) fit a GMM criterion function over a finite  
064 set of unconditional moments. Similarly, Bennett et al. (2019) introduced deep learning models to  
065 GMM estimation; they named their method DEEPGMM. DEEPGMM differs from AGMM in using  
066 a weighted norm to define the objective function. The experiments in (Bennett et al., 2019) showed  
067 that DEEPGMM outperformed AGMM for IV analysis, and both won against DEEPIV. Nonetheless,  
068 to our knowledge, none of these methods has a federated counterpart. Notably, both AGMM and  
069 DEEPGMM translate to a minimax optimization problem corresponding to a smooth zero-sum game.

070 The zero-sum game formulated for GMM estimation is essentially nonconvex-nonconcave (Bennett  
071 et al., 2019). Such a game corresponds to a sequential game as it may have differing maximin and  
072 minimax solutions. Unlike Nash equilibrium, the global minimax points – the Stackelberg equilibria –  
073 are guaranteed to exist for nonconvex-nonconcave games. However, finding a global minimax point  
074 is generally NP-hard, necessitating solving for a surrogate local equilibrium (Jin et al., 2020).

075 Now, considering the federated version of this problem, a fundamental challenge arises in establishing  
076 that a federated minimax optimization algorithm retrieves a local Stackelberg equilibrium of the  
077 federated zero-sum game. Even if it did, it requires showing that the federated equilibrium translates to  
078 the client-local setting under heterogeneity. Finally, it entails proving that the client-local equilibrium  
079 under heterogeneity is a consistent GMM estimator for its data. In this work, we address these  
080 challenges. Our contributions are summarized as follows:

- 081 1. We introduce **FEDIV**: federated IV analysis. To our knowledge, **FEDIV** is the first work on IV  
082 analysis in a federated setting.
- 083 2. We present **FEDDEEPGMM**<sup>1</sup> – a federated adaptation of DEEPGMM of Bennett et al. (2019) to  
084 solve FEDIV. FEDDEEPGMM is implemented as a federated smooth zero-sum game.
- 085 3. We characterize an **approximate local equilibrium solution for federated zero-sum game**. We  
086 show that the limit points of a federated gradient descent ascent (FEDGDA) algorithm include the  
087 equilibria of the zero-sum game.
- 088 4. We show that an equilibrium solution of the federated game obtained at the server consistently  
089 estimates the **moment conditions of every client**. An important insight derived from our results is  
090 that the consistency of the GMM-estimators on clients directly depend on the heterogeneity bias.
- 091 5. We experimentally validate that even for non-i.i.d. data, FEDDEEPGMM has convergent dynamics  
092 analogous to the centralized DEEPGMM algorithm.

093 This work focuses on the existence results of federated equilibrium solutions, federated consistent  
094 GMM estimators, and thereby structurally solving the federated IV analysis problem. The existence  
095 of approximate client-local equilibria via federated solution has applications beyond the GMM and  
096 IV analysis, to problems such as federated generative adversarial networks (FedGAN) (Rasouli et al.,  
097 2020), where a Nash Equilibrium may not exist (Farnia & Ozdaglar, 2020). However, the scope of  
098 our discussion does not include FedGAN, or a new federated minimax algorithm, or, for that matter,  
099 the convergence theory and scalability. We leave an open problem to characterize and recover a  
100 federated mixed-strategy Nash equilibrium, which has enormous applications to diverse domains  
(Barron, 2024). We compare and contrast our method against related works in Appendix D.

## 102 2 PRELIMINARIES AND MODEL

103  
104 [In this section, we introduce our basic terminologies tailored to a motivating application as described  
105 by Dayan et al. \(2021\) in the context of the global-scale federated learning.](#)  
106

107 <sup>1</sup>Wu et al. (2023) used FEDGMM as an acronym for federated Gaussian mixture models.

**Client-local causal inference.** We begin by adapting (Bennett et al., 2019) for a client-local setting. Consider a distributed system as a set of  $N$  clients  $[N]$  with datasets  $S^i = \{(x_j^i, y_j^i)\}_{j=1}^{n_i}, \forall i \in [N]$ . We assume that for a client  $i \in [N]$ , the treatment and outcome variables  $x_j^i$  and  $y_j^i$ , respectively, are related by the process  $Y^i = g_0^i(X^i) + \epsilon^i, i \in [N]$ .

Referring to (Dayan et al., 2021),  $x_j^i$  and  $y_j^i$  represent the clinical features (CBC, D-dimer, oxygen devices, etc.) and the outcome (death, duration of survival over oxygen administration, etc.), respectively, on a client  $i \in [N]$ .

We assume that each client-local residual  $\epsilon^i$  has zero mean and finite variance, i.e.  $\mathbb{E}[\epsilon^i] = 0, \mathbb{E}[(\epsilon^i)^2] < \infty$ . Furthermore, we assume that the treatment variables  $X^i$  are endogenous on the clients, i.e.  $\mathbb{E}[\epsilon^i|X^i] \neq 0$ , and therefore,  $g_0^i(X^i) \neq \mathbb{E}[Y^i|X^i]$ . We assume that the treatment variables are influenced by instrumental variables  $Z^i, \forall i \in [N]$  so that

$$P(X^i|Z^i) \neq P(X^i). \quad (1)$$

Furthermore, the instrumental variables do not directly influence the outcome variables  $Y^i, \forall i \in [N]$ :

$$\mathbb{E}[\epsilon^i|Z^i] = 0. \quad (2)$$

Referring to (Dayan et al., 2021),  $Z^i$  would represent the HIF levels in the patients on the client  $i \in [N]$ , which their implementation misses to incorporate.

**Federated causal inference function.** Note that, assumptions 1, 2 are local to the clients, thus, honour the data-privacy requirements of a federated learning task. In this setting, we aim to discover a common or global causal response function that would *fit the data generation processes of each client without centralizing the data*. More specifically, we learn a parametric function  $g_0(\cdot) \in G := \{g(\cdot, \theta) | \theta \in \Theta\}$  expressed as  $g_0 := g(\cdot, \theta_0)$  for  $\theta_0 \in \Theta$ , defined by

$$g(\cdot, \theta_0) = \frac{1}{N} \sum_{i=1}^N g^i(\cdot, \theta_0). \quad (3)$$

**The generalized method of moments (GMM)** estimates the parameters of the causal response function (3) using a certain number of *moment conditions*. Define the *moment function* on a client  $i \in [N]$  as a vector-valued function  $f^i : \mathbb{R}^{|Z^i|} \rightarrow \mathbb{R}^m$  with components  $f_1^i, f_2^i, \dots, f_m^i$ . Based on equation (2), we define the moment conditions using parametrized functions  $\{f_j^i\}_{j=1}^m, \forall i \in [N]$  as

$$\mathbb{E}[f_j^i(Z^i)\epsilon^i] = 0, \forall j \in [m], \forall i \in [N], \quad (4)$$

We assume that  $m$  moment conditions  $\{f_j^i\}_{j=1}^m$  at each client  $i \in [N]$  are sufficient to identify a unique federated estimate  $\hat{\theta}$  to  $\theta_0$ . With (4), we define the moment conditions on a client  $i \in [N]$  as

$$\psi(f_j^i; \theta) = 0, \forall j \in [m], \quad (5)$$

where  $\psi(f^i; \theta) := \mathbb{E}[f^i(Z^i)\epsilon^i] = \mathbb{E}[f^i(Z^i)(Y^i - g^i(X^i; \theta))]$ . In empirical terms, the sample moments for the  $i$ -th client with  $n_i$  samples are given by

$$\psi_{n_i}(f^i; \theta) := \mathbb{E}_{n_i}[f^i(Z)\epsilon^i] = \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i)(Y_k^i - g^i(X_k^i; \theta)), \quad (6)$$

where  $\psi_{n_i}(f^i; \theta) = (\psi_{n_i}(f_1^i; \theta), \psi_{n_i}(f_2^i; \theta), \dots, \psi_{n_i}(f_m^i; \theta))$  is the moment condition vector, and  $\psi_{n_i}(f_j^i; \theta) = \frac{1}{n_i} \sum_{k=1}^{n_i} f_j^i(Z_k^i)(Y_k^i - g^i(X_k^i; \theta))$ . Thus, for empirical estimation of the causal response function  $g_0^i$  at client  $i \in [N]$ , it needs to satisfy

$$\psi_{n_i}(f_j^i; \theta_0) = 0, \forall i \in [N] \text{ and } j \in [m] \text{ at } \theta = \theta_0. \quad (7)$$

**The optimization problem.** Equation (7) is reformulated as an optimization problem given by

$$\min_{\theta \in \Theta} \|\psi_{n_i}(f_1^i; \theta), \psi_{n_i}(f_2^i; \theta), \dots, \psi_{n_i}(f_m^i; \theta)\|^2, \quad (8)$$

where we use the Euclidean norm  $\|w\|^2 = w^T w$ . Drawing inspiration from Hansen (1982), DEEP-GMM used a weighted norm, which yields minimal asymptotic variance for a consistent estimator

162  $\tilde{\theta}$ , to cater to the cases of (finitely) large number of moment conditions. We adapt their weighted  
 163 norm  $\|w\|_{\tilde{\theta}}^2 = w^T \mathcal{C}_{\tilde{\theta}}^{-1} w$ , to a client-local setting via the [positive semi-definite](#) covariance matrix  $\mathcal{C}_{\tilde{\theta}}$   
 164 defined by

$$165 \quad [\mathcal{C}_{\tilde{\theta}}]_{jl} = \frac{1}{n_i} \sum_{k=1}^{n_i} f_j^i(Z_k^i) f_l^i(Z_k^i) (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2. \quad (9)$$

166 Now considering the vector space  $\mathcal{V}$  of real-valued functions  $f_i$  of  $Z$ ,  $\psi_{n_i}(f^i; \theta) =$   
 167  $(\psi_{n_i}(f_1^i; \theta), \psi_{n_i}(f_2^i; \theta), \dots, \psi_{n_i}(f_m^i; \theta))$  is a linear operator on  $\mathcal{V}$  and

$$171 \quad \mathcal{C}_{\tilde{\theta}}(f^i, h^i) = \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i) h^i(Z_k^i) (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2 \quad (10)$$

172 is a bilinear form. With that, for any subset  $\mathcal{F}^i \subset \mathcal{V}$ , we define a function

$$175 \quad \Psi_{n_i}(\theta, \mathcal{F}^i, \tilde{\theta}) = \sup_{f^i \in \mathcal{F}^i} \psi_{n_i}(f^i; \theta) - \frac{1}{4} \mathcal{C}_{\tilde{\theta}}(f^i, f^i),$$

176 which leads to the following client-local optimization problem:

$$177 \quad \theta^{\text{GMM}} \in \arg \min_{\theta \in \Theta} \Psi_{n_i}(\theta, \mathcal{F}^i, \tilde{\theta}), \quad (11)$$

178 where  $\mathcal{F}^i = \text{span}(\{f_j^i\}_{j=1}^m)$ ,  $\Psi_{n_i}(\theta, \mathcal{F}^i, \tilde{\theta}) = \|\psi_{n_i}(f_1^i; \theta), \psi_{n_i}(f_2^i; \theta), \dots, \psi_{n_i}(f_m^i; \theta)\|_{\tilde{\theta}}^2$ , and the  
 179 the weighted norm  $\|\cdot\|_{\tilde{\theta}}$  defined by equation (9).

180 **The zero-sum game for deep generalized method of moments.** As the data-dimension grows,  
 181 the function class  $\mathcal{F}^i$  is replaced with a class of neural networks of a certain architecture, i.e.  
 182  $\mathcal{F}^i = \{f^i(z, \tau) : \tau \in \mathcal{T}\}$  with varying weights  $\tau$ . Similarly, let  $\mathcal{G}^i = \{g^i(x, \theta) : \theta \in \Theta\}$  be another  
 183 class of neural networks with varying weights  $\theta$ . With that, define

$$184 \quad U_{\tilde{\theta}}^i(\theta, \tau) := \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i, \tau) (Y_k^i - g^i(X_k^i; \theta)) - \frac{1}{4n_i} \sum_{k=1}^{n_i} (f^i(Z_k^i, \tau))^2 (Y_k^i - g^i(X_k^i; \theta))^2 \quad (12)$$

185 Then for a client  $i$ , (11) is reformulated as the following

$$186 \quad \theta^{\text{DGMM}} \in \arg \min_{\theta \in \Theta} \sup_{\tau \in \mathcal{T}} U_{\tilde{\theta}}^i(\theta, \tau). \quad (13)$$

187 Equation (13) forms a zero-sum game, whose equilibrium solution is shown to be a true estimator to  
 188  $\theta_0$  under a set of standard assumptions; see Theorem 2 in (Bennett et al., 2019).

### 189 3 FEDERATED DEEP GMM VIA FEDERATED EQUILIBRIUM SOLUTIONS

#### 190 3.1 FEDERATED DEEP GENERALIZED METHOD MOMENT (FEDDEEPGMM)

191 We need to find the global moment estimators for the causal response function to fit data on each  
 192 client. Thus, the federated counterpart of equation (5) is given by

$$193 \quad \psi(f; \theta) := \mathbb{E}_i[\mathbb{E}[f^i(Z^i)(Y_k^i - g^i(X_k^i; \theta))] = 0, \quad (14)$$

194 where the expectation  $\mathbb{E}_i$  is over the clients. In this work, we consider *full client participation*. Thus,  
 195 for the empirical federated moment estimation, we formulate:

$$196 \quad \psi_n(f; \theta) := \frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f^i; \theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i) (Y_k^i - g^i(X_k^i; \theta)) \quad (15)$$

197 With that, the federated problem for GMM following (11) is formulated as:

$$198 \quad \theta^{\text{FedDeepGMM}} \in \arg \min_{\theta \in \Theta} \|\psi_n(f; \theta)\|_{\tilde{\theta}}^2, \quad (16)$$

199 where  $\|w\|_{\tilde{\theta}} = w^T \mathcal{C}_{\tilde{\theta}}^{-1} w$  is the previously defined weighted-norm with inverse covariance as weights.  
 200 We propose FEDDEEPGMM, a “deep” reformulation of the federated optimization problem based  
 201 on the neural networks of a given architecture shared among clients and is shown to have the same  
 202 solution as the federated GMM problem formulated earlier.

**Lemma 1.** Let  $\mathcal{F} = \text{span}\{f_j^i \mid i \in [N], j \in [m]\}$ . An equivalent objective function for the federated moment estimation optimization problem (16) is given by:

$$\|\psi_N(f; \theta)\|_{\bar{\theta}}^2 = \sup_{\substack{f^i \in \mathcal{F} \\ \forall i \in [N]}} \frac{1}{N} \sum_{i=1}^N \left( \psi_{n_i}(f^i; \theta) - \frac{1}{4} \mathcal{C}_{\bar{\theta}}(f^i; f^i) \right), \text{ where}$$

$$\psi_{n_i}(f^i; \theta) := \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i) (Y_k^i - g^i(X_k^i; \theta)), \text{ and } \mathcal{C}_{\bar{\theta}}(f^i; f^i) := \frac{1}{n_i} \sum_{k=1}^{n_i} (f^i(Z_k^i))^2 (Y_k^i - g^i(X_k^i; \bar{\theta}))^2.$$

The proof of Lemma 1 is given in Appendix B.1. The federated zero-sum game is then defined by:

$$\hat{\theta}^{\text{FedDeepGMM}} \in \arg \min_{\theta \in \Theta} \sup_{\tau \in \mathcal{T}} U_{\bar{\theta}}(\theta, \tau) := \frac{1}{N} \sum_{i=1}^N U_{\bar{\theta}}^i(\theta, \tau), \quad (17)$$

where  $U_{\bar{\theta}}^i(\theta, \tau)$  is defined in equation (12). The federated DEEPGMM formulation as a zero-sum game defined by a federated minimax optimization problem (17) provides a framework to recover the global estimator as a federated equilibrium solution.

Referring to (Dayan et al., 2021),  $\hat{\theta}^{\text{FedDeepGMM}}$  is the federated GMM estimator that consistently estimates the moment conditions of the clients under an approximation error as described later in Definition 3. These moment conditions are then employed on each client to analyse the impact of the instrumental variable  $Z^i$ .

### 3.2 FEDERATED SEQUENTIAL GAMES AND THEIR EQUILIBRIUM SOLUTIONS

As minimax is not equal to maximin in general for a non-convex-non-concave problem, it is important to model the federated game as a sequential game (Jin et al., 2020) whose outcome would depend on what move – maximization or minimization – is taken first. We start with the following assumptions:

**Assumption 1.** Client-local objective  $U_{\bar{\theta}}^i(\theta, \tau) \forall i \in [N]$  is twice continuously differentiable for both  $\theta$  and  $\tau$ . Thus, the global objective  $U_{\bar{\theta}}(\theta, \tau)$  is also a twice continuously differentiable function.

**Assumption 2 (Smoothness).** The gradient of each client’s local objective,  $\nabla U_{\bar{\theta}}^i(\theta, \tau)$ , is Lipschitz continuous with respect to both  $\theta$  and  $\tau$ . For all  $i \in [N]$ , there exist constants  $L > 0$  such that:

$$\begin{aligned} \|\nabla_{\theta} U_{\bar{\theta}}^i(\theta_1, \tau_1) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta_2, \tau_2)\| &\leq L \|(\theta_1, \tau_1) - (\theta_2, \tau_2)\|, \text{ and} \\ \|\nabla_{\tau} U_{\bar{\theta}}^i(\theta_1, \tau_1) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta_2, \tau_2)\| &\leq L \|(\theta_1, \tau_1) - (\theta_2, \tau_2)\|, \end{aligned}$$

$\forall (\theta_1, \tau_1), (\theta_2, \tau_2)$ . Thus,  $U_{\bar{\theta}}(\theta, \tau)$  is  $L$ -Lipschitz smooth.

**Assumption 3 (Bounded Gradient Dissimilarity).** The heterogeneity of the local gradients with respect to (w.r.t.)  $\theta$  and  $\tau$  is bounded as follows:

$$\|\nabla_{\theta} U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\theta} U_{\bar{\theta}}(\theta, \tau)\| \leq \zeta_{\theta}^i \quad \|\nabla_{\tau} U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\tau} U_{\bar{\theta}}(\theta, \tau)\| \leq \zeta_{\tau}^i,$$

where  $\zeta_{\theta}^i, \zeta_{\tau}^i \geq 0$  are the bounds that quantify the degree of gradient dissimilarity at client  $i \in [N]$ .

**Assumption 4 (Bounded Hessian Dissimilarity).** The heterogeneity in terms of hessian w.r.t.  $\theta$  and  $\tau$  is bounded as follows:

$$\begin{aligned} \|\nabla_{\theta\theta}^2 U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\theta\theta}^2 U_{\bar{\theta}}(\theta, \tau)\|_{\sigma} &\leq \rho_{\theta}^i, & \|\nabla_{\tau\tau}^2 U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)\|_{\sigma} &\leq \rho_{\tau}^i, \\ \|\nabla_{\theta\tau}^2 U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\theta\tau}^2 U_{\bar{\theta}}(\theta, \tau)\|_{\sigma} &\leq \rho_{\theta\tau}^i, & \|\nabla_{\tau\theta}^2 U_{\bar{\theta}}^i(\theta, \tau) - \nabla_{\tau\theta}^2 U_{\bar{\theta}}(\theta, \tau)\|_{\sigma} &\leq \rho_{\tau\theta}^i, \end{aligned}$$

where  $\rho_{\theta}^i, \rho_{\tau}^i, \rho_{\theta\tau}^i$ , and  $\rho_{\tau\theta}^i \geq 0$  quantify the degree of hessian dissimilarity at client  $i \in [N]$  by spectral norm  $\|\cdot\|_{\sigma}$ .

Assumptions 3 and 4 provide a measure of data heterogeneity across clients in a federated setting. In the special case, when  $\zeta$  and  $\rho$ ’s are all 0, then the data is homogeneous across clients.

We adopt the Stackelberg equilibrium for pure strategies (Jin et al., 2020) to characterize the solution of the minimax federated optimization problem for a non-convex non-concave function  $U_{\bar{\theta}}(\theta, \tau)$  for the sequential game where min-player goes first and the max-player goes second. To avoid ambiguity between the adjectives of the terms global/local objective functions in federated learning and the global/local nature of minimax points in optimization, we refer to a global objective as the federated objective and a local objective as the client’s objective.

**Definition 1 (Local minimax point).** [Definition 14 of (Jin et al., 2020)] Let  $U(\theta, \tau)$  be a function defined over  $\Theta \times \mathcal{T}$  and let  $h$  be a function satisfying  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ . There exists a  $\delta_0$ , such that for any  $\delta \in (0, \delta_0]$ , and any  $(\theta, \tau)$  such that  $\|\theta - \hat{\theta}\| \leq \delta$  and  $\|\tau - \hat{\tau}\| \leq \delta$ , then a point  $(\hat{\theta}, \hat{\tau})$  is a local minimax point of  $U$ , if  $\forall (\theta, \tau) \in \Theta \times \mathcal{T}$ , it satisfies:

$$U_{\hat{\theta}}(\hat{\theta}, \tau) \leq U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\hat{\theta}}(\theta, \tau'). \quad (18)$$

With that, the first-order & second-order necessary conditions for local minimax points are as below.

**Lemma 2** (Propositions 18, 19, 20 of (Jin et al., 2020)). *Under assumption 1, any local minimax point satisfies the following conditions:*

- **First-order Necessary Condition:** A local minimax point  $(\theta, \tau)$  satisfies:  $\nabla_{\theta} U_{\hat{\theta}}(\theta, \tau) = 0$  and  $\nabla_{\tau} U_{\hat{\theta}}(\theta, \tau) = 0$ .
- **Second-order Necessary Condition:** A local minimax point  $(\theta, \tau)$  satisfies:  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau) \preceq \mathbf{0}$ .  
Moreover, if  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau) \prec 0$ , then  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}} - \nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right] (\theta, \tau) \succeq 0$ .
- **Second-order Sufficient Condition:** A stationary point  $(\theta, \tau)$  that satisfies  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau) \prec \mathbf{0}$ , and  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}} - \nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right] (\theta, \tau) \succ 0$  guarantees that  $(\theta, \tau)$  is a strict local minimax.

Now, in order to define the federated approximate equilibrium solutions, we first define an approximate local minimax point.

**Definition 2 (Approximate Local minimax point).** [An adaptation of definition 34 of (Jin et al., 2020)] Let  $U(\theta, \tau)$  be a function defined over  $\Theta \times \mathcal{T}$  and let  $h$  be a function satisfying  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ . There exists a  $\delta_0$ , such that for any  $\delta \in (0, \delta_0]$ , and any  $(\theta, \tau)$  such that  $\|\theta - \hat{\theta}\| \leq \delta$  and  $\|\tau - \hat{\tau}\| \leq \delta$ , then a point  $(\hat{\theta}, \hat{\tau})$  is an  $\varepsilon$ -approximate local minimax point of  $U$ , if it satisfies:

$$U_{\hat{\theta}}(\hat{\theta}, \tau) - \varepsilon \leq U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\hat{\theta}}(\theta, \tau') + \varepsilon, \quad (19)$$

We aim to achieve approximate local minimax points for every client as a solution of the federated minimax optimization. With that, we characterize the federated solution as the following.

**Definition 3 ( $\mathcal{E}$ -Approximate Federated Equilibrium Solutions).** Let  $\mathcal{E} = \{\varepsilon^i\}_{i=1}^N$  be the approximation error vector for clients  $i \in [N]$ . Let  $U_{\hat{\theta}}^i(\theta, \tau)$  be a function defined over  $\Theta \times \mathcal{T}$  for a client  $i \in [N]$  and  $U_{\hat{\theta}}(\theta, \tau) := \frac{1}{N} \sum_{i=1}^N U_{\hat{\theta}}^i(\theta, \tau)$ . An  $\mathcal{E}$ -approximate federated equilibrium point  $(\hat{\theta}, \hat{\tau})$  (that is an  $\varepsilon^i$ -approximate local minimax point for each client's objective  $U_{\hat{\theta}}^i$ ), must follow the conditions below:

1.  **$\varepsilon^i$ -First-order Necessary Condition:** The point  $(\hat{\theta}, \hat{\tau})$  must be an  $\varepsilon^i$  stationary point for every client  $i \in [N]$ , i.e.,  $\|\nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})\| \leq \varepsilon^i$ , and  $\|\nabla_{\tau} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})\| \leq \varepsilon^i$ .
2. **Second-Order  $\varepsilon^i$  Necessary Condition:** The point  $(\hat{\theta}, \hat{\tau})$  must satisfy the second-order conditions:  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \preceq -\varepsilon^i I$ , and  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right] (\hat{\theta}, \hat{\tau}) \succeq \varepsilon^i I$ .
3. **Second-Order  $\varepsilon^i$  Sufficient Condition:** An  $\varepsilon^i$  stationary point  $(\theta, \tau)$  that satisfies  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \prec -\varepsilon^i I$ , and  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right] (\hat{\theta}, \hat{\tau}) \succ \varepsilon^i I$  guarantees that  $(\hat{\theta}, \hat{\tau})$  is a strict local minimax point  $\forall i \in [N]$  that satisfies  $\varepsilon^i$  approximate equilibrium as in definition 2.

We now state the main theoretical result of our work in this theorem.

**Theorem 1.** Under assumptions 1, 2, 3 and 4, a minimax solution  $(\hat{\theta}, \hat{\tau})$  of federated optimization problem (17) that satisfies the equilibrium condition as in definition 1:  $U_{\hat{\theta}}(\hat{\theta}, \tau) \leq U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\hat{\theta}}(\theta, \tau')$ , is an  $\mathcal{E}$ -approximate federated equilibrium solution as defined in 3, where the approximation error  $\varepsilon^i$  for each client  $i \in [N]$  lies in:  $\max\{\zeta_{\theta}^i, \zeta_{\tau}^i\} \leq \varepsilon^i \leq \min\{\alpha - \rho_{\tau}^i, \beta - B^i\}$  for  $\rho_{\tau}^i < \alpha$  and  $B^i > \beta$ , such that  $\alpha := \left| \lambda_{\max} \left( \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \right) \right|$ ,  $\beta := \lambda_{\min} \left( \left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}} - \nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right] (\hat{\theta}, \hat{\tau}) \right)$  and  $B^i := \rho_{\theta}^i + L\rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L\rho_{\tau\theta}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L^2\rho_{\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i) \cdot \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|}$ .

The proof of Theorem 1 is given in Appendix B.2. Note that when data is homogeneous (i.e., for each client  $i$ ,  $\zeta_\theta^i, \zeta_\tau^i, \rho_\tau^i$  and  $B^i$  are all zeroes), each client satisfies an exact local minimax equilibrium.

**Remark 1.** In Theorem 1, note that if the interval  $[\max\{\zeta_\theta^i, \zeta_\tau^i\}, \min\{\alpha - \rho_\tau^i, \beta - B^i\}]$  is empty, i.e.  $\max\{\zeta_\theta^i, \zeta_\tau^i\} > \min\{\alpha - \rho_\tau^i, \beta - B^i\}$ , then no such  $\varepsilon^i$  exists and  $(\hat{\theta}, \hat{\tau})$  fails to be a local  $\varepsilon^i$  approximate equilibrium point for that clients. It may happen in two cases:

1. The gradient dissimilarity  $\zeta_\theta^i, \zeta_\tau^i$  is too large, indicating high heterogeneity, then  $(\hat{\theta}, \hat{\tau})$ - the solution to the federated objective would fail to become an approximate equilibrium point for the clients. It is a practical consideration for a federated convergence facing difficulty against high heterogeneity.
2. If  $\alpha \approx \rho_\tau^i$  or  $\beta \approx B^i$ , this indicates that the client's local curvature structure significantly differs from the global curvature. In this case, the client's objective may be flatter or even oppositely curved compared to the global model, reflecting high heterogeneity.

Now we state the result on the per-client consistency of the FEDGMM estimator.

**Theorem 2 (Consistency).** [Adaptation of Theorem 2 of (Bennett et al., 2019)] Let  $\tilde{\theta}_n$  be a data-dependent choice for the federated objective that has a limit in probability. Let  $h$  be a function satisfying  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ . For each client  $i \in [N]$ , define  $m^i(\theta, \tau, \tilde{\theta}) := f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4}f^i(Z^i; \tau)^2(Y^i - g(X^i; \tilde{\theta}))^2$ ,  $M^i(\theta) = \sup_{\tau \in \mathcal{T}} \mathbb{E}[m^i(\theta, \tau, \tilde{\theta})]$  and  $\eta^i(\epsilon) := \inf_{d(\theta, \theta_0) \geq \epsilon} M^i(\theta) - M^i(\theta_0)$  for every  $\epsilon > 0$ . Fix some  $\delta_0$ , for any  $\delta \in (0, \delta_0]$  and any  $(\theta, \tau)$  such that  $\|\theta - \hat{\theta}\| \leq \delta$  and  $\|\tau - \hat{\tau}\| \leq \delta$ , let  $(\hat{\theta}_n, \hat{\tau}_n)$  be a solution that satisfies the approximate equilibrium for each of the client  $i \in [N]$  as

$$\sup_{\tau \in \mathcal{T}} U_{\hat{\theta}}^i(\hat{\theta}_n, \tau) - \varepsilon^i - o_p(1) \leq U_{\hat{\theta}}^i(\hat{\theta}_n, \hat{\tau}_n) \leq \inf_{\theta \in \Theta} \max_{\tau: \|\tau - \hat{\tau}_n\| \leq h(\delta)} U_{\hat{\theta}}^i(\theta, \tau) + \varepsilon^i + o_p(1).$$

Then, under similar assumptions as in Assumptions 1 to 5 of (Bennett et al., 2019), the global solution  $\hat{\theta}_n$  is a consistent estimator to the true parameter  $\theta_0$ , i.e.  $\hat{\theta}_n \xrightarrow{P} \theta_0$  when the approximate error  $\varepsilon^i < \frac{\eta^i(\epsilon)}{2}$  for every  $\epsilon > 0$  for each client  $i \in [N]$ .

The assumptions and the proof of Theorem 2 are included in Appendix B.3.

**Remark 2.** Theorem 2 formalizes a tradeoff between data heterogeneity and the consistency of the global estimator in federated learning for each client. If the approximation error  $\varepsilon^i$  is large for a client  $i \in [N]$ , then the solution  $\hat{\theta}_n$  may fail to consistently estimate the true parameter of client  $i$ . In contrast, when data across clients have similar distribution (i.e., case for low heterogeneity), the federated optimal model  $\hat{\theta}_n$  is consistent across clients.

### 3.3 FEDERATED GRADIENT DESCENT ASCENT ALGORITHM AND IT'S LIMIT POINTS

Bennett et al. (2019) used Optimistic Adam (OADAM), a variant of Adam (Kingma, 2015) based stochastic gradient descent ascent (SGDA) algorithm (Daskalakis et al., 2018). However, it is known that a well-tuned SGD outperforms Adam in overparametrized settings (Wilson et al., 2017). As our experiments show in Section (4), that gradient descent ascent updates are competitive to OADAM for minimax optimization in centralized setting. Considering this, we employ an adaptation of the standard gradient descent ascent algorithm to federated (FEDGDA) setting.

FEDGDA is well-explored in the literature: (Deng & Mahdavi, 2021; Sharma et al., 2022; Shen et al., 2024; Wu et al., 2024). The clients run the gradient descent ascent algorithm for several local updates and then the orchestrating server synchronizes them by collecting the model states, averaging them, and broadcasting it to the clients. A detailed description is included as a pseudocode in Appendix A.

Similar to (Bennett et al., 2019), we note that the federated minimax optimization problem (17) is not convex-concave on  $(\theta, \tau)$ . The convergence results of variants of FEDGDA (Sharma et al., 2022; Shen et al., 2024; Wu et al., 2024) assume that  $U_{\hat{\theta}}(\theta, \tau)$  is non-convex on  $\theta$  and satisfies a  $\mu$ -Polyak Łojasiewicz (PL) inequality on  $\tau$ , see assumption 4 in (Sharma et al., 2022). PL condition is known to be satisfied by over-parametrized neural networks (Charles & Papailiopoulos, 2018; Liu et al., 2022). The convergence results of FEDGDA will follow (Sharma et al., 2022). We include a formal statement in Appendix A. However, beyond convergence, we primarily aim to show that an optimal solution will consistently estimate the moment conditions of the clients, which we do next.

For Algorithm 1 in Appendix A, let  $\alpha_1 = \frac{\eta}{\gamma}$ ,  $\alpha_2 = \eta$  be the learning rates for gradient updates to  $\theta$  and  $\tau$ , respectively. Without loss of generality the FEDGDA updates are:

$$\theta_{t+1} = \theta_t - \eta \frac{1}{\gamma} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \text{ and } \tau_{t+1} = \tau_t + \eta \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i)$$

We call it  $\gamma$ -FEDGDA, where  $\gamma$  is the ratio of  $\alpha_1$  to  $\alpha_2$ . As  $\eta \rightarrow 0$  corresponds to FEDGDA-flow, under the smoothness of  $U_{\bar{\theta}}^i$ , bounded gradient heterogeneity (assumption 3) and for fixed local rounds  $R$ , FEDGDA-flow becomes:

$$\frac{d\theta}{dt} = -\frac{1}{\gamma} R \nabla_{\theta} U_{\bar{\theta}}(\theta, \tau) + \mathcal{O}\left(\frac{R}{\gamma} \zeta_{\theta}\right), \text{ and } \frac{d\tau}{dt} = R \nabla_{\tau} U_{\bar{\theta}}(\theta, \tau) + \mathcal{O}(R \zeta_{\tau}).$$

We further elaborate on FEDGDA-flow in Appendix C.1. We aim to find out the relationship between stable equilibrium and local minimax points of the federated optimization problem. For that, we now define a strictly linearly stable equilibrium of the  $\gamma$ -FEDGDA flow.

**Proposition 1.** *Given the Jacobian matrix for  $\gamma$ -FEDGDA flow as  $\mathbf{J} = \begin{pmatrix} -\frac{1}{\gamma} R \nabla_{\theta\theta}^2 U_{\bar{\theta}}(\theta, \tau) & -\frac{1}{\gamma} R \nabla_{\theta\tau}^2 U_{\bar{\theta}}(\theta, \tau) \\ R \nabla_{\tau\theta}^2 U_{\bar{\theta}}(\theta, \tau) & R \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) \end{pmatrix}$ , a point  $(\theta, \tau)$  is a strictly linearly stable equilibrium of the  $\gamma$ -FEDGDA flow if and only if the real parts of all eigenvalues of  $\mathbf{J}$  are negative, i.e.,  $\text{Re}(\Lambda_j) < 0$  for all  $j$ .*

The proof follows a strategy similar to (Jin et al., 2020).

Let  $\gamma$ -FGDA be the set of strictly linearly stable points of the  $\gamma$ -FEDGDA flow, and  $\text{LocMinimax}$  be the set of local minimax points of the federated zero-sum game. Define

$$\begin{aligned} \overline{\infty - \text{FGDA}} &:= \limsup_{\gamma \rightarrow \infty} \gamma - \text{FGDA} := \bigcap_{\gamma_0 > 0} \bigcup_{\gamma > \gamma_0} \gamma - \text{FGDA}, \text{ and} \\ \underline{\infty - \text{FGDA}} &:= \liminf_{\gamma \rightarrow \infty} \gamma - \text{FGDA} := \bigcup_{\gamma_0 > 0} \bigcap_{\gamma > \gamma_0} \gamma - \text{FGDA}. \end{aligned}$$

We now state the theorem that establishes that the stable limit points of  $\infty$ -FGDA are the local minimax points, up to some degenerate cases.

**Theorem 3.** *Under Assumption 1,  $\text{LocMinimax} \subset \underline{\infty - \text{FGDA}} \subset \overline{\infty - \text{FGDA}} \subset \text{LocMinimax} \cup \mathcal{A}$ , where  $\mathcal{A} := \{(\theta, \tau) | (\theta, \tau) \text{ is stationary and } \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) \text{ is degenerate}\}$ . Moreover, if the hessian  $\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)$  is smooth, then  $\mathcal{A}$  has measure zero in  $\Theta \times \mathcal{T} \subset \mathbb{R}^d \times \mathbb{R}^k$ .*

Essentially, Theorem 3 states that the limit points of FEDGDA are the local minimax solutions, and thereby the equilibrium solution of the federated zero-sum game, up to some degenerate case. The proof of Theorem 3 is included in Appendix C.2. Theorems 1, 2, and 3 together complete the theoretical foundation of the pipeline in our work.

## 4 EXPERIMENTS

We extend the experimental evaluations of DEEPGMM (Bennett et al., 2019) to a federated setting. We further discuss this benchmark structure in Appendix E. More specifically, we evaluate the ability of FEDDEEPGMM to fit low- and high- dimensional data to demonstrate its convergence. Similar to DEEPGMM, we assess two scenarios in regards to  $((X, Y), Z)$ :

- (a) **The instrumental and treatment variables  $Z$  and  $X$  are both low-dimensional.** In this case, we use 1-dimensional synthetic datasets corresponding to the following functions: (a) **Absolute:**  $g_0(x) = |x|$ , (b) **Step:**  $g_0(x) = 1_{\{x \geq 0\}}$ , (c) **Linear:**  $g_0(x) = x$ . To generate the synthetic data, similar to (Bennett et al., 2019; Lewis & Syrgkanis, 2018) we apply the following generation process:

$$Y = g_0(X) + e + \delta \quad \text{and } X = Z^{(1)} + Z^{(2)} + e + \gamma \quad (20)$$

$$(Z^{(1)}, Z^{(2)}) \sim \text{Uniform}([-3, 3]^2) \quad \text{and } e \sim \mathcal{N}(0, 1), \quad \gamma, \delta \sim \mathcal{N}(0, 0.1) \quad (21)$$



(b)  $Z$  and  $X$  are low-dimensional or high-dimensional or both. First,  $Z$  and  $X$  are generated as in (20,21). Then for high-dimensional data, we map  $Z$  and  $X$  to an image using the mapping:

$$\text{Image}(x) = \text{Dataset}(\text{round}(\min(\max(1.5x + 5, 0), 9))),$$

where  $\text{round}(\min(\max(1.5x + 5, 0), 9))$  returns an integer between 0 and 9. Essentially, the function  $\text{Dataset}(\cdot)$  randomly selects an image following its index. We use datasets FEMNIST (Federated Extended MNIST) and CIFAR10 (Caldas et al., 2018) for images of size  $28 \times 28$  and  $3 \times 32 \times 32$ , respectively. Thus, we have the following cases: (a)  $\text{Dataset}_z$ :  $X = X^{\text{low}}, Z = \text{Image}(Z^{\text{low}})$ , (b)  $\text{Dataset}_x$ :  $Z = Z^{\text{low}}, X = \text{Image}(X^{\text{low}})$ , and (c)  $\text{Dataset}_{x,z}$ :  $Z = \text{Image}(Z^{\text{low}}), X = \text{Image}(X^{\text{low}})$ , where  $\text{Dataset}$  takes values FEMNIST and CIFAR10.

We implemented and benchmarked FEDGDA and FEDSGDA to solve the FEDDEEPGMM problem. For reference, we implemented OADAM, GDA, and SGDA to solve the DEEPGMM in centralized setting. For high-dimensional scenarios, we implement a CNN architecture to process images, while for low-dimensional scenarios, we use a multilayer perceptron (MLP). Code is available at <https://anonymous.4open.science/r/FederatedDeepGMM-417C>.

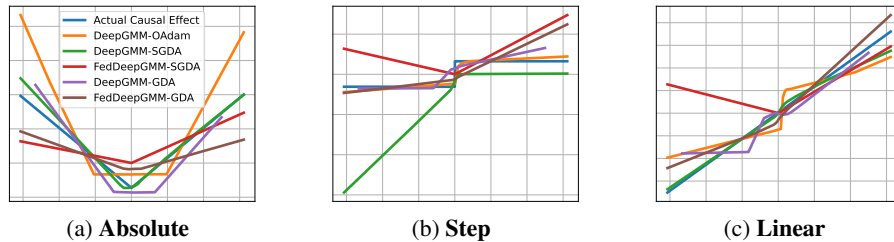


Figure 1: Estimated  $\hat{g}$  compared to true  $g$  in low-dimensional scenarios

Estimations	DEEPGMM-OAdam	DEEPGMM-GDA	FDEEPGMM-GDA	DEEPGMM-SGDA	FDEEPGMM-SGDA
<b>Absolute</b>	$0.03 \pm 0.01$	$0.013 \pm .01$	$0.4 \pm 0.01$	$0.009 \pm 0.01$	$0.2 \pm 0.00$
<b>Step</b>	$0.3 \pm 0.00$	$0.03 \pm 0.00$	$0.04 \pm 0.01$	$0.112 \pm 0.00$	$0.23 \pm 0.01$
<b>Linear</b>	$0.01 \pm 0.00$	$0.02 \pm 0.00$	$0.01 \pm 0.00$	$0.03 \pm 0.00$	$0.04 \pm 0.00$
<b>FEMNIST<sub>x</sub></b>	$0.50 \pm 0.00$	$1.11 \pm 0.01$	$0.21 \pm 0.02$	$0.40 \pm 0.01$	$0.19 \pm 0.01$
<b>FEMNIST<sub>x,z</sub></b>	$0.24 \pm 0.00$	$0.46 \pm 0.09$	$0.19 \pm 0.03$	$0.14 \pm 0.02$	$0.20 \pm 0.00$
<b>FEMNIST<sub>z</sub></b>	$0.10 \pm 0.00$	$0.42 \pm 0.01$	$0.24 \pm 0.01$	$0.11 \pm 0.02$	$0.23 \pm 0.01$
<b>CIFAR10<sub>x</sub></b>	$0.55 \pm 0.30$	$0.19 \pm 0.01$	$0.25 \pm 0.03$	$0.20 \pm 0.08$	$0.22 \pm 0.08$
<b>CIFAR10<sub>x,z</sub></b>	$0.40 \pm 0.11$	$0.24 \pm 0.00$	$0.24 \pm 0.03$	$0.19 \pm 0.03$	$0.22 \pm 0.02$
<b>CIFAR10<sub>z</sub></b>	$0.13 \pm 0.03$	$0.13 \pm 0.01$	$1.70 \pm 2.60$	$0.24 \pm 0.01$	$0.52 \pm 0.60$

Table 1: The averaged Test MSE with standard deviation on the low- and high-dimensional scenarios.

**Non-i.i.d. data.** To set up a non-i.i.d. distribution of data between clients, samples were divided amongst the clients using a Dirichlet distribution  $Dir_S(\alpha)$  (Hsu et al., 2019), where  $\alpha$  determines the degree of heterogeneity across  $S$  clients. We used  $Dir_S(\alpha) = 0.3$  for each train, test, and validation samples. Given the non-i.i.d. data, for the low-dimensional scenario, we sample  $n = 20000$  points for each train, validation, and test set, while, for the high-dimensional scenario, we have  $n = 20000$  for the train set and  $n = 10000$  for the validation and test set.

**Hyperparameters.** We perform extensive grid-search to tune the learning rate. For FEDSGDA, we use a minibatch-size of 256. To avoid numerical instability, we standardize the observed  $Y$  values by removing the mean and scaling to unit variance. We perform five runs of each experiment and present the mean and standard deviation of the results.

**Observations and Discussion.** In figure (1), we first observe that SGDA and GDA algorithms perform at par with OADAM to fit the DEEPGMM estimator. It establishes that hyperparameter tuning is effective. With that, we further observe that the federated algorithms efficiently fit the estimated function to the true data-generating process even though the data is decentralized and non-i.i.d. Thus, it shows that the federated algorithm converges effectively. In Table 1 we present the test mean squared error (MSE) values. In many cases, the federated MSE values are close or better than the centralized results, which sufficiently demonstrate that our federated implementation achieves a convergent dynamics. We include additional experimental results in Appendix E that investigate the effects of heterogeneity. These experiments establish the efficacy of our method.

---

## EXISTENCE OF FEDERATED MIXED-STRATEGY EQUILIBRIUM AND ITS IMPLICATIONS

In this work, we presented the equilibrium solutions of federated zero-sum games through federated local minimax solutions for non-convex non-concave minimax optimization problems. The translation of the federated equilibrium as an approximately consistent GMM estimator for the clients was obtained through the gradient and Hessian dissimilarities across the clients, see Theorem 1, Theorem 2, and Definition 3. We note that our minimax optimization solution provides a federated pure strategy equilibrium. However, a pure strategy equilibrium can correspond to only full gradients and a full client participation setting. To elaborate,

- Firstly, with stochastic gradients on the clients, there will be no guarantee of descent (correspondingly, ascent) at an optimization step, which is available only in expectation in this case. However, in a pure strategy zero-sum game, the minimizing player (correspondingly, the maximizing player) takes a step to minimize (correspondingly, maximize) the game objective at each step.
- Secondly, the path to the saddle-point of a player in a pure strategy game should be retraceable/deterministic, which can not be possible with minimax optimization with stochastic gradients and/or partial client participation, considering the true random sampling.

Allowing for stochasticity, whether arising from stochastic gradients or client sampling for each communication round, would necessitate accommodating a distribution over multiple actions. Whereby the game ceases to be a pure strategy game, as the actions become non-deterministic, essentially, resulting in a mixed-strategy zero-sum game. It is well understood that, regardless of the analytical assumptions regarding the objective, mixed strategy solutions for zero-sum games exist (Jin et al., 2020).

However, for federated mixed strategy solutions, recovering a GMM estimator for a client is not immediate. To elaborate, there are no analogous necessary and sufficient conditions – the first-order and second-order necessary and sufficient conditions that we have for federated pure strategy solutions in Lemma 2 – for the mixed strategy solutions, which would correspond to a distribution over a set of the global model states synchronized across clients. Therefore, we can not directly apply Theorem 1. Still, we note here that a federated mixed strategy equilibrium will provide a robust federated GMM estimator compared to pure-strategy solutions, as it will output a probability distribution over a set of model states that accounts for the uncertainty across clients.

We leave the algorithm and characterization of a federated mixed strategy equilibrium solution for a robust federated GMM estimator as an open problem.

## REFERENCES

- Alejandro Almodóvar, Juan Parras, and Santiago Zazo. Propensity weighted federated learning for treatment effect estimation in distributed imbalanced environments. *Computers in Biology and Medicine*, 178:108779, 2024.
- Joshua D Angrist and Alan B Krueger. Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic perspectives*, 15(4):69–85, 2001.
- Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics: An empiricist’s companion*. Princeton university press, 2009.
- Emmanuel N Barron. *Game theory: an introduction*. John Wiley & Sons, 2024.
- Andrew Bennett, Nathan Kallus, and Tobias Schnabel. Deep generalized method of moments for instrumental variable analysis. *Advances in neural information processing systems*, 32, 2019.
- Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.

---

540 Zachary Charles and Dimitris Papailiopoulos. Stability and generalization of learning algorithms that  
541 converge to global optima. In *International conference on machine learning*, pp. 745–754. PMLR,  
542 2018.

543 Bapi Chatterjee, Vyacheslav Kungurtsev, and Dan Alistarh. Federated sgd with local asynchrony.  
544 In *2024 IEEE 44th International Conference on Distributed Computing Systems (ICDCS)*, pp.  
545 857–868. IEEE, 2024.

546 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training gans with  
547 optimism. In *International Conference on Learning Representations*, 2018.

548 Ittai Dayan, Holger R Roth, Aoxiao Zhong, Ahmed Harouni, Amilcare Gentili, Anas Z Abidin,  
549 Andrew Liu, Anthony Beardsworth Costa, Bradford J Wood, Chien-Sung Tsai, et al. Federated  
550 learning for predicting clinical outcomes in patients with covid-19. *Nature medicine*, 27(10):  
551 1735–1743, 2021.

552 Yuyang Deng and Mehrdad Mahdavi. Local stochastic gradient descent ascent: Convergence analysis  
553 and communication efficiency. In *International Conference on Artificial Intelligence and Statistics*,  
554 pp. 1387–1395. PMLR, 2021.

555 Farzan Farnia and Asuman Ozdaglar. Do gans always have nash equilibria? In *International*  
556 *Conference on Machine Learning*, pp. 3029–3039. PMLR, 2020.

557 Melike Gecer and Benoit Garbinato. Federated learning for mobility applications. *ACM Computing*  
558 *Surveys*, 56(5):1–28, 2024.

559 Lars Peter Hansen. Large sample properties of generalized method of moments estimators. *Econo-*  
560 *metrica: Journal of the econometric society*, pp. 1029–1054, 1982.

561 Jason Hartford, Greg Lewis, Kevin Leyton-Brown, and Matt Taddy. Deep iv: A flexible approach  
562 for counterfactual prediction. In *International Conference on Machine Learning*, pp. 1414–1423.  
563 PMLR, 2017.

564 Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational*  
565 *and Graphical Statistics*, 20(1):217–240, 2011.

566 Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2 edition,  
567 2012.

568 Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data  
569 distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.

570 Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave  
571 minimax optimization? In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*  
572 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*  
573 *Research*, pp. 4880–4889. PMLR, 07 2020. URL [https://proceedings.mlr.press/  
574 v119/jin20e.html](https://proceedings.mlr.press/v119/jin20e.html).

575 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and  
576 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In  
577 *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.

578 Diederik P Kingma. Adam: A method for stochastic optimization. *ICLR*, 2015.

579 Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations*  
580 *and Trends® in Machine Learning*, 12(4):307–392, 2019.

581 Greg Lewis and Vasilis Syrgkanis. Adversarial generalized method of moments, 2018. URL  
582 <https://arxiv.org/abs/1803.07164>.

583 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.  
584 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*,  
585 2:429–450, 2020.

- 
- 594 Yafen Liang, Wei Ruan, Yandong Jiang, Richard Smalling, Xiaoyi Yuan, and Holger K Eltzschig.  
595 Interplay of hypoxia-inducible factors and oxygen therapy in cardiovascular medicine. *Nature*  
596 *Reviews Cardiology*, 20(11):723–737, 2023.
- 597  
598 Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized  
599 non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:  
600 85–116, 2022.
- 601 Guodong Long, Yue Tan, Jing Jiang, and Chengqi Zhang. Federated learning for open banking. In  
602 *Federated learning: privacy and incentive*, pp. 240–254. Springer, 2020.
- 603  
604 Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal  
605 effect inference with deep latent-variable models. *Advances in neural information processing*  
606 *systems*, 30, 2017.
- 607  
608 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas.  
609 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*  
610 *gence and statistics*, pp. 1273–1282. PMLR, 2017.
- 611  
612 Maxence Noble, Aurélien Bellet, and Aymeric Dieuleveut. Differentially private federated learning  
613 on heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp.  
614 10110–10145. PMLR, 2022.
- 615  
616 Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D Lee, and Meisam Razaviyayn. Solving a  
617 class of non-convex min-max games using iterative first order methods. In *Advances in Neural*  
618 *Information Processing Systems*, volume 32, pp. 14934–14942, 2019.
- 619  
620 Wonsuk Oh and Girish N Nadkarni. Federated learning in health care using structured medical data.  
621 *Advances in kidney disease and health*, 30(1):4–16, 2023.
- 622  
623 Judea Pearl. Causal inference in statistics: An overview. 2009.
- 624  
625 Mohammad Rasouli, Tao Sun, and Ram Rajagopal. Fedgan: Federated generative adversarial  
626 networks for distributed data. *arXiv preprint arXiv:2006.07228*, 2020.
- 627  
628 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,  
629 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint*  
630 *arXiv:2003.00295*, 2020.
- 631  
632 Olav Reiersøl. *Confluence analysis by means of instrumental sets of variables*. PhD thesis, Almqvist  
633 & Wiksell, 1945.
- 634  
635 Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: general-  
636 ization bounds and algorithms. In *International conference on machine learning*, pp. 3076–3085.  
637 PMLR, 2017.
- 638  
639 Pranay Sharma, Rohan Panda, Gauri Joshi, and Pramod Varshney. Federated minimax optimization:  
640 Improved convergence analyses and algorithms. In *International Conference on Machine Learning*,  
641 pp. 19683–19730. PMLR, 2022.
- 642  
643 Wei Shen, Minhui Huang, Jiawei Zhang, and Cong Shen. Stochastic smoothed gradient descent  
644 ascent for federated minimax optimization. In *International Conference on Artificial Intelligence*  
645 *and Statistics*, pp. 3988–3996. PMLR, 2024.
- 646  
647 Alison Etheridge Steif, Jianqing Fan, Xiao-Li Meng, Bin Yu, David Madigan, and Juan Romo  
Manteiga. Nobel prize in economics. *IMS Bulletin*, 43(1), 2014.
- Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes.  
*Advances in neural information processing systems*, 33:21394–21405, 2020.
- Thanh Vinh Vo, Arnab Bhattacharyya, Young Lee, and Tze-Yun Leong. An adaptive kernel approach  
to federated learning of heterogeneous causal effects. *Advances in Neural Information Processing*  
*Systems*, 35:24459–24473, 2022a.

---

648 Thanh Vinh Vo, Young Lee, Trong Nghia Hoang, and Tze-Yun Leong. Bayesian federated estimation  
649 of causal effects from observational data. In *Uncertainty in Artificial Intelligence*, pp. 2024–2034.  
650 PMLR, 2022b.

651 Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal  
652 value of adaptive gradient methods in machine learning. *Advances in neural information processing*  
653 *systems*, 30, 2017.

654 Jeffrey M Wooldridge. Applications of generalized method of moments estimation. *Journal of*  
655 *Economic perspectives*, 15(4):87–100, 2001.

656 Philip Green Wright. *The tariff on animal and vegetable oils*. Number 26. Macmillan, 1928.

657 Xidong Wu, Jianhui Sun, Zhengmian Hu, Aidong Zhang, and Heng Huang. Solving a class of non-  
658 convex minimax optimization in federated learning. *Advances in Neural Information Processing*  
659 *Systems*, 36, 2024.

660 Yue Wu, Shuaicheng Zhang, Wenchao Yu, Yanchi Liu, Quanquan Gu, Dawei Zhou, Haifeng Chen,  
661 and Wei Cheng. Personalized federated learning under mixture of distributions. In *International*  
662 *Conference on Machine Learning*, pp. 37860–37879. PMLR, 2023.

663 Ruoxuan Xiong, Allison Koenecke, Michael Powell, Zhu Shen, Joshua T Vogelstein, and Susan  
664 Athey. Federated causal inference in heterogeneous observational data. *Statistics in Medicine*, 42  
665 (24):4418–4439, 2023.

666 Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning:  
667 State-of-the-art and research challenges. *ACM Computing Surveys*, 56(3):1–44, 2023.

668 Mishael Zedek. Continuity and location of zeros of linear combinations of polynomials. *Proceedings*  
669 *of the American Mathematical Society*, 16(1):78–84, 1965. ISSN 00029939, 10886826. URL  
670 <http://www.jstor.org/stable/2034005>.

671 Weijia Zhang, Lin Liu, and Jiuyong Li. Treatment effect estimation with disentangled latent factors.  
672 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10923–10930,  
673 2021.

674 Miaoxi Zhu, Li Shen, Bo Du, and Dacheng Tao. Stability and generalization of the decentralized  
675 stochastic gradient descent ascent algorithm. *Advances in Neural Information Processing Systems*,  
676 36, 2024.

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

---

# APPENDIX

702		
703		
704		
705	<b>A Federated Gradient Descent Ascent Algorithm Description</b>	<b>15</b>
706	A.1 Convergence of FEDGDA . . . . .	15
707		
708		
709	<b>B Proofs</b>	<b>22</b>
710	B.1 Proof of Lemma 1 . . . . .	22
711	B.2 Proof of Theorem 1 . . . . .	24
712	B.3 Consistency . . . . .	27
713	B.3.1 Assumptions . . . . .	27
714	B.3.2 Proof of Theorem 2 . . . . .	27
715		
716		
717		
718	<b>C Limit Points of FEDGDA</b>	<b>30</b>
719	C.1 FEDGDA Flow . . . . .	31
720	C.1.1 Proof of Proposition 1 . . . . .	32
721	C.2 Proof of Theorem 3 . . . . .	33
722		
723		
724	<b>D Related Work</b>	<b>34</b>
725		
726		
727	<b>E Benchmark Considerations and Additional Experiments</b>	<b>35</b>
728	E.1 The Experimental Benchmark Design . . . . .	35
729	E.2 Additional Experiments . . . . .	36
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

---

## A FEDERATED GRADIENT DESCENT ASCENT ALGORITHM DESCRIPTION

**Algorithm 1** FEDGDA running on a federated learning server to solve the minimax problem (17)

**Server Input:** initial global estimate  $\theta_1, \tau_1$ ; constant local learning rate  $\alpha_1, \alpha_2$ ; total  $N$  clients

**Output:** global model states  $\theta_{T+1}, \tau_{T+1}$

```

1: for synchronization round  $t = 1, \dots, T$  do
2:   server sends  $\theta_t, \tau_t$  to all clients
3:   for each  $i \in [N]$  in parallel do
4:      $\theta_{t,1}^i \leftarrow \theta_t, \tau_{t,1}^i \leftarrow \tau_t$ 
5:     for  $r = 1, 2, \dots, R$  do
6:        $\theta_{t,r+1}^i = \theta_{t,r}^i - \alpha_1 \nabla_{\theta} U_{\hat{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i)$ 
7:        $\tau_{t,r+1}^i = \tau_{t,r}^i + \alpha_2 \nabla_{\tau} U_{\hat{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i)$ 
8:     end for
9:      $(\Delta \theta_t^i, \Delta \tau_t) \leftarrow (\theta_{t,R+1}^i - \theta_t, \tau_{t,R+1}^i - \tau_t)$ 
10:   end for
11:    $(\Delta \theta_t, \Delta \tau_t) \leftarrow \frac{1}{N} \sum_{i \in [N]} (\Delta \theta_t^i, \Delta \tau_t^i)$ 
12:    $\theta_{t+1} \leftarrow (\theta_t + \Delta \theta_t), \tau_{t+1} \leftarrow (\tau_t + \Delta \tau_t)$ 
13: end for
14: return  $\theta_{T+1}; \tau_{T+1}$ 

```

---

### A.1 CONVERGENCE OF FEDGDA

We adapt the proof of Theorem 1 in (Sharma et al., 2022) for the SGDA algorithm proposed in (Deng & Mahdavi, 2021) for the FEDGDA algorithm 1 for smooth non-convex- PL problems.

**Assumption 5** (Polyak Łojasiewicz (PL) condition in  $\tau$ ). The function  $U_{\hat{\theta}}$  satisfies  $\mu - PL$  condition in  $\tau$ ,  $\mu > 0$ , if for any fixed  $\theta$ ,  $\arg \max_{\tau'} U_{\hat{\theta}}(\theta, \tau') \neq \phi$  and  $\|\nabla_{\tau} U_{\hat{\theta}}(\theta, \tau)\|^2 \geq 2\mu (\max_{\tau'} U_{\hat{\theta}}(\theta, \tau') - U_{\hat{\theta}}(\theta, \tau))$ .

We use the following result about the smoothness of  $\Phi(\cdot)$ .

**Lemma 3.** (Nouiehed et al., 2019) *If the function  $U_{\hat{\theta}}(\theta, \cdot)$  satisfies Assumptions 2, 5 ( $L$ -smoothness and  $\mu$ -PL condition in  $\tau$ ), then  $\Phi(\theta)$  is  $L_{\Phi}$ -smooth with  $L_{\Phi} = \kappa L/2 + L$ , where  $\kappa = L/\mu$  is the condition number.*

**Lemma 4** (One-Step Envelope Descent). (Deng & Mahdavi, 2021) *Suppose the local client loss functions  $\{U_{\hat{\theta}}^i(\theta, \tau)\}$  satisfy Assumptions 2, 5. Then the iterates generated by FEDGDA satisfy:*

$$\begin{aligned} \Phi(\theta_{t+1}) &\leq \Phi(\theta_t) - \frac{\alpha_1}{2} \|\nabla \Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_{\Phi} \alpha_1) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\hat{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\ &\quad + \frac{2\alpha_1 L^2}{\mu} (\Phi(\theta_t) - U_{\hat{\theta}}(\theta_t, \tau_t)) + 2\alpha_1 L^2 R \Delta_t^{\theta, \tau} \end{aligned}$$

where the synchronization error is defined as:

$$\Delta_t^{\theta, \tau} := \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R (\|\theta_{t,r}^i - \theta_t\|^2 + \|\tau_{t,r}^i - \tau_t\|^2).$$

*Proof.* Using Lemma 3,  $\Phi(\cdot)$  is  $L_{\Phi} = \kappa L/2 + L$ -smooth, and together with the updating rule, we have:

$$\begin{aligned} \Phi(\theta_{t+1}) &\leq \Phi(\theta_t) + \langle \nabla \Phi(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{L_{\Phi}}{2} \|\theta_{t+1} - \theta_t\|^2 \\ &\leq \Phi(\theta_t) - \alpha_1 \left\langle \nabla \Phi(\theta_t), \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\hat{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\rangle + \frac{L_{\Phi}}{2} \alpha_1^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\hat{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \end{aligned}$$

Using the identity  $\langle \mathbf{a}, \mathbf{b} \rangle = -\frac{1}{2}\|\mathbf{a} - \mathbf{b}\|^2 + \frac{1}{2}\|\mathbf{a}\|^2 + \frac{1}{2}\|\mathbf{b}\|^2$ , we have:

$$\begin{aligned}
& \Phi(\theta_{t+1}) - \Phi(\theta_t) \\
& \leq -\frac{\alpha_1}{2} \|\nabla\Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + \alpha_1 \left\| \nabla\Phi(\theta_t) - \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t) \right\|^2 \\
& \quad + \alpha_1 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t) \right\|^2 + \frac{L_{\Phi}}{2} \alpha_1^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
& \leq -\frac{\alpha_1}{2} \|\nabla\Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_{\Phi} \alpha_1) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + \alpha_1 L^2 \|\phi(\theta_t) - \tau_t\|^2 \\
& \quad + \alpha_1 L^2 \frac{R}{N} \sum_{i=1}^N \sum_{r=1}^R (2\|\theta_{t,r}^i - \theta_t\|^2 + 2\|\tau_{t,r}^i - \tau_t\|^2) \tag{22}
\end{aligned}$$

$$\begin{aligned}
& \leq -\frac{\alpha_1}{2} \|\nabla\Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_{\Phi} \alpha_1) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
& \quad + \alpha_1 L^2 \|\phi(\theta_t) - \tau_t\|^2 + 2\alpha_1 L^2 R \Delta_t^{\theta, \tau} \tag{23}
\end{aligned}$$

Using the quadratic growth property of  $\mu$ -PL function  $U_{\bar{\theta}}(\theta, \cdot)$ , i.e.,  $\frac{\mu}{2} \|\tau - \Phi(\theta)\|^2 \leq \max_{\tau'} U_{\bar{\theta}}(\theta, \tau') - U_{\bar{\theta}}(\theta, \tau)$ ,  $\forall \theta, \tau$ , where  $\Phi(\theta) := \arg \max_{\tau'} U_{\bar{\theta}}(\theta, \tau')$ , we have

$$\begin{aligned}
& \Phi(\theta_{t+1}) - \Phi(\theta_t) \\
& \leq -\frac{\alpha_1}{2} \|\nabla\Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_{\Phi} \alpha_1) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
& \quad + \frac{2\alpha_1 L^2}{\mu} (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) + 2\alpha_1 L^2 R \Delta_t^{\theta, \tau} \tag{24}
\end{aligned}$$

□

**Lemma 5.** (Sharma et al., 2022) Suppose the local loss functions  $\{U_{\bar{\theta}}^i\}$  satisfy Assumptions 2 and 3. Further, in Algorithm 1, we choose step-sizes  $\alpha_1, \alpha_2$  satisfying  $\alpha_2 \leq 1/\mu$ ,  $\frac{\alpha_1}{\alpha_2} \leq \frac{1}{8\kappa^2}$ . Then the following inequality holds.

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) \\
& \leq \frac{2(\Phi(\theta_1) - U_{\bar{\theta}}(\theta_1, \tau_1))}{\alpha_2 \mu T} + \frac{2L^2 R}{\mu \alpha_2} (2\alpha_1(1 - \alpha_2 \mu) + \alpha_2) \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau} + (1 - \alpha_2 \mu) \frac{\alpha_1}{\alpha_2 \mu} \frac{1}{T} \sum_{t=1}^T \|\nabla\Phi(\theta_t)\|^2 \\
& \quad + \left[ (1 - \alpha_2 \mu) \frac{\alpha_1^2}{2} (L + L_{\Phi}) + \alpha_2 L^2 \alpha_1^2 \right] \frac{2}{\alpha_2 \mu T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2
\end{aligned}$$

*Proof.* Using  $L$ -smoothness of  $U_{\bar{\theta}}(\theta, \cdot)$

$$U_{\bar{\theta}}(\theta_{t+1}, \tau_t) + \langle \nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t), \tau_{t+1} - \tau_t \rangle - \frac{L}{2} \|\tau_{t+1} - \tau_t\|^2 \leq U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1})$$



Using the update rule in Algorithm 1

$$\begin{aligned}
U_{\bar{\theta}}(\theta_{t+1}, \tau_t) &\leq U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1}) - \alpha_2 \left\langle \nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t), \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\rangle \\
&\quad + \frac{\alpha_2^2 L}{2} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&= U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1}) - \frac{\alpha_2}{2} \|\nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t)\|^2 - \frac{\alpha_2}{2} (1 - \alpha_2 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + \frac{\alpha_2}{2} \left\| \nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t) - \nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t) + \nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t) - \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\leq U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1}) - \frac{\alpha_2}{2} \|\nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t)\|^2 - \frac{\alpha_2}{2} (1 - \alpha_2 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + \alpha_2 L^2 \|\theta_{t+1} - \theta_t\|^2 + \alpha_2 L^2 R \Delta_t^{\theta, \tau}, \tag{25}
\end{aligned}$$

Note that

$$\|\theta_{t+1} - \theta_t\|^2 = \alpha_1^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2. \tag{26}$$

Also, using Assumption 5,

$$\|\nabla_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau_t)\|^2 \geq 2\mu \left( \max_{\tau} U_{\bar{\theta}}(\theta_{t+1}, \tau) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right) = 2\mu \left( \Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right). \tag{27}$$

Substituting (26), (27) in (25), and rearranging the terms, we get

$$\begin{aligned}
&\alpha_2 \mu \left( \Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right) \\
&\leq U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) - \frac{\alpha_2}{2} (1 - \alpha_2 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + \alpha_2 L^2 \left[ \alpha_1^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \right] + \alpha_2 L^2 R \Delta_t^{\theta, \tau} \\
&\Rightarrow \left( \Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1}) \right) \\
&\leq (1 - \alpha_2 \mu) \left( \Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right) - \frac{\alpha_2}{2} (1 - \alpha_2 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + \alpha_1^2 \alpha_2 L^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + \alpha_2 L^2 R \Delta_t^{\theta, \tau}. \tag{28}
\end{aligned}$$

Next, we bound  $\left( \Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right)$ .

$$\begin{aligned}
&\Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \\
&= \underbrace{\left( \Phi(\theta_{t+1}) - \Phi(\theta_t) \right)}_{T_1} + \left( \Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t) \right) + \underbrace{\left( U_{\bar{\theta}}(\theta_t, \tau_t) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \right)}_{T_2} \tag{29}
\end{aligned}$$

$T_1$  is bounded in Lemma 4. We next bound  $T_2$ . Using  $L$ -smoothness of  $U_{\bar{\theta}}(\cdot, \tau_t)$ ,

$$\begin{aligned}
& U_{\bar{\theta}}(\theta_t, \tau_t) + \langle \nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t), \theta_{t+1} - \theta_t \rangle - \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \leq U_{\bar{\theta}}(\theta_{t+1}, \tau_t) \\
\Rightarrow T_2 &= (U_{\bar{\theta}}(\theta_t, \tau_t) - U_{\bar{\theta}}(\theta_{t+1}, \tau_t)) \\
&\leq \alpha_1 \left\langle \nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t), \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\rangle + \frac{\alpha_1^2 L}{2} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\leq \frac{\alpha_1}{2} \left( \|\nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t)\|^2 + \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \right) + \frac{\alpha_1^2 L}{2} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\leq \alpha_1 \left( \|\nabla \Phi(\theta_t)\|^2 + \|\nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t) - \nabla \Phi(\theta_t)\|^2 \right) + \frac{\alpha_1}{2} (1 + \alpha_1 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\stackrel{(a)}{\leq} \alpha_1 \|\nabla \Phi(\theta_t)\|^2 + \alpha_1 L^2 \|\tau_t - \tau^*(\theta_t)\|^2 + \frac{\alpha_1}{2} (1 + \alpha_1 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\leq \alpha_1 \|\nabla \Phi(\theta_t)\|^2 + \frac{2\alpha_1 L^2}{\mu} (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) + \frac{\alpha_1}{2} (1 + \alpha_1 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2.
\end{aligned} \tag{30}$$

where (a) follows from Assumption 2 and Assumption 3. Also, recall that  $\tau^*(\theta) \in \arg \max_{\tau'} U_{\bar{\theta}}(\theta, \tau')$ . (30) follows from the quadratic growth property of  $\mu$ -PL functions. Substituting the bounds on  $T_1, T_2$  from Lemma 4 and (30) respectively, in (28), we get

$$\begin{aligned}
& (\Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1})) \\
&\leq (1 - \alpha_2 \mu) \left( 1 + \frac{4\alpha_1 L^2}{\mu} \right) (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) \\
&\quad + (1 - \alpha_2 \mu) \left[ -\frac{\alpha_1}{2} \|\nabla \Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_{\Phi} \alpha_1) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + 2\alpha_1 L^2 R \Delta_t^{\theta, \tau} \right] \\
&\quad + (1 - \alpha_2 \mu) \left[ \alpha_1 \|\nabla \Phi(\theta_t)\|^2 + \frac{\alpha_1}{2} (1 + \alpha_1 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \right] \\
&\quad - \frac{\alpha_2}{2} (1 - \alpha_2 L) \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + \alpha_1^2 \alpha_2 L^2 \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + \alpha_2 L^2 R \Delta_t^{\theta, \tau} \\
&\leq \left( 1 - \frac{\alpha_2 \mu}{2} \right) (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) + \alpha_2 L^2 R \Delta_t^{\theta, \tau} \\
&\quad + \left[ (1 - \alpha_2 \mu) \frac{\alpha_1^2}{2} (L + L_{\Phi}) + \alpha_2 L^2 \alpha_1^2 \right] \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 \\
&\quad + (1 - \alpha_2 \mu) \left[ \frac{\alpha_1}{2} \|\nabla \Phi(\theta_t)\|^2 + 2\alpha_1 L^2 R \Delta_t^{\theta, \tau} \right],
\end{aligned} \tag{31}$$

where we choose  $\alpha_1$  such that  $(1 - \alpha_2\mu) \left(1 + \frac{4\alpha_1 L^2}{\mu}\right) \leq (1 - \frac{\alpha_2\mu}{2})$ . This holds if  $\frac{4\alpha_1 L^2}{\mu} \leq \frac{\alpha_2\mu}{2} \Rightarrow \alpha_1 \leq \frac{\alpha_2}{8\kappa^2}$ . Summing (31) over  $t = 1, \dots, T$ , and rearranging the terms, we get

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T (\Phi(\theta_{t+1}) - U_{\bar{\theta}}(\theta_{t+1}, \tau_{t+1})) \\ & \leq \left(1 - \frac{\alpha_2\mu}{2}\right) \frac{1}{T} \sum_{t=1}^T (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) + L^2 R (2\alpha_1(1 - \alpha_2\mu) + \alpha_2) \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau} \\ & \quad + \left[ (1 - \alpha_2\mu) \frac{\alpha_1^2}{2} (L + L_\Phi) + \alpha_2 L^2 \alpha_1^2 \right] \frac{1}{T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + (1 - \alpha_2\mu) \frac{\alpha_1}{2} \frac{1}{T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 \end{aligned}$$

Rearranging the terms, we get

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T (\Phi(\theta_t) - U_{\bar{\theta}}(\theta_t, \tau_t)) \\ & \leq \frac{2}{\alpha_2\mu} \left[ \frac{\Phi(\theta_1) - U_{\bar{\theta}}(\theta_1, \tau_1)}{T} - \frac{\Phi(\theta_{T+1}) - U_{\bar{\theta}}(\theta_{T+1}, \tau_{T+1})}{T} \right] + \frac{2L^2 R}{\alpha_2\mu} (2\alpha_1(1 - \alpha_2\mu) + \alpha_2) \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau} \\ & \quad + \left[ (1 - \alpha_2\mu) \frac{\alpha_1^2}{2} (L + L_\Phi) + \alpha_2 L^2 \alpha_1^2 \right] \frac{2}{\alpha_2\mu T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + (1 - \alpha_2\mu) \frac{\alpha_1}{\alpha_2\mu T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 \\ & \leq \frac{2(\Phi(\theta_1) - U_{\bar{\theta}}(\theta_1, \tau_1))}{\alpha_2\mu T} + \frac{2L^2 R}{\alpha_2\mu} (2\alpha_1(1 - \alpha_2\mu) + \alpha_2) \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau} \\ & \quad + \left[ (1 - \alpha_2\mu) \frac{\alpha_1^2}{2} (L + L_\Phi) + \alpha_2 L^2 \alpha_1^2 \right] \frac{2}{\alpha_2\mu T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2 + (1 - \alpha_2\mu) \frac{\alpha_1}{\alpha_2\mu T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 \end{aligned}$$

since  $\Phi(\theta_T) := \arg \max_{\tau} U_{\bar{\theta}}(\theta_T, \tau)$ , which concludes the proof.  $\square$

**Lemma 6.** Suppose the local loss functions  $\{U_{\bar{\theta}}^i\}$  satisfy Assumptions 2 and 3. Further, in Algorithm 1, using bounded gradient assumption, i.e.,  $\|\nabla U_{\bar{\theta}}^i(\theta, \tau)\| \leq G$ , we choose step-sizes  $\alpha_1, \alpha_2 \leq \frac{1}{8RL}$ . Then, the iterates  $\{\theta_t, \tau_t\}$  generated by Algorithm 1 satisfy

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau} & := \frac{1}{T} \sum_{t=1}^T \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \left( \|\theta_{t,r}^i - \theta_t\|^2 + \|\tau_{t,r}^i - \tau_t\|^2 \right) \\ & \leq 6R(R-1)^2 [(\zeta_{\theta}^2 \alpha_1^2 + \zeta_{\tau}^2 \alpha_2^2) + (\alpha_1^2 + \alpha_2^2) G^2]. \end{aligned}$$

*Proof of Lemma 6.* We define the separate synchronization errors for  $\theta$  and  $\tau$

$$\Delta_t^{\theta} := \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \|\theta_{t,r}^i - \theta_t\|^2, \quad \Delta_t^{\tau} := \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \|\tau_{t,r}^i - \tau_t\|^2,$$

such that  $\Delta_t^{\theta, \tau} = \Delta_t^{\theta} + \Delta_t^{\tau}$ . We first bound the  $\theta$ -synchronization error  $\Delta_t^{\theta}$ . Then,

$$\begin{aligned} \Delta_t^{\theta} & := \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \|\theta_{t,r}^i - \theta_t\|^2 \\ & = \alpha_1^2 \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \left\| \sum_{j=1}^r \nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,j}^i, \tau_{t,j}^i) \right\|^2 \\ & \leq \alpha_1^2 \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R (r-1) \sum_{j=1}^r \|\nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,j}^i, \tau_{t,j}^i)\|^2. \end{aligned}$$

This can be written as

$$\begin{aligned}
\Delta_t^\theta &\stackrel{(a)}{\leq} \alpha_1^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{R-1} \|\nabla_\theta U_{\bar{\theta}}^i(\theta_{t,j}^i, \tau_{t,j}^i)\|^2 \sum_{r=j+1}^R (r-1) \\
&\stackrel{(b)}{\leq} \alpha_1^2 \frac{(R-1)^2}{2} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{R-1} \|\nabla_\theta U_{\bar{\theta}}^i(\theta_{t,j}^i, \tau_{t,j}^i)\|^2 \\
&\stackrel{(c)}{\leq} \alpha_1^2 \frac{(R-1)^2}{2} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^R \|\nabla_\theta U_{\bar{\theta}}^i(\theta_{t,j}^i, \tau_{t,j}^i)\|^2,
\end{aligned}$$

where (a) follows from rewriting the sum, (b) follows since  $\sum_{r=j+1}^R (r-1) \leq \frac{R^2}{2}$  and (c) follows because a positive quantity is being added. Now,

$$\begin{aligned}
\Delta_t^\theta &\leq \alpha_1^2 \frac{(R-1)^2}{2} \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \|\nabla_\theta U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_\theta U_{\bar{\theta}}^i(\theta_t, \tau_t) + \nabla_\theta U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_\theta U_{\bar{\theta}}(\theta_t, \tau_t) + \nabla_\theta U_{\bar{\theta}}(\theta_t, \tau_t)\|^2 \\
&\leq \alpha_1^2 \frac{(R-1)^2}{2} \frac{1}{N} \sum_{i=1}^N \sum_{r=1}^R \left[ 3 \|\nabla_\theta U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_\theta U_{\bar{\theta}}^i(\theta_t, \tau_t)\|^2 + 3 \|\nabla_\theta U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_\theta U_{\bar{\theta}}(\theta_t, \tau_t)\|^2 \right] \\
&\quad + 3\alpha_1^2 \frac{R(R-1)^2}{2} \|\nabla_\theta U_{\bar{\theta}}(\theta_t, \tau_t)\|^2 \\
&\leq 6L^2\alpha_1^2 \frac{(R-1)^2}{2} \Delta_t^{\theta,\tau} + 3\zeta_\theta^2\alpha_1^2 \frac{R(R-1)^2}{2} + 3\alpha_1^2 \frac{R(R-1)^2}{2} \|\nabla_\theta U_{\bar{\theta}}(\theta_t, \tau_t)\|^2. \quad (32)
\end{aligned}$$

Similarly, for the synchronization error  $\Delta_t^\tau$ , we have

$$\Delta_t^\tau \leq 6L^2\alpha_2^2 \frac{(R-1)^2}{2} \Delta_t^{\theta,\tau} + 3\zeta_\tau^2\alpha_2^2 \frac{R(R-1)^2}{2} + 3\alpha_2^2 \frac{R(R-1)^2}{2} \|\nabla_\tau U_{\bar{\theta}}(\theta_t, \tau_t)\|^2. \quad (33)$$

Using bounded gradient assumption, i.e.,  $\|\nabla U_{\bar{\theta}}^i(\theta, \tau)\| \leq G$ , and adding (32) and (33), we obtain

$$\Delta_t^{\theta,\tau} \leq 6L^2(\alpha_1^2 + \alpha_2^2)(R-1)^2 \Delta_t^{\theta,\tau} + 3(\zeta_\theta^2\alpha_1^2 + \zeta_\tau^2\alpha_2^2)R(R-1)^2 + 3(\alpha_1^2 + \alpha_2^2)R(R-1)^2G^2.$$

For our choice of  $\alpha_1$  and  $\alpha_2$ , we have  $6L^2(\alpha_1^2 + \alpha_2^2)(R-1)^2 \leq \frac{1}{2}$ , thus

$$\Delta_t^{\theta,\tau} \leq 6R(R-1)^2 [(\zeta_\theta^2\alpha_1^2 + \zeta_\tau^2\alpha_2^2) + (\alpha_1^2 + \alpha_2^2)G^2].$$

Averaging across all the communication rounds  $t = 1, \dots, T$  proves the lemma.  $\square$

**Theorem 4** (Convergence of FEDGDA). *Suppose the local loss functions  $\{U_{\bar{\theta}}^i\}_i$  satisfy Assumptions 2,3 and have bounded gradients, and the global function  $U_{\bar{\theta}}$  satisfies 5. Suppose the step-sizes  $\alpha_1, \alpha_2$  are chosen such that  $\alpha_2 \leq \frac{1}{8LR}$ ,  $\frac{\alpha_1}{\alpha_2} = \frac{1}{8\kappa^2}$ , where  $\kappa = \frac{L}{\mu}$  is the condition number. Then for the output  $\bar{\theta}_T$  of Algorithm 1, the following holds.*

$$\begin{aligned}
\|\nabla\Phi(\bar{\theta}_T)\|^2 &= \frac{1}{T} \sum_{t=1}^T \|\nabla\Phi(\theta_t)\|^2 \\
&\leq \underbrace{\mathcal{O}\left(\kappa^2 \frac{\Delta_\Phi}{\alpha_2 T}\right)}_{\text{Error with full synchronization}} + \underbrace{\mathcal{O}\left(L^2\kappa^2 R(R-1)^2 [\alpha_2^2 (G^2 + \zeta_\tau^2) + \alpha_1^2 \zeta_\theta^2]\right)}_{\text{Error due to local updates}}, \quad (34)
\end{aligned}$$

where  $\Phi(\theta) := \max_\tau U_{\bar{\theta}}(\theta, \tau)$  is the envelope function,  $\Delta_\Phi := \Phi(\theta_1) - \min_\theta \Phi(\theta)$ . Using  $\alpha_2 = \sqrt{\frac{N}{LT}}$  and  $\alpha_1 = \frac{1}{8\kappa^2} \sqrt{\frac{N}{LT}}$ , we get

$$\|\nabla\Phi(\bar{\theta}_T)\|^2 \leq \mathcal{O}\left(\frac{\kappa^2 \Delta_\Phi}{\sqrt{NT}} + \kappa^2 R(R-1)^2 \frac{N(\sigma^2 + \zeta_\theta^2 + \zeta_\tau^2)}{T}\right).$$

1080 *Proof.* We start by summing the expression in Lemma 4 over  $t = 1, \dots, T$ .

$$1081$$

$$1082 \frac{1}{T} \sum_{t=1}^T (\Phi(\theta_{t+1}) - \Phi(\theta_t)) \leq -\frac{\alpha_1}{2} \frac{1}{T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 - \frac{\alpha_1}{2} (1 - L_\Phi \alpha_1) \frac{1}{T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^n \sum_{r=1}^R \nabla_\theta U_\theta^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2$$

$$1083$$

$$1084$$

$$1085 + \frac{2\alpha_1 L^2}{\mu} \frac{1}{T} \sum_{t=1}^T [\Phi(\theta_t) - U_{\hat{\theta}}(\theta_t, \tau_t)] + 2\alpha_1 L^2 R \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau}. \quad (35)$$

$$1086$$

$$1087$$

1088 Substituting the bound on  $\frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau}$  from Lemma 6, and the bound on  
1089  $\frac{1}{T} \sum_{t=1}^T (\Phi(\theta_t) - U_{\hat{\theta}}(\theta_t, \tau_t))$  from Lemma 5, and rearranging the terms in (35), we get

$$1090$$

$$1091 \frac{1}{T} (\Phi(\theta_T) - \Phi(\theta_1))$$

$$1092$$

$$1093 \leq - \underbrace{\left( \frac{\alpha_1}{2} - (1 - \alpha_2 \mu) \frac{2\alpha_1^2 L^2}{\alpha_2 \mu^2} \right)}_{\geq \alpha_1/4} \frac{1}{T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2$$

$$1094$$

$$1095 - \underbrace{\frac{\alpha_1}{2} \left( 1 - L_\Phi \alpha_1 - \frac{8L^2}{\mu^2 \alpha_2} \left[ (1 - \alpha_2 \mu) \frac{\alpha_1^2}{2} (L + L_\Phi) + \alpha_2 L^2 \alpha_1^2 \right] \right)}_{\geq 0} \frac{1}{T} \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^n \nabla_\theta U_\theta^i(\theta_{t,r}^i, \tau_{t,r}^i) \right\|^2$$

$$1096$$

$$1097$$

$$1098 + \left[ \frac{2\alpha_1 L^2}{\mu} \left( \frac{2L^2 R}{\mu} + \frac{4\alpha_1 R L^2 (1 - \alpha_2 \mu)}{\mu \alpha_2} \right) + 2\alpha_1 L^2 R \right] \frac{1}{T} \sum_{t=1}^T \Delta_t^{\theta, \tau}$$

$$1099$$

$$1100 + \frac{4\alpha_1 L^2}{\mu} \frac{(\Phi(\theta_1) - U_{\hat{\theta}}(\theta_1, \tau_1))}{\alpha_2 \mu T}. \quad (36)$$

$$1101$$

$$1102$$

$$1103$$

$$1104$$

$$1105$$

$$1106$$

$$1107$$

1108 Here,  $\frac{\alpha_1}{2} - \frac{2\alpha_1^2(1-\mu\alpha_2)L^2}{\mu^2\alpha_2} \geq \frac{\alpha_1}{4}$  holds since  $\frac{\alpha_1}{\alpha_2} \leq \frac{1}{8\kappa^2}$ . Also,  $1 - L_\Phi \alpha_1 -$   
1109  $\frac{8L^2}{\mu^2 \alpha_2} \left[ (1 - \alpha_2 \mu) \frac{\alpha_1^2}{2} (L + L_\Phi) + \alpha_2 L^2 \alpha_1^2 \right] \geq 0$  follows from the bounds on  $\alpha_1, \alpha_2$ . Rearranging the  
1110 terms in Equation (36) and using lemma (6), we get

$$1111$$

$$1112 \frac{1}{T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 \leq \frac{4(\Phi(\theta_1) - \Phi(\theta_T))}{\alpha_1 T}$$

$$1113$$

$$1114 + \frac{4}{\alpha_1} 2\alpha_1 L^2 R \left[ 1 + 2\kappa^2 + 4\kappa^2 \frac{\alpha_1}{\alpha_2} \right] 6R(R-1)^2 [(\alpha_1^2 + \alpha_2^2) G^2 + (\alpha_1^2 \zeta_\theta^2 + \alpha_2^2 \zeta_\tau^2)]$$

$$1115$$

$$1116 + \frac{4}{\alpha_1} \frac{4\alpha_1 \kappa^2}{\alpha_2} \frac{(\Phi(\theta_1) - U_{\hat{\theta}}(\theta_1, \tau_1))}{T}$$

$$1117$$

$$1118 \stackrel{(a)}{\leq} \frac{4\Delta_\Phi}{\alpha_1 T} + 8L^2 R [2 + 2\kappa^2] 6R(R-1)^2 [(\alpha_1^2 + \alpha_2^2) G^2 + (\alpha_1^2 \zeta_\theta^2 + \alpha_2^2 \zeta_\tau^2)] + \frac{16\kappa^2 \Delta_\Phi}{\alpha_2 T}$$

$$1119$$

$$1120 \stackrel{(b)}{\leq} \frac{4\Delta_\Phi}{\alpha_1 T} + 192L^2 \kappa^2 R(R-1)^2 [(\alpha_1^2 + \alpha_2^2) G^2 + \alpha_1^2 \zeta_\theta^2 + \alpha_2^2 \zeta_\tau^2] + \frac{16\kappa^2 \Delta_\Phi}{\alpha_2 T}$$

$$1121$$

$$1122 = \mathcal{O} \left( \frac{\Delta_\Phi}{\alpha_1 T} + \kappa^2 \frac{\Delta_\Phi}{\alpha_2 T} + L^2 \kappa^2 R(R-1)^2 [(\alpha_1^2 + \alpha_2^2) G^2 + \alpha_1^2 \zeta_\theta^2 + \alpha_2^2 \zeta_\tau^2] \right).$$

$$1123$$

$$1124 = \underbrace{\mathcal{O} \left( \kappa^2 \frac{\Delta_\Phi}{\alpha_2 T} \right)}_{\text{Error with full synchronization}} + \underbrace{\mathcal{O} \left( L^2 \kappa^2 R(R-1)^2 [\alpha_2^2 (G^2 + \zeta_\tau^2) + \alpha_1^2 \zeta_\theta^2] \right)}_{\text{Error due to local updates}}. \quad (\text{since } \kappa \geq 1)$$

$$1125$$

$$1126$$

$$1127$$

$$1128$$

$$1129$$

$$1130$$

1131 where, we denote  $\Delta_\Phi := \Phi(\theta_1) - \min_\theta \Phi(\theta)$ . (a) follows from  $\frac{\alpha_1}{\alpha_2} \leq \frac{1}{8\kappa^2}$ ; (b) follows since  $\kappa \geq 1$   
1132 and  $L_\Phi \geq L$ . Therefore,  $\frac{8\kappa^2 \alpha_1}{\alpha_2} \frac{\alpha_1 \sigma^2}{N} (L + L_\Phi) \leq \frac{\alpha_1 \sigma^2}{N} (L + L_\Phi) \leq \frac{2L_\Phi \alpha_1 \sigma^2}{N}$ , which results in  
1133 Equation (34).

Using  $\alpha_2 = \sqrt{\frac{N}{LT}}$  and  $\alpha_1 = \frac{1}{8\kappa^2} \sqrt{\frac{N}{LT}} \leq \frac{\alpha_2}{8\kappa^2}$ , and since  $\kappa \geq 1$ , we get

$$\frac{1}{T} \sum_{t=1}^T \|\nabla \Phi(\theta_t)\|^2 \leq \mathcal{O} \left( \frac{\kappa^2 \Delta_\Phi}{\sqrt{NT}} + \kappa^2 R(R-1)^2 \frac{N}{T} \left[ G^2 + \frac{\zeta_\theta^2}{\kappa^4} + \zeta_\tau^2 \right] \right).$$

□

## B PROOFS

### B.1 PROOF OF LEMMA 1

**Lemma 7** (Restatement of Lemma 1). *Let  $\mathcal{F} = \text{span}\{f_j^i \mid i \in [N], j \in [m]\}$ . An equivalent objective function for the federated moment estimation optimization problem (16) is given by:*

$$\|\psi_N(f; \theta)\|_{\bar{\theta}}^2 = \sup_{\substack{f^i \in \mathcal{F} \\ \forall i \in [N]}} \frac{1}{N} \sum_{i=1}^N \left( \psi_{n_i}(f^i; \theta) - \frac{1}{4} \mathcal{C}_{\bar{\theta}}(f^i; f^i) \right), \text{ where} \quad (37)$$

$$\psi_{n_i}(f^i; \theta) := \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i) (Y_k^i - g^i(X_k^i; \theta)), \text{ and } \mathcal{C}_{\bar{\theta}}(f^i, f^i) := \frac{1}{n_i} \sum_{k=1}^{n_i} (f^i(Z_k^i))^2 (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2.$$

*Proof.* Let  $\psi = (\frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f_1^i; \theta), \frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f_2^i; \theta), \dots, \frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f_m^i; \theta))$ .

We know that  $\|v\|^2 = v^\top C_{\bar{\theta}}^{-1} v$  and the associated dual norm is obtained as  $\|v\|_*^2 = \sup_{\|v\| \leq 1} v^\top v = v^\top C_{\bar{\theta}} v$ .

Using the definition of the dual norm,

$$\begin{aligned} \|\psi\| &= \sup_{\|v\|_* \leq 1} v^\top \psi \\ \|\psi\|^2 &= \sup_{\|v\|_* \leq \|\psi\|} v^\top \psi \\ \|\psi\|^2 &= \sup_{v^\top C_{\bar{\theta}} v \leq \|\psi\|^2} v^\top \psi. \end{aligned} \quad (38)$$

We now find the equivalent dual optimization problem for (38).

The Lagrangian of the constrained maximization problem (38) is given as

$$\mathcal{L}(v, \lambda) = v^\top \psi + \lambda(v^\top C_{\bar{\theta}} v - \|\psi\|^2), \text{ where } \lambda \leq 0.$$

To maximize  $\mathcal{L}(v, \lambda)$  w.r.t.  $v$ , put  $\frac{\partial \mathcal{L}}{\partial v} = \psi + 2\lambda C_{\bar{\theta}} v = 0$  to obtain  $v = \frac{-1}{2\lambda} C_{\bar{\theta}}^{-1} \psi$ .

When  $\|\psi\| > 0$ ,  $v = 0$  satisfies Slater's condition as a strictly feasible interior point of the constraint  $v^\top C_{\bar{\theta}} v - \|\psi\|^2 \leq 0$ . Since  $C_{\bar{\theta}} \succeq 0$ , the quadratic form  $v^\top C_{\bar{\theta}} v$  is convex in  $v$ , and the objective  $v^\top \psi$  is linear. Hence, for this convex optimization problem, the Slater's condition applies. Thus, strong duality holds. Substituting  $v = \frac{-1}{2\lambda} C_{\bar{\theta}}^{-1} \psi$  in the Lagrangian gives

$$\begin{aligned} \mathcal{L}^*(\lambda) &= \frac{-1}{2\lambda} \psi^\top C_{\bar{\theta}}^{-1} \psi + \frac{1}{4\lambda} \psi^\top C_{\bar{\theta}}^{-1} \psi - \lambda \|\psi\|^2 \\ &= -\frac{\|\psi\|^2}{4\lambda} - \lambda \|\psi\|^2. \end{aligned}$$

Hence, the dual becomes  $\|\psi\|^2 = \inf_{\lambda < 0} \{\mathcal{L}^*(\lambda)\}$ . Thus, the equivalent dual optimization problem for (38) is given as

$$\|\psi\|^2 = \inf_{\lambda < 0} \left\{ -\frac{\|\psi\|^2}{4\lambda} - \lambda \|\psi\|^2 \right\}. \quad (39)$$

1188 Putting  $\frac{\partial \mathcal{L}}{\partial \lambda} = \frac{\|\psi\|^2}{4\lambda^2} - \|\psi\|^2 = 0$  gives  $\lambda = \frac{-1}{2}$ . Thus, due to strong duality  $\|\psi\|^2 = \sup_v \mathcal{L}(v, \frac{-1}{2}) =$   
 1189  $\sup_v v^\top \psi - \frac{1}{2}(v^\top C_{\tilde{\theta}} v - \|\psi\|^2)$ .  
 1190

1191 Rewriting it  $\frac{1}{2}\|\psi\|^2 = \sup_v v^\top \psi - \frac{1}{2}v^\top C_{\tilde{\theta}} v$  and substituting  $u = 2v$

$$1192 \quad \|\psi\|^2 = \sup_u u^\top \psi - \frac{1}{4}u^\top C_{\tilde{\theta}} u.$$

1193 Using the change of variables  $u \rightarrow v$

$$1194 \quad \|\psi\|^2 = \sup_v v^\top \psi - \frac{1}{4}v^\top C_{\tilde{\theta}} v.$$

1195 Now, we want to find a function form for the optimization problem mentioned above.

1196 Consider a finite-dimensional functional spaces  $\mathcal{F}^i = \text{span}\{f_1^i, f_2^i, \dots, f_m^i\}$  for each client  $i$ . Hence,  
 1200 for  $f^i \in \mathcal{F}^i$

$$1201 \quad f^i = \sum_{j=1}^m v_j f_j^i.$$

1202 Since all the clients share the same neural network architecture, we define a global functional space  
 1204  $\mathcal{F}$  as

$$1205 \quad \mathcal{F} = \text{span}\{f_j^i \mid i \in [N], j \in [m]\}.$$

1206 Therefore,  $v$  corresponds to  $f^i$  such that

$$1207 \quad f^i = \sum_{c=1}^N \sum_{j=1}^m v_j^i f_j^c, \text{ where } v_j^i = \begin{cases} v_j & \text{if } c = i \\ 0 & \text{if } c \neq i \end{cases}$$

1208 Hence,

$$1209 \quad v^\top \psi = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m v_j \psi_{n_i}(f_j^i; \theta)$$

$$1210 \quad = \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z_k^i)(Y_k^i - g^i(X_k^i; \theta)).$$

1211 Similarly,

$$1212 \quad v^\top C_{\tilde{\theta}} v = \sum_{p=1}^m \sum_{q=1}^m v_p v_q [C_{\tilde{\theta}}]_{pq}$$

$$1213 \quad = \sum_{p=1}^m \sum_{q=1}^m v_p v_q \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{k=1}^{n_i} f_p^i(Z_k^i) f_q^i(Z_k^i) (Y_k^i - g^i(X_k^i; \tilde{\theta}))$$

$$1214 \quad = \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{k=1}^{n_i} \sum_{p=1}^m v_p f_p^i(Z_k^i) \sum_{q=1}^m v_q f_q^i(Z_k^i) (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2$$

$$1215 \quad = \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{k=1}^{n_i} (f^i(Z_k^i))^2 (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2$$

$$1216 \quad = \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\tilde{\theta}}(f^i, f^i).$$

1217 Thus, using the linear isomorphism between  $\mathbb{R}^m$  and  $\text{span}\{f_1^i, f_2^i, \dots, f_m^i\}$ , using  $v^\top \psi =$   
 1218  $\frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f^i; \theta)$  and  $v^\top C_{\tilde{\theta}} v = \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\tilde{\theta}}(f^i, f^i)$ , we can write the objective in functional  
 1219 form as

$$1220 \quad \|\psi\|^2 = \sup_{\substack{f^i \in \mathcal{F} \\ \forall i \in [N]}} \frac{1}{N} \sum_{i=1}^N \left( \psi_{n_i}(f^i; \theta) - \frac{1}{4} \mathcal{C}_{\tilde{\theta}}(f^i, f^i) \right).$$

1221 This gives us the desired result.

1222

□

## B.2 PROOF OF THEOREM 1

**Theorem 5** (Restatement of Theorem 1). *Under assumptions 1, 2, 3 and 4, a minimax solution  $(\hat{\theta}, \hat{\tau})$  of federated optimization problem (17) that satisfies the equilibrium condition as in definition 1:  $U_{\hat{\theta}}(\hat{\theta}, \tau) \leq U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\hat{\theta}}(\hat{\theta}, \tau')$ , is an  $\mathcal{E}$ -approximate federated equilibrium solution as defined in 3, where the approximation error  $\varepsilon^i$  for each client  $i \in [N]$  lies in:  $\max\{\zeta_{\hat{\theta}}^i, \zeta_{\hat{\tau}}^i\} \leq \varepsilon^i \leq \min\{\alpha - \rho_{\hat{\tau}}^i, \beta - B^i\}$  for  $\rho_{\hat{\tau}}^i < \alpha$  and  $B^i > \beta$ , such that  $\alpha := \left| \lambda_{\max} \left( \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \right) \right|$ ,  $\beta := \lambda_{\min} \left( \left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}} - \nabla_{\theta\tau}^2 U_{\hat{\theta}} \left( \nabla_{\tau\tau}^2 U_{\hat{\theta}} \right)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right] (\hat{\theta}, \hat{\tau}) \right)$  and  $B^i := \rho_{\hat{\theta}}^i + L\rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L\rho_{\tau\theta}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L^2\rho_{\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i) \cdot \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}})|}$ .*

*Proof.* The pure-strategy Stackelberg equilibrium for the federated objective is:

$$U_{\hat{\theta}}(\hat{\theta}, \tau) \leq U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\hat{\theta}}(\hat{\theta}, \tau'), \quad (40)$$

We want to show that the  $\varepsilon^i$ -approximate equilibrium for each client's objective  $U_{\hat{\theta}}^i$  also hold individually.

The first-order necessary condition for (40) to hold is  $\nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) = 0$  and  $\nabla_{\tau} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) = 0$ . Thus,  $\left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \right\|^2 = 0$ .

Consider

$$\begin{aligned} \left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \right\|^2 &= \left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) + \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 \\ &= \left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 + \left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 \\ &\quad + 2 \left( \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right)^{\top} \left( \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right) \end{aligned}$$

Rearranging

$$\begin{aligned} 2 \left( \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \right)^{\top} \left( \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right) - \left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 &= \left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 \\ 2 \left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 - \left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 &= \left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 \end{aligned}$$

Using gradient heterogeneity assumption (3) on R.H.S.

$$\left\| \nabla_{\theta} U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) - \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 \leq (\zeta_{\hat{\theta}}^i)^2$$

Thus, we obtain  $\left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\| \leq \zeta_{\hat{\theta}}^i$ . Similarly,  $\left\| \nabla_{\tau} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\| \leq \zeta_{\hat{\tau}}^i$ .

In the special case, when  $\zeta_{\hat{\theta}}^i = 0$  and  $\zeta_{\hat{\tau}}^i = 0$ , thus we will have  $\left\| \nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 = \left\| \nabla_{\tau} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \right\|^2 = 0$  for all  $i \in [N]$ , which gives  $\nabla_{\theta} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) = \nabla_{\tau} U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) = 0$  for all clients  $i$ .

Next, we prove that each client satisfies the second-order necessary condition approximately. Since  $(\hat{\theta}, \hat{\tau})$  satisfy the equilibrium condition (40), the second-order necessary condition holds for the global function  $U_{\hat{\theta}}$ , i.e.  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \preceq \mathbf{0}$ . We now prove that  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \preceq \mathbf{0}$ .

Using assumption 1, the hessian is symmetric. Thus,  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \preceq \mathbf{0}$  implies  $\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) \leq 0$ , where  $\lambda_{\max}$  is the largest eigenvalue of the hessian. Suppose,  $\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) = -\alpha$ , for some  $\alpha \geq 0$ .

We can write  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) = \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) + \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})$ .

Using a corollary of Weyl's theorem (Horn & Johnson, 2012) for real symmetric matrices  $A$  and  $B$ ,  $\lambda_{\max}(A + B) \leq \lambda_{\max}(A) + \lambda_{\max}(B)$ . Hence,

$$\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})) \leq \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) + \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})).$$



1296 Thus,  $\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})) \leq \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) - \alpha$ .

1297 Since the spectral norm of a real symmetric matrix  $A$  is given as  $\|A\|_{\sigma} =$   
 1298  $\max\{|\lambda_{\max}(A)|, |\lambda_{\min}(A)|\}$ .

1300 Under hessian heterogeneity assumption 4

$$\begin{aligned} 1301 \quad & \|\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})\|_{\sigma} = \max\{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\theta, \tau) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau))|, \\ 1302 \quad & |\lambda_{\min}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\theta, \tau) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau))|\} \\ 1303 \quad & \leq \rho_{\tau}^i. \end{aligned}$$

1306 Thus, we have

$$\begin{aligned} 1307 \quad & \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) \leq \max\left\{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}))|, \right. \\ 1308 \quad & \left. |\lambda_{\min}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}))|\right\} \\ 1309 \quad & \leq \rho_{\tau}^i. \end{aligned}$$

1313 Thus,  $\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})) \leq \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) - \nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})) - \alpha \leq \rho_{\tau}^i - \alpha$ , where  $\rho_{\tau}^i \geq 0$ .  
 1314 Hence,

$$\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \preceq (\rho_{\tau}^i - \alpha)\mathbf{I}.$$

1316 When  $\rho_{\tau}^i \leq \alpha$ , then  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \preceq 0$ .

1318 Now, since  $(\hat{\theta}, \hat{\tau})$  satisfy the equilibrium condition (40), thus  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \prec 0$  and the Schur  
 1319 complement of  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\hat{\theta}, \hat{\tau})$  is positive semi-definite. Now when  $\rho_{\tau}^i < \alpha$ , it follows from above that  
 1320  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}) \prec 0$ , hence  $(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau}))^{-1}$  exists. Now, we need to show that Schur complement of  
 1321  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i(\hat{\theta}, \hat{\tau})$  is positive semi-definite.

1324 Since,  $S(\hat{\theta}, \hat{\tau}) := \left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}} - \nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right] (\hat{\theta}, \hat{\tau}) \succ 0$ .

1326 Define  $S^i := \left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right]$ . We aim to prove  $\lambda_{\min}(S^i) \geq 0$  to show  $S^i$   
 1327 is positive semidefinite (PSD).

1329 Analogous to the above part, using corollary to Weyl's theorem, we have

$$\lambda_{\min}(S^i - S) + \lambda_{\min}(S) \leq \lambda_{\min}(S^i).$$

1332 Let  $\lambda_{\min}(S) = \beta$ , where  $\beta \geq 0$ . Moreover,  $\|S^i - S\|_{\sigma} = \max\{|\lambda_{\max}(S^i - S)|, |\lambda_{\min}(S^i - S)|\}$ ,  
 1333 thus  $\lambda_{\min}(S^i - S) \geq -\|S^i - S\|_{\sigma}$ .

1334 Thus, we have

$$-\|S^i - S\|_{\sigma} + \beta \leq \lambda_{\min}(S^i).$$

1336 We can write  $S^i - S$  as

$$\begin{aligned} 1337 \quad & S^i - S = (\nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\theta}^2 U_{\hat{\theta}}) - \left[ (\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}) (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right. \\ 1338 \quad & \left. + \nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} (\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i - \nabla_{\tau\theta}^2 U_{\hat{\theta}}) + \nabla_{\theta\tau}^2 U_{\hat{\theta}} \left( (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} - (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \right) \nabla_{\tau\theta}^2 U_{\hat{\theta}} \right]. \end{aligned}$$

1341 Hence,

$$\begin{aligned} 1342 \quad & \|S^i - S\|_{\sigma} \leq \|\nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\theta}^2 U_{\hat{\theta}}\|_{\sigma} + \underbrace{\|(\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}) (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma}}_{T_1} \\ 1343 \quad & + \underbrace{\|\nabla_{\theta\tau}^2 U_{\hat{\theta}} (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} (\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i - \nabla_{\tau\theta}^2 U_{\hat{\theta}})\|_{\sigma}}_{T_2} \\ 1344 \quad & + \underbrace{\|\nabla_{\theta\tau}^2 U_{\hat{\theta}} \left( (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} - (\nabla_{\tau\tau}^2 U_{\hat{\theta}})^{-1} \right) \nabla_{\tau\theta}^2 U_{\hat{\theta}}\|_{\sigma}}_{T_3}. \end{aligned}$$

Note that the eigenvalue of  $(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}$  is  $\lambda\left((\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\right) = \frac{1}{\lambda(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)}$ , hence  $\|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} = \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|}$  as  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i$  is negative definite. By Assumption 2, each client's function  $U^i$  is  $L$ -Lipschitz thus  $\|\nabla^2 U_{\hat{\theta}}^i\|_{\sigma} \leq L$ . Since the Hessian  $\nabla^2 U_{\hat{\theta}}^i$  is a block matrix of the form:

$$\nabla^2 U_{\hat{\theta}}^i = \begin{bmatrix} \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i & \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i \\ \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i & \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i \end{bmatrix},$$

The norm of Hessian is at least the norm of one of its components

$$\|\nabla_{\theta\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \leq L, \quad \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \leq L, \quad \|\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \leq L, \quad \|\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \leq L.$$

Thus, each Hessian block is individually bounded by  $L$ . Additionally,  $U$  is  $L$ -Lipschitz too. Using Assumption 4, bounding  $T_1$

$$\begin{aligned} T_1 &= \|(\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i)(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &\leq \|(\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i)\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &\leq L \rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} \end{aligned}$$

Similarly, bounding  $T_2$

$$\begin{aligned} T_2 &= \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} (\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i - \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i)\|_{\sigma} \\ &\leq \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|(\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i - \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i)\|_{\sigma} \\ &\leq L \rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} \end{aligned}$$

Lastly we bound  $T_3$ , it is easy to verify that  $\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}$

$$\begin{aligned} T_3 &= \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i \left( (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} - (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} \right) \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &\leq \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} - (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &= \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1} (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i - \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i) (\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &\leq \|\nabla_{\theta\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i - \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\theta}^2 U_{\hat{\theta}}^i\|_{\sigma} \\ &\leq L^2 \rho_{\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i) \cdot \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} \end{aligned}$$

Using bounds for  $T_1$ ,  $T_2$  and  $T_3$ , we can obtain a bound on  $\|S^i - S\|_{\sigma} \leq B^i$ , where  $B^i = \rho_{\hat{\theta}}^i + L \rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L \rho_{\tau\theta}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|} + L^2 \rho_{\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i) \cdot \lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)|}$ . Consider  $\rho^i = \max\{\rho_{\hat{\theta}}^i, \rho_{\tau\theta}^i, \rho_{\theta\tau}^i, \rho_{\tau}^i\}$ . Hence,  $B^i \leq \rho^i \left( 1 + \frac{L}{\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)} \left( 2 + \frac{1}{\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\hat{\theta}}^i)} \right) \right)$ . Hence, we obtain

$$\lambda_{\min}(S^i) \geq -B^i + \beta,$$

where  $\lambda_{\max}(S) = \beta$  such that  $\beta \geq 0$ . Hence, we obtain  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i \left( \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i \right)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right] (\hat{\theta}, \hat{\tau}) \succeq (\beta - B^i)I$ . When  $\beta \geq B^i$ , then  $S^i$  is positive semi-definite. When  $B^i = 0$ , hence  $\left[ \nabla_{\theta\theta}^2 U_{\hat{\theta}}^i - \nabla_{\theta\tau}^2 U_{\hat{\theta}}^i \left( \nabla_{\tau\tau}^2 U_{\hat{\theta}}^i \right)^{-1} \nabla_{\tau\theta}^2 U_{\hat{\theta}}^i \right] (\hat{\theta}, \hat{\tau}) \succeq \beta I$ , thus it will be positive semidefinite. When  $\rho_{\tau}^i < \alpha$  and  $\beta > B^i$ , then the sufficient condition for  $\varepsilon^i$ -approximate equilibrium is satisfied. And we obtain the result.

Thus, for each client  $i$ , any approximation error  $\varepsilon^i$  that satisfies:

$$\max\{\zeta_{\hat{\theta}}^i, \zeta_{\hat{\tau}}^i\} \leq \varepsilon^i \leq \min\{\alpha - \rho_{\tau}^i, \beta - B^i\}.$$

for  $\rho_{\tau}^i < \alpha$  and  $B^i > \beta$ , then  $(\hat{\theta}, \hat{\tau})$  is an  $\varepsilon^i$ -approximate local equilibrium point for client  $i$ .  $\square$

1404 B.3 CONSISTENCY

1405  
1406 B.3.1 ASSUMPTIONS  
1407

1408 We first state the assumptions that are necessary to establish the consistency of the estimated parameter.  
1409

1410 **Assumption 6** (Identification).  $\theta_0$  is the unique  $\theta \in \Theta$  such that  $\psi(f^i; \theta) = 0$  for all  $f^i \in \mathcal{F}$ , where  
1411  $i \in [n]$ .  
1412

1413 **Assumption 7** (Absolutely Star Shaped). For every  $f^i \in \mathcal{F}^i$  and  $|c| \leq 1$ , we have  $cf^i \in \mathcal{F}^i$ .  
1414

1415 **Assumption 8** (Continuity). For any  $x$ ,  $g^i(x; \theta)$ ,  $f^i(x; \tau)$  are continuous in  $\theta$  and  $\tau$ , respectively for  
1416 all  $i \in [N]$ .

1417 **Assumption 9** (Boundedness).  $Y^i$ ,  $\sup_{\theta \in \Theta} |g^i(X; \theta)|$ ,  $\sup_{\tau \in \mathcal{T}} |f^i(Z; \tau)|$  are bounded random  
1418 variables for all  $i \in [N]$ .  
1419

1420 **Assumption 10** (Bounded Complexity).  $\mathcal{F}^i$  and  $\mathcal{G}^i$  have bounded Rademacher complexities:  
1421

$$1422 \frac{1}{2^{n_i}} \sum_{\xi_i \in \{-1, +1\}^{n_i}} \mathbb{E} \sup_{\tau \in \mathcal{T}} \frac{1}{n_i} \sum_{k=1}^{n_i} \xi_i f^i(Z_k; \tau) \rightarrow 0, \quad \frac{1}{2^{n_i}} \sum_{\xi_i \in \{-1, +1\}^{n_i}} \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n_i} \sum_{k=1}^{n_i} \xi_i g^i(X_k; \theta) \rightarrow 0.$$

1425  
1426 B.3.2 PROOF OF THEOREM 2

1427 **Theorem 6** (Restatement of Theorem 2). Let  $\tilde{\theta}_n$  be a data-dependent choice for the federated objective that has a limit in probability. For each client  $i \in [N]$ , define  $m^i(\theta, \tau, \tilde{\theta}) :=$   
1428  $f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4} f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2$ ,  $M^i(\theta) = \sup_{\tau \in \mathcal{T}} \mathbb{E}[m^i(\theta, \tau, \tilde{\theta})]$  and  
1429  $\eta^i(\epsilon) := \inf_{d(\theta, \theta_0) \geq \epsilon} M^i(\theta) - M^i(\theta_0)$  for every  $\epsilon > 0$ . Let  $(\hat{\theta}_n, \hat{\tau}_n)$  be a solution that satisfies the  
1430 approximate equilibrium for each of the client  $i \in [N]$  as  
1431  
1432  
1433

$$1434 \sup_{\tau \in \mathcal{T}} U_{\tilde{\theta}}^i(\hat{\theta}_n, \tau) - \epsilon^i - o_p(1) \leq U_{\tilde{\theta}}^i(\hat{\theta}_n, \hat{\tau}_n) \leq \inf_{\theta \in \Theta} \max_{\tau: \|\tau - \hat{\tau}_n\| \leq h(\delta)} U_{\tilde{\theta}}^i(\theta, \tau) + \epsilon^i + o_p(1),$$

1435  
1436 for some  $\delta_0$ , such that for any  $\delta \in (0, \delta_0]$ , and any  $\theta, \tau$  such that  $\|\theta - \hat{\theta}\| \leq \delta$  and  $\|\tau - \hat{\tau}\| \leq \delta$  and  
1437 a function  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ . Then, under similar assumptions as in Assumptions 1 to 5 of (Bennett  
1438 et al., 2019), the global solution  $\hat{\theta}_n$  is a consistent estimator to the true parameter  $\theta_0$ , i.e.  $\hat{\theta}_n \xrightarrow{p} \theta_0$   
1439 when the approximate error  $\epsilon^i < \frac{\eta^i(\epsilon)}{2}$  for every  $\epsilon > 0$  for each client  $i \in [N]$ .  
1440  
1441  
1442  
1443

1444 *Proof.* The proof follows from the result of Bennett et al. (2019) that established the consistency of  
1445 the DEEPGMM estimator.  
1446

1447 First, we define the following terms for the ease of analysis:  
1448

$$1449 m^i(\theta, \tau, \tilde{\theta}) = f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4} f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2$$

$$1450 M^i(\theta) = \sup_{\tau \in \mathcal{T}} \mathbb{E}[m^i(\theta, \tau, \tilde{\theta})]$$

$$1451 M_{n_i}(\theta) = \sup_{\tau \in \mathcal{T}} \mathbb{E}_{n_i}[m^i(\theta, \tau, \tilde{\theta}_n)]$$

1452  
1453  
1454  
1455 Note that  $\tilde{\theta}_n$  is a data-dependent sequence for the global model. Practically, the previous global  
1456 iterate is used as  $\tilde{\theta}$ . Thus, we can define for the federated setting  $\tilde{\theta}_n = \frac{1}{N} \sum_{i=1}^N \tilde{\theta}_{n_i}$ . Let's assume  
1457  $\tilde{\theta}_n \xrightarrow{p} \tilde{\theta}$ .

**Claim 1:**  $\sup_{\theta} |M_{n_i}(\theta) - M^i(\theta)| \xrightarrow{P} 0$ .

$$\begin{aligned}
\sup_{\theta} |M_{n_i}(\theta) - M^i(\theta)| &= \sup_{\theta} \left| \sup_{\tau \in \mathcal{T}} \mathbb{E}_{n_i} [m^i(\theta, \tau, \tilde{\theta}_n)] - \sup_{\tau \in \mathcal{T}} \mathbb{E} [m^i(\theta, \tau, \tilde{\theta})] \right| \\
&\leq \sup_{\theta, \tau} \left| \mathbb{E}_{n_i} [m^i(\theta, \tau, \tilde{\theta}_n)] - \mathbb{E} [m^i(\theta, \tau, \tilde{\theta})] \right| \\
&\leq \sup_{\theta, \tau} \left| \mathbb{E}_{n_i} [m^i(\theta, \tau, \tilde{\theta}_n)] - \mathbb{E} [m^i(\theta, \tau, \tilde{\theta}_n)] \right| + \sup_{\theta, \tau} \left| \mathbb{E} [m^i(\theta, \tau, \tilde{\theta}_n)] - \mathbb{E} [m^i(\theta, \tau, \tilde{\theta})] \right| \\
&\leq \sup_{\theta_1, \theta_2, \tau} \left| \mathbb{E}_{n_i} [m^i(\theta_1, \tau, \theta_2)] - \mathbb{E} [m^i(\theta_1, \tau, \theta_2)] \right| + \sup_{\theta, \tau} \left| \mathbb{E} [m^i(\theta, \tau, \tilde{\theta}_n)] - \mathbb{E} [m^i(\theta, \tau, \tilde{\theta})] \right|
\end{aligned}$$

We will now handle the two terms in the above equation separately.

We will take the first term and call it  $B_1$ . For  $m^i(\theta, \tau, \tilde{\theta}_n)$ , we constitute its empirical counterpart  $m_k^i(\theta, \tau, \tilde{\theta}_n) = f^i(Z_k^i; \tau)(Y_k^i - g^i(X_k^i; \theta)) - \frac{1}{4} f^i(Z_k^i; \tau)^2 (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2$  and using  $m_k^i(\theta, \tau, \tilde{\theta}_n)$  with ghost variables  $\tilde{\theta}'_n$  for symmetrization and  $\epsilon_k$  as  $k$  i.i.d. Rademacher random variables, we obtain

$$\begin{aligned}
\mathbb{E}[B_1] &= \mathbb{E} \left[ \sup_{\theta_1, \theta_2, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} m_k^i(\theta_1, \tau, \theta_2) - \mathbb{E} [m_k^i(\theta_1, \tau, \theta_2)] \right| \right] \\
&\leq \mathbb{E} \left[ \sup_{\theta_1, \theta_2, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} (m_k^i(\theta_1, \tau, \theta_2) - m_k^i(\theta_1, \tau, \theta_2')) \right| \right] \\
&\leq \mathbb{E} \left[ \sup_{\theta_1, \theta_2, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k (m_k^i(\theta_1, \tau, \theta_2) - m_k^i(\theta_1, \tau, \theta_2')) \right| \right] \\
&\leq 2\mathbb{E} \left[ \sup_{\theta_1, \theta_2, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k m_k^i(\theta_1, \tau, \theta_2) \right| \right] \\
&\leq 2\mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k f^i(Z_k^i; \tau)(Y_k^i - g^i(X_k^i; \theta)) \right| \right] \\
&\quad + \frac{1}{2} \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k f^i(Z_k^i; \tau)^2 (Y_k^i - g^i(X_k^i; \tilde{\theta}))^2 \right| \right] \\
&\leq 2\mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k \left( \frac{1}{2} f^i(Z_k^i; \tau)^2 + \frac{1}{2} (Y_k^i - g^i(X_k^i; \theta))^2 \right) \right| \right] \\
&\quad + \frac{1}{2} \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k \left( \frac{1}{2} f^i(Z_k^i; \tau)^4 + \frac{1}{2} (Y_k^i - g^i(X_k^i; \tilde{\theta}))^4 \right) \right| \right] \\
&\leq \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k f^i(Z_k^i; \tau)^2 \right| \right] + \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k (Y_k^i - g^i(X_k^i; \theta))^2 \right| \right] \\
&\quad + \frac{1}{4} \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k f^i(Z_k^i; \tau)^4 \right| \right] + \frac{1}{4} \mathbb{E} \left[ \sup_{\theta, \tau} \left| \frac{1}{n_i} \sum_{k=1}^{n_i} \epsilon_k (Y_k^i - g^i(X_k^i; \tilde{\theta}))^4 \right| \right]
\end{aligned}$$

Using boundedness assumption 9, we consider the mapping from  $f^i(Z_k^i; \tau)$  and  $g^i(X_k^i; \tilde{\theta})$  to the summation terms in the last inequality as Lipschitz functions, hence for any functional class  $\mathcal{F}^i$  and  $L$ -Lipschitz function  $\phi$ ,  $\mathcal{R}_{n_i}(\phi \circ f^i) \leq L\mathcal{R}_{n_i}(\mathcal{F}^i)$ , where  $\mathcal{R}_{n_i}(\mathcal{F}^i)$  is the Rademacher complexity of class  $\mathcal{F}^i$ . Hence,  $\mathbb{E}[B_1] \leq L(\mathcal{R}_{n_i}(\mathcal{G}^i) + \mathcal{R}_{n_i}(\mathcal{F}^i))$ . Using assumption 10,  $\mathbb{E}[B_1] \rightarrow 0$ . Let  $B'_1$  be a modified value of  $B$ , after changing the  $j$ -th value of  $X^i, Z^i$  and  $Y^i$  values, using assumption 9 on

boundedness, we obtain the bounded difference inequality:

$$\sup_{X_{1:n_i}, Z_{1:n_i}, Y_{1:n_i}, X'_j, Z'_j, Y'_j} |B_1 - B'_1| \leq \sup_{\theta_1, \theta_2, \tau, X_{1:n_i}, Z_{1:n_i}, Y_{1:n_i}, X'_j, Z'_j, Y'_j} \left| \frac{1}{n_i} (m_j^i(\theta_1, \tau, \theta_2) - m_j^{i'}(\theta_1, \tau, \theta_2)) \right| \leq \frac{b}{n_i},$$

where  $b$  is some constant. Using McDiarmid's Inequality, we have  $P(|B_1 - \mathbb{E}[B_1]| \geq \epsilon_0) \leq 2 \exp\left(\frac{-2n_i \epsilon_0^2}{c^2}\right)$ . And  $\mathbb{E}[B_1] \rightarrow 0$ , we have  $B_1 \xrightarrow{p} 0$ .

Now, we will handle  $B_2$ . For that

$$\begin{aligned} B_2 &= \sup_{\theta, \tau} \left| \mathbb{E} \left[ m^i(\theta, \tau, \tilde{\theta}_n) \right] - \mathbb{E} \left[ m^i(\theta, \tau, \tilde{\theta}) \right] \right| \\ &= \sup_{\theta, \tau} \left| \mathbb{E} \left[ f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4} f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}_n))^2 \right] \right. \\ &\quad \left. - \mathbb{E} \left[ f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4} f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \right| \\ &= \sup_{\theta, \tau} \frac{1}{4} \left| \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}_n))^2 \right] - \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \right| \\ &= \sup_{\theta, \tau} \frac{1}{4} \left| \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}_n))^2 \right] + \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \right. \\ &\quad \left. - \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] - \mathbb{E} \left[ f^i(Z^i; \tau)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \right| \\ &\leq \frac{1}{4} \sup_{\tau} |\mathbb{E} [f^i(Z^i; \tau)^2 \omega_n]| \end{aligned}$$

Here,  $\omega_n = \left| (Y^i - g(X^i; \tilde{\theta}_n))^2 - (Y^i - g(X^i; \tilde{\theta}))^2 \right|$ . Due to our assumption,  $\tilde{\theta}_n \xrightarrow{p} \tilde{\theta}$ , thus  $\omega_n \xrightarrow{p} 0$  due to Slutsky's and continuous mapping theorem. Since,  $f^i(Z; \tau)$  is uniformly bounded, thus for some constant  $b' > 0$ , we have

$$\begin{aligned} B_2 &\leq \frac{b'}{4} \sup_{\tau} \frac{1}{N} \sum_{i=1}^N |\mathbb{E}[\omega_n]| \\ &\leq \frac{b'}{4} \sup_{\tau} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[|\omega_n|] \end{aligned}$$

Based on the boundedness assumption, we can verify that  $\omega_n$  is bounded, hence using Lebesgue Dominated Convergence Theorem, we can conclude that  $\mathbb{E}[|\omega_n|] \rightarrow 0$ .

Thus, using the convergence of  $B_1$  and  $B_2$ , we have  $\sup_{\theta} |M_{n_i}(\theta) - M^i(\theta)| \xrightarrow{p} 0$  for each  $i \in [N]$ .

**Claim 2:** for every  $\epsilon > 0$ , we have  $\inf_{d(\theta, \theta_0) \geq \epsilon} M^i(\theta) > M^i(\theta_0)$ .

$M^i(\theta_0)$  is the unique minimizer of  $M^i(\theta)$ . By assumption (6) and (7),  $\theta_0$  is the unique minimizer of  $\sup_{\tau} \mathbb{E}[f^i(Z^i; \tau)(Y^i - g^i(X; \theta))]$  such that  $\sup_{\tau} \mathbb{E}[f^i(Z^i; \tau)(Y^i - g^i(X; \theta))] = 0$ . Thus, any other value of  $\theta$  will have at least one  $\tau$  such that this expectation is strictly positive.  $M(\theta_0) = 0$  and  $M(\theta_0) = \sup_{\tau} -\frac{1}{4} f^i(Z^i; \tau)^2 (Y^i - g^i(X; \theta))^2$ , the function whose supremum is being evaluated is non-positive but can be set to zero by assumption (7) by taking the zero function of  $f^i$ . Let for any other  $\theta' \neq \theta_0$ , let  $f^{i'}$  be a function in  $\mathcal{F}^i$  such that  $\mathbb{E}[f^i(Z)(Y^i - g^i(X; \theta'))] > 0$ . If we have  $\mathbb{E}[f^{i'}(Z)^2 (Y^i - g^i(X; \tilde{\theta}))^2] = 0$ , then  $M^i(\theta') > 0$ . Else, consider  $c f^{i'}$  for any  $c \in (0, 1)$ . Using assumption (7),  $c f^{i'} \in \mathcal{F}^i$ , thus

$$\begin{aligned} M^i(\theta') &= \sup_{f^i \in \mathcal{F}^i} \mathbb{E} \left[ f^i(Z^i)(Y^i - g(X^i; \theta')) - \frac{1}{4} f^i(Z^i)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \\ &\leq c \mathbb{E} \left[ f^{i'}(Z^i)(Y^i - g(X^i; \theta')) \right] - \frac{c^2}{4} \mathbb{E} \left[ f^{i'}(Z^i)^2 (Y^i - g(X^i; \tilde{\theta}))^2 \right] \end{aligned}$$

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

This is quadratic in  $c$  and is positive when  $c$  is sufficiently small, thus  $M^i(\theta') > 0$ .

We now prove claim 2 using contradiction. Let us assume claim 2 is false, i.e. for some  $\epsilon > 0$ , we have  $\inf_{\theta \in B(\theta_0, \epsilon)} M^i(\theta) = M^i(\theta_0)$ , where  $B(\theta_0, \epsilon)^c = \{\theta \mid d(\theta, \theta_0) \geq \epsilon\}$ , since  $\theta_0$  is the unique minimizer of  $M^i(\theta)$  by assumption (6). Thus, there must exist some sequence  $(\theta_1, \theta_2, \dots)$  in  $B(\theta_0, \epsilon)^c$  such that  $M^i(\theta_n) \rightarrow M^i(\theta_0)$ . By construction,  $B(\theta_0, \epsilon)^c$  is closed and the corresponding limit parameters  $\theta^* = \lim_{n \rightarrow \infty} \theta_n \in B(\theta_0, \epsilon)^c$  must satisfy  $M^i(\theta^*) = M^i(\theta_0)$  using assumption (8). But  $d(\theta^*, \theta_0) \geq \epsilon > 0$ , thus  $\theta^* \neq \theta_0$ . This contradicts that  $\theta_0$  is the unique minimizer of  $M^i(\theta)$ ; hence, claim 2 is true.

**Claim 3:** For the third part, we know that  $\hat{\theta}_n$  satisfies the  $\varepsilon^i$ - approximate equilibrium condition, given as:

$$\mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \tau, \tilde{\theta}_n)] - \varepsilon^i \leq \mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \hat{\tau}_n, \tilde{\theta}_n)] \leq \max_{\tau': \|\tau' - \hat{\tau}_n\| \leq h(\delta)} \mathbb{E}_{n_i}[m^i(\theta, \tau', \tilde{\theta}_n)] + \varepsilon^i,$$

for a function  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$  and some  $\delta_0$ , such that for any  $\delta \in (0, \delta_0]$ , and any  $\theta, \tau$  such that  $\|\theta - \hat{\theta}\| \leq \delta$  and  $\|\tau - \hat{\tau}\| \leq \delta$ . Assume that this is true with  $o_p(1)$ , hence

$$\sup_{\tau} \mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \tau, \tilde{\theta}_n)] - \varepsilon^i - o_p(1) \leq \mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \hat{\tau}_n, \tilde{\theta}_n)] \leq \inf_{\theta} \max_{\tau': \|\tau' - \hat{\tau}_n\| \leq h(\delta)} \mathbb{E}_{n_i}[m^i(\theta, \tau', \tilde{\theta}_n)] + \varepsilon^i + o_p(1),$$

Now, since  $M_{n_i}(\hat{\theta}_n) = \sup_{\tau} \mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \tau, \tilde{\theta}_n)]$ . Hence,

$$\inf_{\theta} \max_{\tau': \|\tau' - \hat{\tau}_n\| \leq h(\delta)} \mathbb{E}_{n_i}[m^i(\theta, \tau', \tilde{\theta}_n)] \leq \inf_{\theta} \sup_{\tau} \mathbb{E}_{n_i}[m^i(\theta, \tau, \tilde{\theta}_n)] = \inf_{\theta} M_{n_i}(\theta) \leq M_{n_i}(\theta_0)$$

Thus, we have

$$M_{n_i}(\hat{\theta}_n) - \varepsilon^i - o_p(1) \leq \mathbb{E}_{n_i}[m^i(\hat{\theta}_n, \hat{\tau}_n, \tilde{\theta}_n)] \leq M_{n_i}(\theta_0) + \varepsilon^i + o_p(1).$$

We have proven all three conditions until now. From the first and second condition, since  $|M_{n_i}(\theta_0) - M^i(\theta_0)| \xrightarrow{p} 0$ , hence  $M_{n_i}(\hat{\theta}_n) \leq M^i(\theta_0) + 2\varepsilon^i + o_p(1)$ . Hence, we obtain

$$\begin{aligned} M^i(\hat{\theta}_n) - M^i(\theta_0) &\leq M^i(\hat{\theta}_n) - M_{n_i}(\hat{\theta}_n) + 2\varepsilon^i + o_p(1) \\ &\leq \sup_{\theta} |M^i(\hat{\theta}) - M_{n_i}(\hat{\theta})| + 2\varepsilon^i + o_p(1) \\ &\leq 2\varepsilon^i + o_p(1) \end{aligned}$$

Hence, we obtain

$$\begin{aligned} M^i(\hat{\theta}_n) - M^i(\theta_0) - 2\varepsilon^i &\leq M^i(\hat{\theta}_n) - M_{n_i}(\hat{\theta}_n) + o_p(1) \\ &\leq \sup_{\theta} |M^i(\hat{\theta}) - M_{n_i}(\hat{\theta})| + o_p(1) \\ &\leq o_p(1) \end{aligned}$$

Since, let  $\eta^i(\epsilon) := \inf_{d(\theta, \theta_0) \geq \epsilon} M^i(\theta) - M^i(\theta_0)$ . Hence, whenever  $d(\hat{\theta}_n, \theta_0) \geq \epsilon$ , we have  $M^i(\hat{\theta}_n) - M^i(\theta_0) \geq \eta^i(\epsilon)$ . Thus,  $\mathbb{P}[d(\hat{\theta}_n, \theta_0) \geq \epsilon] \leq \mathbb{P}[M^i(\hat{\theta}_n) - M^i(\theta_0) \geq \eta^i(\epsilon)] = \mathbb{P}[M^i(\hat{\theta}_n) - M^i(\theta_0) - 2\varepsilon^i \geq \eta^i(\epsilon) - 2\varepsilon^i]$ . For every  $\epsilon > 0$ , we have  $\eta^i(\epsilon) > 0$  from claim 2, and  $M^i(\hat{\theta}_n) - M^i(\theta_0) - 2\varepsilon^i = o_p(1)$ . Thus,  $\eta^i(\epsilon) - 2\varepsilon^i > 0$  when  $\varepsilon^i < \frac{\eta^i(\epsilon)}{2}$ . We have that for every  $\epsilon > 0$  and  $\varepsilon^i < \frac{\eta^i(\epsilon)}{2}$ , the RHS probability converges to 0, thus  $d(\hat{\theta}_n, \theta_0) = o_p(1)$ , hence  $\hat{\theta}_n$  converges in probability to  $\theta_0$  for each client  $i \in [N]$ . □

## C LIMIT POINTS OF FEDGDA

We first discuss the  $\gamma$ - FEDGDA flow.

---

## C.1 FEDGDA FLOW

The FEDGDA updates can be written as

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \frac{1}{\gamma} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R (\nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t) + (\nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t)) \\ &\quad + (\nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t))) \\ \tau_{t+1} &= \tau_t + \eta \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R (\nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t) + (\nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta_t, \tau_t)) \\ &\quad + (\nabla_{\tau} U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t)))\end{aligned}$$

Rearranging the terms and taking the continuous-time limit as  $\eta \rightarrow 0$

$$\begin{aligned}\lim_{\eta \rightarrow 0} \frac{\theta_{t+1} - \theta_t}{\eta} &= \lim_{\eta \rightarrow 0} -\frac{1}{\gamma} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R (\nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t) + (\nabla_{\theta} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t)) \\ &\quad + (\nabla_{\theta} U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_{\theta} U_{\bar{\theta}}(\theta_t, \tau_t))) \\ \lim_{\eta \rightarrow 0} \frac{\tau_{t+1} - \tau_t}{\eta} &= \lim_{\eta \rightarrow 0} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^R (\nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t) + (\nabla_{\tau} U_{\bar{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta_t, \tau_t)) \\ &\quad + (\nabla_{\tau} U_{\bar{\theta}}^i(\theta_t, \tau_t) - \nabla_{\tau} U_{\bar{\theta}}(\theta_t, \tau_t)))\end{aligned}$$

We obtain the gradient flow equations as

$$\begin{aligned}\frac{d\theta}{dt} &= -\frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}(\theta(t), \tau(t))) - \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta(t), \tau(t))) \\ &\quad - \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\theta} U_{\bar{\theta}}(\theta(t), \tau(t))),\end{aligned}\tag{41}$$

$$\begin{aligned}\frac{d\tau}{dt} &= R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}(\theta(t), \tau(t))) + R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta(t), \tau(t))) \\ &\quad + R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\tau} U_{\bar{\theta}}(\theta(t), \tau(t))).\end{aligned}\tag{42}$$

Using Assumption 3

$$\begin{aligned}\left\| \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\theta} U_{\bar{\theta}}(\theta(t), \tau(t))) \right\| &\leq \frac{R}{\gamma} \zeta_{\theta} \\ \left\| R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\tau} U_{\bar{\theta}}(\theta(t), \tau(t))) \right\| &\leq R \zeta_{\tau}\end{aligned}$$

Thus,

$$\begin{aligned}\frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\theta} U_{\bar{\theta}}(\theta(t), \tau(t))) &= \mathcal{O}\left(\frac{R}{\gamma} \zeta_{\theta}\right) \\ R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta(t), \tau(t)) - \nabla_{\tau} U_{\bar{\theta}}(\theta(t), \tau(t))) &= \mathcal{O}(R \zeta_{\tau})\end{aligned}$$

Since  $U_{\bar{\theta}}^i$  is Lipschitz smooth by assumption 2, we have

$$\begin{aligned} \left\| \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta(t), \tau(t))) \right\| &\leq L \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \|(\theta^i(t), \tau^i(t)) - (\theta(t), \tau(t))\|, \\ \left\| R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta(t), \tau(t))) \right\| &\leq LR \frac{1}{N} \sum_{i \in [N]} \|(\theta^i(t), \tau^i(t)) - (\theta(t), \tau(t))\|. \end{aligned}$$

Substituting these bounds into Equations (41) and (42), we obtain

$$\begin{aligned} \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} (\nabla_{\theta} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\theta} U_{\bar{\theta}}^i(\theta, \tau)) &= \mathcal{O} \left( L \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \|(\theta^i(t), \tau^i(t)) - (\theta(t), \tau(t))\| \right), \\ R \frac{1}{N} \sum_{i \in [N]} (\nabla_{\tau} U_{\bar{\theta}}^i(\theta^i(t), \tau^i(t)) - \nabla_{\tau} U_{\bar{\theta}}^i(\theta, \tau)) &= \mathcal{O} \left( LR \frac{1}{N} \sum_{i \in [N]} \|(\theta^i(t), \tau^i(t)) - (\theta(t), \tau(t))\| \right). \end{aligned}$$

Since the local update follows

$$\begin{aligned} \theta^i(t) &= \theta(t) - \frac{\eta}{\gamma} \sum_{j=1}^R \nabla_{\theta} U_{\bar{\theta}}^i(\theta_j^i(t), \tau_j^i(t)), \\ \tau^i(t) &= \tau(t) + \eta \sum_{j=1}^R \nabla_{\tau} U_{\bar{\theta}}^i(\theta_j^i(t), \tau_j^i(t)), \end{aligned}$$

Using bounded gradient assumption, i.e.  $\|\nabla_{\theta} U_{\bar{\theta}}^i(\theta, \tau)\| \leq G_{\theta}$  and  $\|\nabla_{\tau} U_{\bar{\theta}}^i(\theta, \tau)\| \leq G_{\tau}$  for all  $i$ , as  $\eta \rightarrow 0$  and  $R$  is fixed and finite, the deviation  $\|(\theta^i(t), \tau^i(t)) - (\theta(t), \tau(t))\|$  vanish, leading to

$$\begin{aligned} \frac{d\theta}{dt} &= -\frac{1}{\gamma} R \nabla_{\theta} U_{\bar{\theta}}(\theta(t), \tau(t)) + \mathcal{O} \left( \frac{R}{\gamma} \zeta_{\theta} \right), \\ \frac{d\tau}{dt} &= R \nabla_{\tau} U_{\bar{\theta}}(\theta(t), \tau(t)) + \mathcal{O}(R \zeta_{\tau}). \end{aligned}$$

### C.1.1 PROOF OF PROPOSITION 1

**Proposition** (Restatement of Proposition 1). *Given the Jacobian matrix for  $\gamma$ -FEDGDA flow as*

$$\mathbf{J} = \begin{pmatrix} -\frac{1}{\gamma} R \nabla_{\theta\theta}^2 U_{\bar{\theta}}(\theta, \tau) & -\frac{1}{\gamma} R \nabla_{\theta\tau}^2 U_{\bar{\theta}}(\theta, \tau) \\ R \nabla_{\tau\theta}^2 U_{\bar{\theta}}(\theta, \tau) & R \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) \end{pmatrix},$$

*a point  $(\theta, \tau)$  is a strictly linearly stable equilibrium of the  $\gamma$ -FEDGDA flow if and only if the real parts of all eigenvalues of  $\mathbf{J}$  are negative, i.e.,  $\text{Re}(\Lambda_j) < 0$  for all  $j$ .*

*Proof.* Considering the FEDGDA dynamics with step size  $\eta$ , the Jacobian matrix of this dynamic system is  $\mathbf{I} + \eta \mathbf{J}$ . The eigenvalues of  $\mathbf{J}$  are  $\Lambda_j$ , thus the eigenvalues of  $\mathbf{I} + \eta \mathbf{J}$  are  $\{1 + \eta \Lambda_j\}$ .

By definition, a fixed point  $\mathbf{z}^*$  of a dynamical system  $\mathbf{w}$ , such that  $\mathbf{z}^* = \mathbf{w}(\mathbf{z}^*)$ , is a strict linearly stable point if the spectral radius  $\rho(\mathbf{J}(\mathbf{z}^*)) < 1$ , where  $\mathbf{J}$  is the Jacobian matrix of  $\mathbf{w}$ . Therefore,  $(\theta, \tau)$  is a strict linearly stable point if and only if  $\rho(\mathbf{I} + \eta \mathbf{J}) < 1$ , that is  $|1 + \eta \Lambda_j| < 1$  for all  $j$ . When taking  $\eta \rightarrow 0$ , this is equivalent to  $\text{Re}(\Lambda_j) < 0$  for all  $j$ .  $\square$



C.2 PROOF OF THEOREM 3

*Proof.* Let  $\mathbf{A} = \nabla_{\theta\theta}^2 U_{\bar{\theta}}(\theta, \tau)$ ,  $\mathbf{B} = \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)$  and  $\mathbf{C} = \nabla_{\theta\tau}^2 U_{\bar{\theta}}(\theta, \tau)$ . Consider  $\epsilon = \frac{1}{\gamma}$ , thus for sufficiently small  $\epsilon$  (hence a large  $\gamma$ ), the Jacobian  $\mathbf{J}$  of FEDGDA for a point  $(\theta, \tau)$  is given as:

$$\mathbf{J}_\epsilon = R \begin{pmatrix} -\epsilon\mathbf{A} & -\epsilon\mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}.$$

Using Lemma 9,  $\mathbf{J}_\epsilon$  has  $d_1 + d_2$  complex eigenvalues  $\{\Lambda_j\}_{j=1}^{d_1+d_2}$  such that

$$\begin{aligned} |\Lambda_j + \epsilon\mu_j| &= o(\epsilon) & 1 \leq j \leq d_1 \\ |\Lambda_{j+d_1} - \nu_j| &= o(1), & 1 \leq j \leq d_2, \end{aligned} \quad (43)$$

where  $\{\mu_j\}_{j=1}^{d_1}$  and  $\{\nu_j\}_{j=1}^{d_2}$  are the eigenvalues of matrices  $R(\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top)$  and  $R\mathbf{B}$  respectively.

We now prove the theorem statement:

$$\text{LocMinimax} \subset \overline{\infty - \text{FGDA}} \subset \overline{\infty - \text{FGDA}} \subset \text{LocMinimax} \cup \{(\theta, \tau) | (\theta, \tau) \text{ is stationary and } \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) \text{ is degenerate}\}.$$

By definition of lim sup and lim inf, we know that  $\overline{\infty - \text{FGDA}} \subset \overline{\infty - \text{FGDA}}$ .

Now we show  $\text{LocMinimax} \subset \overline{\infty - \text{FGDA}}$ . Consider a strict local minimax point  $(\theta, \tau)$ , then by sufficient condition it follows that:

$$\mathbf{B} \prec 0, \quad \text{and} \quad \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \succ 0.$$

Thus,  $R\mathbf{B} \prec 0$ , and  $R(\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top) \succ 0$ , where  $R$  is always positive. Hence,  $\{\nu_j\}_{j=1}^{d_2} < 0$  and  $\{\mu_j\}_{j=1}^{d_1} < 0$ . Using equations 43, for some small  $\epsilon_0 < \epsilon$ ,  $\text{Re}(\Lambda_j) < 0$  for all  $j$ . Thus,  $(\theta, \tau)$  is a strict linearly stable point of  $\frac{1}{\epsilon}$ -FEDGDA.

Now, we show  $\overline{\infty - \text{FGDA}} \subset \text{LocMinimax} \cup \{(\theta, \tau) | (\theta, \tau) \text{ is stationary and } \nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) \text{ is degenerate}\}$ . Consider  $(\theta, \tau)$  a strict linearly stable point of  $\frac{1}{\epsilon}$ -FEDGDA, such that for some small  $\epsilon$ ,  $\text{Re}(\Lambda_j) < 0$  for all  $j$ . By equation 43, assuming  $\mathbf{B}^{-1}$  exists

$$R\mathbf{B} \prec 0, \quad \text{and} \quad R(\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top) \succeq 0.$$

Since,  $R$  is positive, thus  $\mathbf{B} \prec 0$ , and  $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \succeq 0$ . Let's assume  $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$  has 0 as an eigenvalue. Thus, there exists a unit eigenvector  $\mathbf{w}$  such that  $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \mathbf{w} = 0$ . Then,

$$\mathbf{J}_\epsilon \cdot (\mathbf{w}, -\mathbf{B}^{-1}\mathbf{C}^\top \mathbf{w})^\top = R \begin{pmatrix} -\epsilon\mathbf{A} & -\epsilon\mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{w} \\ -\mathbf{B}^{-1}\mathbf{C}^\top \mathbf{w} \end{pmatrix} = \mathbf{0}.$$

Thus,  $\mathbf{J}_\epsilon$  has 0 as its eigenvalue, which is a contradiction because for strict linearly stable point  $\text{Re}(\Lambda_j) < 0$  for all  $j$ . Thus,  $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \succ 0$ . Hence,  $(\theta, \tau)$  is a strict local minimax point.

Let  $G : \mathbb{R}^d \times \mathbb{R}^k \rightarrow \mathbb{R}$  be the function defined as:  $G(\theta, \tau) = \det(\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau))$ . Let's assume that  $\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)$  is smooth, thus the determinant function is a polynomial in the entries of the Hessian, which implies that  $G$  is a smooth function. Since  $\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau) = 0$  implies at least one eigenvalue of  $\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)$  is zero, thus  $\det(\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)) = 0$ .

We aim to show that the set

$$\mathcal{A} = \{(\theta, \tau) \mid (\theta, \tau) \text{ is stationary and } \det(\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)) = 0\}$$

has measure zero in  $\mathbb{R}^d \times \mathbb{R}^k$ .

A point  $q \in \mathbb{R}^d \times \mathbb{R}^k$  is a *regular value* of  $G$  if for every  $(\theta, \tau) \in G^{-1}(q)$ , the differential  $dG(\theta, \tau)$  is surjective. Otherwise,  $q$  is a *critical value*.

The differential of  $G$  is given by:  $\nabla G(\theta, \tau) = \text{Tr}(\text{Adj}(\nabla_{\tau\tau}^2 U_{\bar{\theta}}) \cdot \nabla(\nabla_{\tau\tau}^2 U_{\bar{\theta}}))$ . If  $\det(\nabla_{\tau\tau}^2 U_{\bar{\theta}}(\theta, \tau)) = 0$ , then the Hessian  $\nabla_{\tau\tau}^2 U_{\bar{\theta}}$  is singular. This causes its adjugate matrix to lose rank, leading to a degeneracy in  $\nabla G(\theta, \tau)$ , making  $dG(\theta, \tau)$  *not surjective*.

Thus, every  $(\theta, \tau)$  satisfying  $G(\theta, \tau) = 0$  is a critical point of  $G$ , meaning that 0 is a *critical value* of  $G$ .

By Sard’s theorem, the set of critical values of a smooth function has measure zero in the codomain. Since  $G$  is smooth, the set of critical values of  $G$  in  $\mathbb{R}$  has measure zero. In particular, since 0 is a critical value of  $G$ , the set:  $G^{-1}(0) = \{(\theta, \tau) \mid \det(\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau)) = 0\}$  has measure zero in  $\mathbb{R}^{d+k}$ .

Since the set of degenerate  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau)$  is precisely  $G^{-1}(0)$ , we conclude that Lebesgue measure( $\mathcal{A}$ ) = 0. Thus, the set of stationary points where the Hessian  $\nabla_{\tau\tau}^2 U_{\hat{\theta}}(\theta, \tau)$  is singular has measure zero in  $\mathbb{R}^d \times \mathbb{R}^k$ .  $\square$

**Lemma 8.** (Zedek, 1965) *Given a polynomial  $p_n(z) := \sum_{k=0}^n a_k z^k$ , where  $a_n \neq 0$ , an integer  $m \geq n$  and a number  $\epsilon > 0$ , there exists a number  $\delta > 0$  such that whenever the  $m + 1$  complex numbers  $b_k, 0 \leq k \leq m$ , satisfy the inequalities*

$$|b_k - a_k| < \delta \quad \text{for } 0 \leq k \leq n, \quad \text{and} \quad |b_k| < \delta \quad \text{for } n + 1 \leq k \leq m,$$

*then the roots  $\beta_k, 1 \leq k \leq m$ , of the polynomial  $q_m(z) := \sum_{k=0}^m b_k z^k$  can be labeled in such a way as to satisfy, with respect to the zeros  $\alpha_k, 1 \leq k \leq n$ , of  $p_n(z)$ , the inequalities*

$$|\beta_k - \alpha_k| < \epsilon \quad \text{for } 1 \leq k \leq n, \quad \text{and} \quad |\beta_k| > 1/\epsilon \quad \text{for } n + 1 \leq k \leq m.$$

**Lemma 9.** *For any symmetric matrix  $\mathbf{A} \in \mathbb{R}^{d_1 \times d_1}$ ,  $\mathbf{B} \in \mathbb{R}^{d_2 \times d_2}$ , any rectangular matrix  $\mathbf{C} \in \mathbb{R}^{d_1 \times d_2}$  and a scalar  $R$ , assume that  $\mathbf{B}$  is non-degenerate. Then, matrix*

$$R \begin{pmatrix} -\epsilon \mathbf{A} & -\epsilon \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}$$

*has  $d_1 + d_2$  complex eigenvalues  $\{\Lambda_j\}_{j=1}^{d_1+d_2}$  with following form for sufficiently small  $\epsilon$ :*

$$\begin{aligned} |\Lambda_j + \epsilon \mu_j| &= o(\epsilon) & 1 \leq j \leq d_1 \\ |\Lambda_{j+d_1} - \nu_j| &= o(1), & 1 \leq j \leq d_2, \end{aligned}$$

*where  $\{\frac{1}{R}\mu_j\}_{j=1}^{d_1}$  and  $\{\frac{1}{R}\nu_j\}_{j=1}^{d_2}$  are the eigenvalues of matrices  $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$  and  $\mathbf{B}$  respectively.*

The proof follows from Lemma 8 by a similar argument as in (Jin et al., 2020) with  $\{\mu_j\}_{j=1}^{d_1}$  and  $\{\nu_j\}_{j=1}^{d_2}$  as the eigenvalues of matrices  $R(\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top)$  and  $R\mathbf{B}$ , respectively, and is thus omitted.

## D RELATED WORK

The federated supervised learning has received algorithmic advancements guided by factors such as tackling the system and statistical heterogeneities, better sample and communication complexities, model personalization, differential privacy, etc. An incomplete list includes FEDPROX (Li et al., 2020), SCAFFOLD (Karimireddy et al., 2020), FEDOPT (Reddi et al., 2020), LPP-SGD (Chatterjee et al., 2024), PFEDME (T Dinh et al., 2020), DP-SCAFFOLD (Noble et al., 2022), and others.

By contrast, federated learning with confounders in a causal learning setting is a relatively under-explored research area. Vo et al. (2022a) presented a method to learn the similarities among the data sources translating a structural causal model (Pearl, 2009) to federated setting. They transform the loss function by utilizing Random Fourier Features into components associated with the clients. Thereby they compute individual treatment effects (ITE) and average treatment effects (ATE) by a federated maximization of evidence lower bound (ELBO). Vo et al. (2022b) presented another federated Bayesian method to estimate the posterior distributions of the ITE and ATE using a non-parametric approach.

Xiong et al. (2023) presented maximum likelihood estimator (MLE) computation in a federated setting for ATE estimation. They showed that the federated MLE consistently estimates the ATE parameters considering the combined data across clients. However, it is not clear if this approach is applicable to consistent local moment conditions estimation for the participating clients. Almodóvar et al. (2024) applied FedAvg to variational autoencoder (Kingma et al., 2019) based treatment effect estimation TEDVAE (Zhang et al., 2021). However, their work mainly focused on comparing the

1836 performance of vanilla FedAvg with a propensity score-weighted FedAvg in the context of federated  
1837 implementation of TEDVAE.

1838  
1839 Our work differs from the above related works in the following:

- 1840 (a) we introduce IV analysis in federated setting, and, we introduce federated GMM estimators,  
1841 which has applications for various empirical research (Wooldridge, 2001),
- 1842 (b) specifically, we adopt a non-Bayesian approach based on a federated zero-sum game, wherein  
1843 we focus on analysing the dynamics of the federated minimax optimization and characterize the  
1844 global equilibria as a consistent estimator of the clients’ moment conditions.

1845 Our work also differs from federated minimax optimization algorithms: Sharma et al. (2022); Shen  
1846 et al. (2024); Wu et al. (2024); Zhu et al. (2024), where the motivation is to analyse and improve the  
1847 non-asymptotic convergence under various analytical assumptions on the objective functions. We  
1848 primarily focus on deriving the equilibrium via the limit points of the federated GDA algorithm.

1849

## 1850 E BENCHMARK CONSIDERATIONS AND ADDITIONAL EXPERIMENTS

1851

### 1852 E.1 THE EXPERIMENTAL BENCHMARK DESIGN

1853

1854 As stated, our experiments take the Bennett et al. (2019)’s experiments as a centralized-setting  
1855 baseline. Therefore, we have used the same synthetic dataset as DEEPGMM, which they use in their  
1856 experiments to benchmarks against the baselines therein such as DEEPIV (Hartford et al., 2017). It is  
1857 standard to perform experimental analysis on synthetic datasets for unavailability of ground truth  
1858 for causal inference; for example see Section 4.1.1 of Vo et al. (2022b). As the learning process  
1859 essentially involves estimating the true parameter  $\theta_0$  by  $\hat{\theta}$ , to measure the performance of the learning  
1860 procedure, we use the MSE of the estimate  $\hat{g} := g(\cdot, \hat{\theta})$  against the true  $g_0$  averaged over the clients.  
1861 Nonetheless, an experimental comparison of our work with recent works on federated Bayesian  
1862 methods for causal effect estimations does not apply directly. We discuss that below.

1863 The two works in the domain of federated Bayesian methods for causal effect estimations are  
1864 CAUSALRFF (Vo et al., 2022a) and FEDCI (Vo et al., 2022b). The aim of CAUSALRFF (Vo et al.,  
1865 2022a) is to estimate the conditional average treatment effect (CATE) and average treatment effect  
1866 (ATE), whereas FEDCI (Vo et al., 2022b) aims to estimate individual treatment effect (ITE) and  
1867 ATE. For this, (Vo et al., 2022a) consider a setting of  $Y$ ,  $W$ , and  $X$  to be random variables denoting  
1868 the outcome, treatment, and proxy variable, respectively. Along with that, they also consider a  
1869 confounding variable  $Z$ . However, their causal dependency builds on the dependence of each of  $Y$ ,  
1870  $W$ , and  $X$  on  $Z$  besides dependency of  $Y$  on  $W$ . Consequently, to compute CATE and ATE, they  
1871 need to estimate the conditional probabilities  $p(w^i|x^i)$ ,  $p(y^i|x^i, w^i)$ ,  $p(z^i|x^i, y^i, w^i)$ ,  $p(y^i|w^i, z^i)$ ,  
1872 where the superscript  $i$  represents a client. Their experiments compare the estimates of CATE and  
1873 ATE with the Bayesian baselines (Hill, 2011), (Shalit et al., 2017), (Louizos et al., 2017), etc. in  
1874 a centralized setting without any consideration of data decentralization or heterogeneity native to  
1875 federated learning. Further, they compare against the same baselines in a *one-shot federated* setting,  
1876 where at the end of training on separate data sources independently, the predicted treatment effects  
1877 are averaged. Similar is the experimental evaluation of (Vo et al., 2022b). Similarly, Xiong et al.  
1878 (2023) address a fundamentally different causal setting from ours, as they target ATE/ATT estimation  
1879 in observational studies under the unconfoundedness assumption ( $\{Y(0), Y(1)\} \perp W | X$ ), which  
1880 implies that treatment assignment  $W$  is exogenous given observed covariates  $X$ .

1881 By contrast, the setting of IV analysis as in our work does not consider dependency of the outcome  
1882 variable  $Y$  on the confounder  $Z$ , though the treatment variable  $X$  could be endogenous and depend on  
1883  $Z$ . In our synthetic data generation, an unobserved confounder explicitly enters both the treatment and  
1884 outcome equations, inducing correlation between ( $X$ ) and the residual ( $Y - g_0(X)$ ) and therefore  
1885 violating unconfoundedness by construction. For us, computing the treatment effects and thereby  
1886 comparing it against these works is not direct. Furthermore, it is unclear, if the approach of (Vo  
1887 et al., 2022a) and (Vo et al., 2022b), where the predicted inference over a number of datasets is  
1888 averaged as the final result, would be comparable to our approach where the problem is solved  
1889 using a federated maximin optimization with multiple synchronization rounds among the clients.  
For us, the federated optimization subsumes the experimental of comparing the average predicted  
values after independent training with the predicted value over the entire data. This is the reason that

our centralized counterpart i.e. DEEPGMM (Bennett et al., 2019), do not experimentally compare against the baselines of (Vo et al., 2022a) and (Vo et al., 2022b). In summary, for us the experimental benchmarks were guided by showing the efficient fit of the GMM estimator in a federated setting.

## E.2 ADDITIONAL EXPERIMENTS

Estimations	$Dir_S(\alpha) = 0.1$		$Dir_S(\alpha) = 1.0$	
	FDEEPGMM-GDA	FDEEPGMM-SGDA	FDEEPGMM-GDA	FDEEPGMM-SGDA
FEMNIST <sub>x</sub>	0.27 ± 0.04	0.23 ± 0.02	0.17 ± 0.01	0.19 ± 0.03
FEMNIST <sub>x,z</sub>	0.21 ± 0.01	0.24 ± 0.04	0.16 ± 0.03	0.18 ± 0.02
FEMNIST <sub>z</sub>	0.29 ± 0.02	0.25 ± 0.03	0.20 ± 0.04	0.23 ± 0.01
CIFAR10 <sub>x</sub>	0.26 ± 0.01	0.27 ± 0.01	0.18 ± 0.01	0.15 ± 0.02
CIFAR10 <sub>x,z</sub>	0.29 ± 0.02	0.30 ± 0.01	0.21 ± 0.02	0.13 ± 0.01
CIFAR10 <sub>z</sub>	1.73 ± 0.01	0.67 ± 0.02	0.37 ± 0.05	0.35 ± 0.02

Table 2: The averaged Test MSE with standard deviation in the high-dimensional scenarios with varying levels of heterogeneity.

The experimental results included in Section 4 were conducted setting  $Dir_S(\alpha) = 0.3$ , which corresponds to the case wherein a dataset with 10 classes, such as MNIST and CIFAR10, samples of 3 classes on average will be distributed to each client (Hsu et al., 2019). To further investigate the effect of heterogeneity on the performance of FEDDEEPGMM, we conducted experiments with  $Dir_S(\alpha) = 0.1$  and  $Dir_S(\alpha) = 1$ .  $Dir_S(\alpha) = 0.1$  would correspond to the case when every client would have samples from one class on average from a dataset with 10 classes, which represents a high heterogeneity setting. Whereas, setting  $Dir_S(\alpha) = 1$ , the data distribution across clients with regards to samples from different classes becomes roughly uniform representing a near homogeneous scenario. The experimental results are presented in Table 2.

The results presented in Table 2 indicate that on decreasing  $Dir_S(\alpha)$  from 0.3 to 0.1, i.e. increasing heterogeneity, the Test MSE achieved increases marginally. Whereas, on increasing  $Dir_S(\alpha)$  from 0.3 to 1.0, i.e. decreasing heterogeneity, the Test MSE achieved decreases. This set of observations corroborate our theoretical insight that the consistency of the GMM estimator depends on the heterogeneity bias. The change in the MSE values being only marginal can be attributed to the overparametrized setting offered by the CNN on a small-sized data on each client as well as hyperparameter tuning.