FEDERATED EQUILIBRIUM SOLUTIONS FOR GENERAL-IZED METHOD OF MOMENTS APPLIED TO INSTRUMEN-TAL VARIABLE ANALYSIS

Anonymous authorsPaper under double-blind review

ABSTRACT

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

1 Introduction

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

A classical approach to address such limitations is Instrumental variables (IV) analysis, which assumes conditional independence between a confounding variable and the outcome while considering its causal effect on the treatment variable. IV analysis can very practically apply to training Dayan et al. (2021)'s ML model wherein the patients' HIF levels work as an IV that influences the effective organ-level oxygen consumption (a treatment variable) but does not directly affect the mortality of the COVID-19 patients (the outcome). IV analysis has been comprehensively explored in econometrics (Angrist & Krueger, 2001; Angrist & Pischke, 2009) with several decades of history, such as works of Wright (1928) and Reiersøl (1945). Its efficiency is now accepted for learning even high-dimensional 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.

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

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

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

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

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

This work focuses on the existence results of federated equilibrium solutions, federated consistent GMM estimators, and thereby structurally solving the federated IV analysis problem. The existence of approximate client-local equilibria via federated solution has applications beyond the GMM and IV analysis, to problems such as federated generative adversarial networks (FedGAN) (Rasouli et al., 2020), where a Nash Equilibrium may not exist (Farnia & Ozdaglar, 2020). However, the scope of our discussion does not include FedGAN, or a new federated minimax algorithm, or, for that matter, the convergence theory and scalability. We leave an open problem to characterize and recover a 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.

2 Preliminaries and Model

Client-local causal inference. We model our basic terminologies after (Bennett et al., 2019) for a client-local setting. Consider a distributed system as a set of N clients [N] with datasets $S^i = \{(x_i^i, y_i^i)\}_{i=1}^{n_i}, \ \forall i \in [N]$. We assume that for a client $i \in [N]$, the treatment and outcome

¹Wu et al. (2023) used FEDGMM as an acronym for federated Gaussian mixture models.

variables x^i_j and y^i_j , respectively, are related by the process $Y^i = g^i_0(X^i) + \epsilon^i, \ i \in [N]$. We assume that each client-local residual ϵ^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, $g^i_0(X^i) \neq \mathbb{E}[Y^i|X^i]$. We assume that the treatment variables are influenced by instrumental variables $Z^i, \forall i \in [N]$ so that

$$P(X^i|Z^i) \neq P(X^i). \tag{1}$$

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

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

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 $g_0(.) \in G := \{g(.,\theta) | \theta \in \Theta\}$ expressed as $g_0 := g(.,\theta_0)$ for $\theta_0 \in \Theta$, defined by

$$g(.,\theta_0) = \frac{1}{N} \sum_{i=1}^{N} g^i(.,\theta_0).$$
 (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 $f^i : \mathbb{R}^{|\mathcal{Z}|} \to \mathbb{R}^m$ with components $f^i_1, f^i_2, \ldots, f^i_m$. Based on equation (2), we define the moment conditions using parametrized functions $\{f^i_j\}_{j=1}^m \ \forall i \in [N]$ as

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

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 θ_0 . With (4), we define the moment conditions on a client $i \in [N]$ as

$$\psi(f_i^i;\theta) = 0, \ \forall j \in [m],\tag{5}$$

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

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

where $\psi_{n_i}(f^i;\theta) = \left(\psi_{n_i}(f^i_1;\theta),\psi_{n_i}(f^i_2;\theta),\ldots,\psi_{n_i}(f^i_m;\theta)\right)$ is the moment condition vector, and $\psi_{n_i}(f^i_j;\theta) = \frac{1}{n_i}\sum_{k=1}^{n_i}f^i_j(Z^i_k)(Y^i_k-g^i(X^i_k;\theta))$. Thus, for empirical estimation of the causal response function g^i_0 at client $i\in[N]$, it needs to satisfy

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

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

$$\min_{\theta \in \Theta} \|\psi_{n_i}(f_1^i; \theta), \psi_{n_i}(f_2^i; \theta), \dots, \psi_{n_i}(f_m^i; \theta)\|^2,$$
(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 $\tilde{\theta}$, to cater to the cases of (finitely) large number of moment conditions. We adapt their weighted norm $\|w\|_{\tilde{\theta}}^2 = w^T \mathcal{C}_{\tilde{\theta}}^{-1} w$, to a client-local setting via the covariance matrix $\mathcal{C}_{\tilde{\theta}}$ defined by

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

Now considering the vector space $\mathcal V$ of real-valued functions f_i of Z, $\psi_{n_i}(f^i;\theta) = (\psi_{n_i}(f^i_1;\theta),\psi_{n_i}(f^i_2;\theta),\ldots,\psi_{n_i}(f^i_m;\theta))$ is a linear operator on $\mathcal V$ and

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

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

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

which leads to the following client-local optimization problem:

$$\theta^{\text{GMM}} \in \underset{\theta \in \Theta}{\operatorname{arg\,min}} \, \Psi_{n_i}(\theta, \mathcal{F}^i, \tilde{\theta}),$$
(11)

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

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

$$U_{\tilde{\theta}}^{i}(\theta,\tau) := \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} f^{i}(Z_{k}^{i},\tau) \left(Y_{k}^{i} - g^{i}(X_{k}^{i};\theta)\right) - \frac{1}{4n_{i}} \sum_{k=1}^{n_{i}} \left(f^{i}(Z_{k}^{i},\tau)\right)^{2} \left(Y_{k}^{i} - g^{i}(X_{k}^{i};\theta)\right)^{2}$$
(12)

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

$$\theta^{\text{DGMM}} \in \underset{\theta \in \Theta}{\arg \min} \sup_{\tau \in \mathcal{T}} U_{\bar{\theta}}^{i}(\theta, \tau).$$
 (13)

Equation (13) forms a zero-sum game, whose equilibrium solution is shown to be a true estimator to θ_0 under a set of standard assumptions; see Theorem 2 in (Bennett et al., 2019).

3 FEDERATED DEEP GMM VIA FEDERATED EQUILIBRIUM SOLUTIONS

3.1 FEDERATED DEEP GENERALIZED METHOD MOMENT (FEDDEEPGMM)

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

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

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

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

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

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

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

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

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

$$\psi_{n_i}(f^i;\theta) := \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z^i_k) (Y^i_k - g^i(X^i_k;\theta)), \text{ and } \mathcal{C}_{\tilde{\theta}}(f^i,f^i) := \frac{1}{n_i} \sum_{k=1}^{n_i} (f^i(Z^i_k))^2 (Y^i_k - g^i(X^i_k;\tilde{\theta}))^2.$$

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

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

where $U_{\hat{\theta}}^{i}(\theta,\tau)$ 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.

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 $U^i_{\tilde{\theta}}(\theta,\tau) \ \forall i \in [N]$ is twice continuously differentiable for both θ and τ . Thus, the global objective $U_{\tilde{\theta}}(\theta,\tau)$ is also a twice continuously differentiable function.

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

$$\begin{split} \|\nabla_{\theta} U^{i}_{\tilde{\theta}}(\theta_{1},\tau_{1}) - \nabla_{\theta} U^{i}_{\tilde{\theta}}(\theta_{2},\tau_{2})\| &\leq L \|(\theta_{1},\tau_{1}) - (\theta_{2},\tau_{2})\|, \text{ and } \\ \|\nabla_{\tau} U^{i}_{\tilde{\theta}}(\theta_{1},\tau_{1}) - \nabla_{\tau} U^{i}_{\tilde{\theta}}(\theta_{2},\tau_{2})\| &\leq L \|(\theta_{1},\tau_{1}) - (\theta_{2},\tau_{2})\|, \end{split}$$

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

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

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

where $\zeta_{\theta}^{i},\ \zeta_{\tau}^{i}\geq0$ 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. θ and τ is bounded as follows:

$$\begin{split} \|\nabla^2_{\theta\theta}U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\theta\theta}U_{\tilde{\theta}}(\theta,\tau)\|_{\sigma} &\leq \rho^i_{\theta}, \\ \|\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\theta,\tau)\|_{\sigma} &\leq \rho^i_{\tau}, \\ \|\nabla^2_{\theta\tau}U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\theta\tau}U_{\tilde{\theta}}(\theta,\tau)\|_{\sigma} &\leq \rho^i_{\theta\tau}, \\ \|\nabla^2_{\tau\theta}U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\tau\theta}U_{\tilde{\theta}}(\theta,\tau)\|_{\sigma} &\leq \rho^i_{\tau\theta}, \end{split}$$

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

Assumptions 3 and 4 provide a measure of data heterogeneity across clients in a federated setting. In the special case, when ζ and ρ '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 $U_{\tilde{\theta}}(\theta,\tau)$ 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 $h(\delta) \to 0$ as $\delta \to 0$. There exists a δ_0 , such that for any $\delta \in (0, \delta_0]$, and any (θ, τ) such that $\|\theta - \hat{\theta}\| \le \delta$ and $\|\tau - \hat{\tau}\| \le \delta$, 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:

$$U_{\tilde{\theta}}(\hat{\theta}, \tau) \le U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \le \max_{\tau \prime : \|\tau \prime - \hat{\tau}\| \le h(\delta)} U_{\tilde{\theta}}(\theta, \tau \prime). \tag{18}$$

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:*

• Second-order Sufficient Condition: A stationary point (θ, τ) that satisfies $\nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta, \tau) \prec 0$, and $\left[\nabla^2_{\theta\theta} U_{\tilde{\theta}} - \nabla^2_{\theta\tau} U_{\tilde{\theta}} \left(\nabla^2_{\tau\tau} U_{\tilde{\theta}}\right)^{-1} \nabla^2_{\tau\theta} U_{\tilde{\theta}}\right](\theta, \tau) \succ 0$ guarantees that (θ, τ) 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 $h(\delta) \to 0$ as $\delta \to 0$. There exists a δ_0 , such that for any $\delta \in (0,\delta_0]$, and any (θ,τ) such that $\|\theta - \hat{\theta}\| \le \delta$ and $\|\tau - \hat{\tau}\| \le \delta$, then a point $(\hat{\theta},\hat{\tau})$ is an ε -approximate local minimax point of U, if it satisfies:

$$U_{\tilde{\theta}}(\hat{\theta},\tau) - \varepsilon \le U_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) \le \max_{\tau : \|\tau t - \hat{\tau}\| \le h(\delta)} U_{\tilde{\theta}}(\theta,\tau t) + \varepsilon, \tag{19}$$

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 $U^i_{\hat{\theta}}(\theta,\tau)$ be a function defined over $\Theta \times \mathcal{T}$ for a client $i \in [N]$ and $U_{\tilde{\theta}}(\theta,\tau) := \frac{1}{N} \sum_{i=1}^N U^i_{\tilde{\theta}}(\theta,\tau)$. An \mathcal{E} -approximate federated equilibrium point $(\hat{\theta},\hat{\tau})$ (that is an ε^i -approximate local minimax point for each client's objective $U^i_{\tilde{\theta}}$), must follow the conditions below:

- 1. ε^i First-order Necessary Condition: The point $(\hat{\theta}, \hat{\tau})$ must be an ε^i stationary point for every client $i \in [N]$, i.e., $\|\nabla_{\theta} U^i_{\hat{\theta}}(\hat{\theta}, \hat{\tau})\| \leq \varepsilon^i$, and $\|\nabla_{\tau} U^i_{\hat{\theta}}(\hat{\theta}, \hat{\tau})\| \leq \varepsilon^i$.
- 2. **Second-Order** ε^{i} **Necessary Condition:** The point $(\hat{\theta}, \hat{\tau})$ must satisfy the second-order conditions: $\nabla^{2}_{\tau\tau}U^{i}_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \preceq -\varepsilon^{i}I$, and $\left[\nabla^{2}_{\theta\theta}U^{i}_{\tilde{\theta}} \nabla^{2}_{\theta\tau}U^{i}_{\tilde{\theta}}\left(\nabla^{2}_{\tau\tau}U_{\tilde{\theta}}\right)^{-1}\nabla^{2}_{\tau\theta}U^{i}_{\tilde{\theta}}\right](\hat{\theta}, \hat{\tau}) \succeq \varepsilon^{i}I$.
- 3. **Second-Order** ε^{i} **Sufficient Condition:** An ε^{i} stationary point (θ, τ) that satisfies $\nabla^{2}_{\tau\tau}U^{i}_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \prec -\varepsilon^{i}I$, and $\left[\nabla^{2}_{\theta\theta}U^{i}_{\hat{\theta}} \nabla^{2}_{\theta\tau}U^{i}_{\hat{\theta}}\left(\nabla^{2}_{\tau\tau}U^{i}_{\hat{\theta}}\right)^{-1}\nabla^{2}_{\tau\theta}U^{i}_{\hat{\theta}}\right](\hat{\theta}, \hat{\tau}) \succ \varepsilon^{i}I$ guarantees that $(\hat{\theta}, \hat{\tau})$ is a strict local minimax point $\forall i \in [N]$ that satisfies ε^{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: $U_{\tilde{\theta}}(\hat{\theta}, \tau) \leq U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \leq \max_{\tau': \|\tau' - \hat{\tau}\| \leq h(\delta)} U_{\tilde{\theta}}(\theta, \tau')$, is an \mathcal{E} -approximate federated equilibrium solution as defined in 3, where the approximation error ε^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 $\alpha := \left|\lambda_{\max}\left(\nabla_{\tau\tau}^2 U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau})\right)\right|$, $\beta := \lambda_{\min}\left(\left[\nabla_{\theta\theta}^2 U_{\tilde{\theta}} - \nabla_{\theta\tau}^2 U_{\tilde{\theta}} \left(\nabla_{\tau\tau}^2 U_{\tilde{\theta}}\right)^{-1} \nabla_{\tau\theta}^2 U_{\tilde{\theta}}\right](\hat{\theta}, \hat{\tau})\right)$ and $B^i := \rho_{\theta}^i + L\rho_{\theta\tau}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2 U_{\tilde{\theta}}^i)|} + L\rho_{\tau\theta}^i \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^2$

The proof of Theorem 1 is given in Appendix B.2. Note that when data is homogeneous (i.e., for each client i, ζ_{θ}^{i} , ζ_{τ}^{i} , ρ_{τ}^{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 ε^i exists and $(\hat{\theta},\hat{\tau})$ fails to be a local ε^i approximate equilibrium point for that clients. It may happen in two cases:

1. The gradient dissimilarity ζ_{θ}^{i} , ζ_{τ}^{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 $h(\delta) \to 0$ as $\delta \to 0$. 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 $\eta^i(\epsilon) := \inf_{d(\theta, \theta_0) \ge \epsilon} M^i(\theta) - M^i(\theta_0)$ for every $\epsilon > 0$. Fix some δ_0 , for any $\delta \in (0, \delta_0]$ and any (θ, τ) such that $\|\theta - \hat{\theta}\| \le \delta$ and $\|\tau - \hat{\tau}\| \le \delta$, 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 θ_0 , i.e. $\hat{\theta}_n \stackrel{p}{\to} \theta_0$ 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 ε^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 (θ,τ) . The convergence results of variants of FEDGDA (Sharma et al., 2022; Shen et al., 2024; Wu et al., 2024) assume that $U_{\tilde{\theta}}(\theta,\tau)$ is non-convex on θ and satisfies a μ -Polyak Łojasiewicz (PL) inequality on τ , 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 θ and τ , respectively. Without loss of generality the FEDGDA updates are:

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

We call it γ -FEDGDA, where γ is the ratio of α_1 to α_2 . As $\eta \to 0$ corresponds to FEDGDA-flow, under the smoothness of $U^i_{\tilde{\theta}}$, bounded gradient heterogeneity (assumption 3) and for fixed local

rounds R, FEDGDA-flow becomes:

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

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 γ -FEDGDA flow.

 $\operatorname{Re}(\Lambda_i) < 0$ for all j.

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

Let γ - \mathcal{FGDA} be the set of strictly linearly stable points of the γ -FEDGDA flow, and \mathcal{L} oc \mathcal{M} inimax be the set of local minimax points of the federated zero-sum game. Define

$$\overline{\infty - \mathcal{FGDA}} := \limsup_{\gamma \to \infty} \gamma - \mathcal{FGDA} := \bigcap_{\gamma_0 > 0} \cup_{\gamma > \gamma_0} \gamma - \mathcal{FGDA}, \text{ and } \\ \underline{\infty - \mathcal{FGDA}} := \liminf_{\gamma \to \infty} \gamma - \mathcal{FGDA} := \cup_{\gamma_0 > 0} \cap_{\gamma > \gamma_0} \gamma - \mathcal{FGDA}.$$

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

Theorem 3. Under Assumption 1, $\mathcal{L}oc\mathcal{M}inimax \subset \infty - \mathcal{FGDA} \subset \overline{\infty - \mathcal{FGDA}} \subset$ $\mathcal{L}oc\mathcal{M}inimax \cup \mathcal{A}$, where $\mathcal{A} := \{(\theta, \tau) | (\theta, \tau) \text{ is stationary an} \overline{d \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta, \tau)} \text{ is degenerate} \}$. Moreover, if the hessian $\nabla^2_{ au au}U_{\tilde{\theta}}(\theta, au)$ is smooth, then \mathcal{A} has measure zero in $\Theta imes\mathcal{T}\subset\mathbb{R}^d imes\mathbb{R}^k$.

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.

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**: $g_0(x) = 1_{\{x>0\}}$, (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:

$$Y = g_0(X) + e + \delta \qquad \text{and } X = Z^{(1)} + Z^{(2)} + e + \gamma$$
 (20)
$$(Z^{(1)}, Z^{(2)}) \sim \text{Uniform}([-3, 3]^2) \qquad \text{and } e \sim \mathcal{N}(0, 1), \quad \gamma, \delta \sim \mathcal{N}(0, 0.1)$$
 (21)

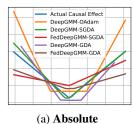
$$(Z^{(1)}, Z^{(2)}) \sim \text{Uniform}([-3, 3]^2)$$
 and $e \sim \mathcal{N}(0, 1), \quad \gamma, \delta \sim \mathcal{N}(0, 0.1)$ (21)

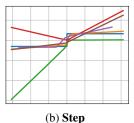
(b) Z and X are low-dimensional or high-dimensional or both. First, Z and X are generated as in (20,21). Then for high-dimensional data, we map Z and X to an image using the mapping:

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

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





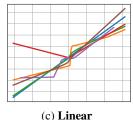


Figure 1: Estimated \hat{g} compared to true g in low-dimensional scenarios

Estimations	DEEPGMM-	DEEPGMM-	FDEEPGMM-	DEEPGMM-	FDEEPGMM-
	OAdam	GDA	GDA	SGDA	SGDA
Absolute	0.03 ± 0.01	$0.013 \pm .01$	0.4 ± 0.01	0.009 ± 0.01	0.2 ± 0.00
Step	0.3 ± 0.00	0.03 ± 0.00	0.04 ± 0.01	0.112 ± 0.00	0.23 ± 0.01
Linear	0.01 ± 0.00	0.02 ± 0.00	0.01 ± 0.00	0.03 ± 0.00	0.04 ± 0.00
$\overline{\text{FEMNIST}_{x}}$	0.50 ± 0.00	1.11 ± 0.01	0.21 ± 0.02	0.40 ± 0.01	0.19 ± 0.01
$\overline{\text{FEMNIST}_{x,z}}$	0.24 ± 0.00	0.46 ± 0.09	0.19 ± 0.03	0.14 ± 0.02	0.20 ± 0.00
$\mathbf{FEMNIST_{z}}$	0.10 ± 0.00	0.42 ± 0.01	0.24 ± 0.01	0.11 ± 0.02	0.23 ± 0.01
$CIFAR10_x$	0.55 ± 0.30	0.19 ± 0.01	0.25 ± 0.03	0.20 ± 0.08	0.22 ± 0.08
$CIFAR10_{x,z}$	0.40 ± 0.11	0.24 ± 0.00	0.24 ± 0.03	0.19 ± 0.03	0.22 ± 0.02
$CIFAR10_z$	0.13 ± 0.03	0.13 ± 0.01	1.70 ± 2.60	0.24 ± 0.01	0.52 ± 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 α 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.

AN OPEN PROBLEM

In this work, we characterized the equilibrium solutions of federated zero-sum games in consideration of local minimax solutions for non-convex non-concave minimax optimization problems. Regardless of the analytical assumptions over the objective, the mixed strategy solutions for zero-sum games exist. However, unlike the pure strategy solutions, where the standard heterogeneity considerations over gradients and Hessians across clients, translates a local minimax solution for the federated objective to approximate local solutions for the clients, it is not immediate how a mixed strategy solution as a probability measure can be translated to that for clients. It leaves an interesting open problem to characterize the mixed startegy solutions for federated zero-sum games.

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.
- Emmanual 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.
- Zachary Charles and Dimitris Papailiopoulos. Stability and generalization of learning algorithms that converge to global optima. In *International conference on machine learning*, pp. 745–754. PMLR, 2018.
- Bapi Chatterjee, Vyacheslav Kungurtsev, and Dan Alistarh. Federated sgd with local asynchrony. In 2024 IEEE 44th International Conference on Distributed Computing Systems (ICDCS), pp. 857–868. IEEE, 2024.
- Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training gans with optimism. In *International Conference on Learning Representations*, 2018.
- Ittai Dayan, Holger R Roth, Aoxiao Zhong, Ahmed Harouni, Amilcare Gentili, Anas Z Abidin, Andrew Liu, Anthony Beardsworth Costa, Bradford J Wood, Chien-Sung Tsai, et al. Federated learning for predicting clinical outcomes in patients with covid-19. *Nature medicine*, 27(10): 1735–1743, 2021.
- Yuyang Deng and Mehrdad Mahdavi. Local stochastic gradient descent ascent: Convergence analysis and communication efficiency. In *International Conference on Artificial Intelligence and Statistics*, pp. 1387–1395. PMLR, 2021.
- Farzan Farnia and Asuman Ozdaglar. Do gans always have nash equilibria? In *International Conference on Machine Learning*, pp. 3029–3039. PMLR, 2020.
- Melike Gecer and Benoit Garbinato. Federated learning for mobility applications. *ACM Computing Surveys*, 56(5):1–28, 2024.
- Lars Peter Hansen. Large sample properties of generalized method of moments estimators. *Econometrica: Journal of the econometric society*, pp. 1029–1054, 1982.
- Jason Hartford, Greg Lewis, Kevin Leyton-Brown, and Matt Taddy. Deep iv: A flexible approach for counterfactual prediction. In *International Conference on Machine Learning*, pp. 1414–1423. PMLR, 2017.
- Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2 edition, 2012.
 - Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv* preprint arXiv:1909.06335, 2019.

- Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 4880–4889. PMLR, 07 2020. URL https://proceedings.mlr.press/v119/jin20e.html.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.
- Diederik P Kingma. Adam: A method for stochastic optimization. ICLR, 2015.

- Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- Greg Lewis and Vasilis Syrgkanis. Adversarial generalized method of moments, 2018. URL https://arxiv.org/abs/1803.07164.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020.
- Yafen Liang, Wei Ruan, Yandong Jiang, Richard Smalling, Xiaoyi Yuan, and Holger K Eltzschig. Interplay of hypoxia-inducible factors and oxygen therapy in cardiovascular medicine. *Nature Reviews Cardiology*, 20(11):723–737, 2023.
- Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59: 85–116, 2022.
- Guodong Long, Yue Tan, Jing Jiang, and Chengqi Zhang. Federated learning for open banking. In *Federated learning: privacy and incentive*, pp. 240–254. Springer, 2020.
- Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. *Advances in neural information processing systems*, 30, 2017.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Maxence Noble, Aurélien Bellet, and Aymeric Dieuleveut. Differentially private federated learning on heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp. 10110–10145. PMLR, 2022.
- Wonsuk Oh and Girish N Nadkarni. Federated learning in health care using structured medical data. *Advances in kidney disease and health*, 30(1):4–16, 2023.
- Judea Pearl. Causal inference in statistics: An overview. 2009.
- Mohammad Rasouli, Tao Sun, and Ram Rajagopal. Fedgan: Federated generative adversarial networks for distributed data. *arXiv preprint arXiv:2006.07228*, 2020.
- Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečnỳ, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
 - Olav Reiersøl. Confluence analysis by means of instrumental sets of variables. PhD thesis, Almqvist & Wiksell, 1945.
 - Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *International conference on machine learning*, pp. 3076–3085. PMLR, 2017.

Pranay Sharma, Rohan Panda, Gauri Joshi, and Pramod Varshney. Federated minimax optimization: Improved convergence analyses and algorithms. In *International Conference on Machine Learning*, pp. 19683–19730. PMLR, 2022.

- Wei Shen, Minhui Huang, Jiawei Zhang, and Cong Shen. Stochastic smoothed gradient descent ascent for federated minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 3988–3996. PMLR, 2024.
- 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.
- Thanh Vinh Vo, Young Lee, Trong Nghia Hoang, and Tze-Yun Leong. Bayesian federated estimation of causal effects from observational data. In *Uncertainty in Artificial Intelligence*, pp. 2024–2034. PMLR, 2022b.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *Advances in neural information processing systems*, 30, 2017.
- Jeffrey M Wooldridge. Applications of generalized method of moments estimation. *Journal of Economic perspectives*, 15(4):87–100, 2001.
- Philip Green Wright. The tariff on animal and vegetable oils. Number 26. Macmillan, 1928.
- Xidong Wu, Jianhui Sun, Zhengmian Hu, Aidong Zhang, and Heng Huang. Solving a class of non-convex minimax optimization in federated learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yue Wu, Shuaicheng Zhang, Wenchao Yu, Yanchi Liu, Quanquan Gu, Dawei Zhou, Haifeng Chen, and Wei Cheng. Personalized federated learning under mixture of distributions. In *International Conference on Machine Learning*, pp. 37860–37879. PMLR, 2023.
- Ruoxuan Xiong, Allison Koenecke, Michael Powell, Zhu Shen, Joshua T Vogelstein, and Susan Athey. Federated causal inference in heterogeneous observational data. *Statistics in Medicine*, 42 (24):4418–4439, 2023.
- Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning: State-of-the-art and research challenges. *ACM Computing Surveys*, 56(3):1–44, 2023.
- Mishael Zedek. Continuity and location of zeros of linear combinations of polynomials. *Proceedings of the American Mathematical Society*, 16(1):78–84, 1965. ISSN 00029939, 10886826. URL http://www.jstor.org/stable/2034005.
- Weijia Zhang, Lin Liu, and Jiuyong Li. Treatment effect estimation with disentangled latent factors. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10923–10930, 2021.
- Miaoxi Zhu, Li Shen, Bo Du, and Dacheng Tao. Stability and generalization of the decentralized stochastic gradient descent ascent algorithm. *Advances in Neural Information Processing Systems*, 36, 2024.

APPENDIX

A Federated Gradient Descent Ascent Algorithm Description						
B Proofs		14				
B.1 Proof of Lemma 1		14				
B.2 Proof of Theorem 1		16				
B.3 Consistency		19				
B.3.1 Assumptions		19				
B.3.2 Proof of Theorem 2		20				
C Limit Points of FEDGDA		23				
C.1 FEDGDA Flow		24				
C.2 Proof of Theorem 3		25				
D Related Work		27				
E Benchmark Considerations and Additional Experiments		27				
E.1 The Experimental Benchmark Design		27				
E.2 Additional Experiments		28				

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

Server Input: initial global estimate θ_1, τ_1 ; constant local learning rate α_1, α_2 ; total N clients **Output**: global model states θ_{T+1}, τ_{T+1}

```
1: for synchronization round t = 1, ..., T do
                 server sends \theta_t, \tau_t to all clients
  2:
  3:
                 for each i \in [N] in parallel do
                          \theta_{t,1}^i \leftarrow \theta_t, \, \tau_{t,1}^i \leftarrow \tau_t
  4:
                         for r=1,2,\ldots,R do \theta^i_{t,r+1}=\theta^i_{t,r}-\alpha_1\nabla_{\theta}f_i(\theta^i_{t,r},\tau^i_{t,r}) \tau^i_{t,r+1}=\tau^i_{t,r}+\alpha_2\nabla_{\tau}f_i(\theta^i_{t,r},\tau^i_{t,r})
  5:
  6:
  7:
  8:
                         (\Delta \theta_t^i, \Delta \tau_t) \leftarrow (\theta_{t,R+1}^i - \theta_t, \tau_{t,R+1}^i - \tau_t)
  9:
10:
                 (\Delta \theta_t, \Delta \tau_t) \leftarrow \frac{1}{N} \sum_{i \in [N]} (\Delta \theta_t^i, \Delta \tau_t^i)
11:
                 \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}
```

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 Łojaisiewicz (PL) condition in τ). The function $U_{\tilde{\theta}}$ satisfyies $\mu - PL$ condition in τ , $\mu > 0$, if for any fixed θ , $\arg\max_{\tau'} U_{\tilde{\theta}}(\theta, \tau') \neq \phi$ and $\|\nabla_{\tau} U_{\tilde{\theta}}(\theta, \tau)\|^2 \geq 2\mu \left(\max_{\tau'} U_{\tilde{\theta}}(\theta, \tau') - U_{\tilde{\theta}}(\theta, \tau)\right)$.

Theorem 4. Let the local loss functions $U^i_{\bar{\theta}}$ for all $i \in \{1, 2, ..., N\}$ satisfy assumption 2 and 3. The federated objective function satisfies assumption 5. Suppose $\alpha_2 \leq \frac{1}{8LR}$, $\frac{\alpha_1}{\alpha_2} \leq \frac{1}{8\kappa^2}$, where $\kappa = \frac{L}{\mu}$ is the condition number. Let $\bar{\theta}_{T+1}$ is drawn uniformly at random from $\{\theta_t\}_{t=1}^{T+1}$, then the following holds:

$$\|\nabla \tilde{\Phi}(\bar{\theta}_{T+1})\|^2 \leq \mathcal{O}\left(\kappa^2 \left[\frac{\Delta_{\tilde{\Phi}}}{\alpha_2 R(T+1)}\right]\right) + \mathcal{O}\left(\kappa^2 (R-1)^2 \left[\alpha_2^2 \zeta_{\tau}^2 + \alpha_1^2 \zeta_{\theta}^2\right]\right)$$

where $\nabla \tilde{\Phi}(.) := \max_{\tau} U_{\tilde{\theta}}(.,\tau)$ is the envelope function, $\Delta_{\tilde{\Phi}} := \tilde{\Phi}(\theta_0) - \min_{\theta} \tilde{\Phi}(\theta)$, and $\zeta_{\theta} := \frac{1}{N} \sum_{i=1}^{N} \zeta_{\theta}^{i}$, $\zeta_{\tau} := \frac{1}{N} \sum_{i=1}^{N} \zeta_{\tau}^{i}$. Using $\alpha_1 = \mathcal{O}\left(\frac{1}{\kappa^2} \sqrt{\frac{N}{R(T+1)}}\right)$, $\alpha_2 = \mathcal{O}\left(\sqrt{\frac{N}{R(T+1)}}\right)$, $\|\nabla \tilde{\Phi}(\bar{\theta}_{T+1})\|^2$ can be bounded as

$$\mathcal{O}\left(\frac{\kappa^2\Delta_{\tilde{\Phi}}}{\sqrt{NR(T+1)}} + \kappa^2(R-1)^2 \frac{NR(\zeta_{\theta}^2 + \zeta_{\tau}^2)}{R(T+1)}\right).$$

Although the original assumption uses the supremum of average squared deviations, say ζ_{θ}' and ζ_{τ}' , we use per-client dissimilarity bounds ζ_{θ}^i , ζ_{τ}^i and upper bound their quantity as ${\zeta_{\theta}'}^2 \leq \frac{1}{N} \sum_{i=1}^N (\zeta_{\theta}^i)^2 := \zeta_{\theta}^2$ and ${\zeta_{\tau}'}^2 \leq \frac{1}{N} \sum_{i=1}^N (\zeta_{\tau}^i)^2 := \zeta_{\tau}^2$. Since there is no stochasticity, we used the bounded variance $\sigma = 0$. For details, refer to proof of Theorem 1 in (Sharma et al., 2022).

B Proofs

B.1 PROOF OF LEMMA 1

Lemma 3 (Restatement of Lemma 1). Let $\mathcal{F} = span\{f_j^i \mid i \in [N], j \in [m]\}$. An equivalent objective function for the federated moment estimation optimization problem (16) is given by:

$$\|\psi_N(f;\theta)\|_{\tilde{\theta}}^2 = \sup_{\substack{f^i \in \mathcal{F} \\ \forall i \in [N]}} \frac{1}{N} \sum_{i=1}^N \left(\psi_{n_i}(f^i;\theta) - \frac{1}{4} \mathcal{C}_{\tilde{\theta}}(f^i;f^i) \right), \text{ where}$$
 (22)

$$\psi_{n_i}(f^i;\theta) := \frac{1}{n_i} \sum_{k=1}^{n_i} f^i(Z^i_k) (Y^i_k - g^i(X^i_k;\theta)), \text{ and } \mathcal{C}_{\tilde{\theta}}(f^i,f^i) := \frac{1}{n_i} \sum_{k=1}^{n_i} (f^i(Z^i_k))^2 (Y^i_k - g^i(X^i_k;\tilde{\theta}))^2.$$

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 $||v||^2 = v^\top C_{\tilde{\theta}}^{-1} v$ and the associated dual norm is obtained as $||v||_* = \sup_{||v|| \le 1} v^\top v = v^\top C_{\tilde{\theta}} v$.

Using the definition of the dual norm,

$$\|\psi\| = \sup_{\|v\|_{*} \leq 1} v^{\top} \psi$$

$$\|\psi\|^{2} = \sup_{\|v\|_{*} \leq \|\psi\|} v^{\top} \psi$$

$$\|\psi\|^{2} = \sup_{v^{\top} C_{\bar{\theta}} v \leq \|\psi\|^{2}} v^{\top} \psi.$$
(23)

We now find the equivalent dual optimization problem for (23).

The Lagrangian of the constrained maximization problem (23) is given as

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

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

When $\|\psi\|>0$, v=0 satisfies the Slater's condition as a strictly feasible interior point of the constraint $v^{\top}C_{\tilde{\theta}}v-\|\psi\|^2\leq 0$. Thus, strong duality holds. Substituting $v=\frac{-1}{2\lambda}C_{\tilde{\theta}}^{-1}\psi$ in the Lagrangian gives

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

Hence, the dual becomes $\|\psi\|^2 = \inf_{\lambda < 0} \{\mathcal{L}^*(\lambda)\}$. Thus, the equivalent dual optimization problem for (23) is given as

$$\|\psi\|^2 = \inf_{\lambda < 0} \left\{ -\frac{\|\psi\|^2}{4\lambda} - \lambda \|\psi\|^2 \right\}. \tag{24}$$

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 $\|\psi\|^2 = \sup_v \mathcal{L}(v, \frac{-1}{2}) = \sup_v v^\top \psi - \frac{1}{2}(v^\top C_{\tilde{\theta}}v - \|\psi\|^2)$.

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

$$\|\psi\|^2 = \sup_{u} u^{\mathsf{T}} \psi - \frac{1}{4} u^{\mathsf{T}} C_{\tilde{\theta}} u.$$

Using change of variables $u \to v$

$$\|\psi\|^2 = \sup_{v} v^{\mathsf{T}} \psi - \frac{1}{4} v^{\mathsf{T}} C_{\tilde{\theta}} v.$$

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

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, for $f^i \in \mathcal{F}^i$

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

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

$$\mathcal{F}=\operatorname{span}\{f^i_j\mid i\in[N],\ j\in[m]\}.$$

Therefore, v corresponds to f^i such that

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

Hence.

$$v^{\top} \psi = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{m} v_{j} \psi_{n_{i}}(f_{j}^{i}; \theta)$$
$$= \frac{1}{N} \sum_{i=1}^{N} \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} f^{i}(Z_{k}^{i})(Y_{k}^{i} - g^{i}(X_{k}^{i}; \theta)).$$

Similarly,

$$v^{\top}C_{\tilde{\theta}}v = \sum_{p=1}^{m} \sum_{q=1}^{m} v_{p}v_{q}[C_{\tilde{\theta}}]pq$$

$$= \sum_{p=1}^{m} \sum_{q=1}^{m} v_{p}v_{q} \frac{1}{N} \sum_{i=1}^{N} \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} f_{p}^{i}(Z_{k}^{i}) f_{q}^{i}(Z_{k}^{i}) (Y_{k}^{i} - g^{i}(X_{k}^{i}; \tilde{\theta}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} \sum_{p=1}^{m} v_{p} f_{p}^{i}(Z_{k}^{i}) \sum_{q=1}^{m} v_{q} f_{q}^{i}(Z_{k}^{i}) (Y_{k}^{i} - g^{i}(X_{k}^{i}; \tilde{\theta}))^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} (f^{i}(Z_{k}^{i}))^{2} (Y_{k}^{i} - g^{i}(X_{k}^{i}; \tilde{\theta}))^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} C_{\tilde{\theta}}(f^{i}, f^{i}).$$

Thus, applying the Riesz Representation theorem using the representations $v^\top \psi = \frac{1}{N} \sum_{i=1}^N \psi_{n_i}(f^i;\theta)$ 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 form as

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

This gives us the desired result.

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: $U_{\tilde{\theta}}(\hat{\theta},\tau) \leq U_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) \leq \max_{\tau: \|\tau - \hat{\tau}\| \leq h(\delta)} U_{\tilde{\theta}}(\theta,\tau')$, is an \mathcal{E} -approximate federated equilibrium solution as defined in 3, where the approximation error ε^i for each client $i \in [N]$ lies in: $\max\{\zeta^i_{\theta},\zeta^i_{\tau}\} \leq \varepsilon^i \leq \min\{\alpha - \rho^i_{\tau},\beta - B^i\}$ for $\rho^i_{\tau} < \alpha$ and $B^i > \beta$, such that $\alpha := \left|\lambda_{\max}\left(\nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})\right)\right|, \ \beta := \lambda_{\min}\left(\left[\nabla^2_{\theta\theta}U_{\tilde{\theta}} - \nabla^2_{\theta\tau}U_{\tilde{\theta}}\left(\nabla^2_{\tau\tau}U_{\tilde{\theta}}\right)^{-1}\nabla^2_{\tau\theta}U_{\tilde{\theta}}\right](\hat{\theta},\hat{\tau})\right)$ and $B^i := \rho^i_{\theta} + L\rho^i_{\theta\tau}\frac{1}{|\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}})|} + L\rho^i_{\tau\theta}\frac{1}{|\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}})|} + L^2\rho^i_{\tau}\frac{1}{|\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}})|}.$

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

$$U_{\tilde{\theta}}(\hat{\theta}, \tau) \le U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \le \max_{\tau : \|\tau_{t-\tau}^{*}\| \le h(\hat{\delta})} U_{\tilde{\theta}}(\theta, \tau'), \tag{25}$$

We want to show that the ϵ^i - approximate equilibrium for each client's objective $U^i_{\tilde{\theta}}$ also hold individually.

The first-order necessary condition for (25) to hold is $\nabla_{\theta}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})=0$ and $\nabla_{\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})=0$. Thus, $\left\|\nabla_{\theta}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})\right\|^2=0$.

Consider

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

Rearranging

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

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

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

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

In the special case, when $\zeta^i_{\theta}=0$ and $\zeta^i_{\tau}=0$, thus we will have $\left\|\nabla_{\theta}U^i_{\bar{\theta}}(\hat{\theta},\hat{\tau})\right\|^2=\left\|\nabla_{\tau}U^i_{\bar{\theta}}(\hat{\theta},\hat{\tau})\right\|^2=0$ for all $i\in[N]$, which gives $\nabla_{\theta}U^i_{\bar{\theta}}(\hat{\theta},\hat{\tau})=\nabla_{\tau}U^i_{\bar{\theta}}(\hat{\theta},\hat{\tau})=0$ 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 (25), the second-order necessary condition holds for the global function $U_{\tilde{\theta}}$, i.e. $\nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \leq \mathbf{0}$. We now prove that $\nabla^2_{\tau\tau}U_{\tilde{\theta}}^i(\hat{\theta}, \hat{\tau}) \leq \mathbf{0}$.

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

We can write $\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) = \nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) + \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\hat{\theta},\hat{\tau}).$

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,

$$\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) \leq \lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) + \lambda_{\max}(\nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})).$$

Thus,
$$\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) \leq \lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) - \alpha$$
.

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

Under hessian heterogeneity assumption 4

$$\begin{split} \|\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})\|_{\sigma} &= \max \left\{ \left| \lambda_{\max} (\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta,\tau)) \right|, \\ \left| \lambda_{\min} (\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\theta,\tau) - \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta,\tau)) \right| \right\} \\ &\leq \rho^i_{\tau}. \end{split}$$

By definition of the spectral norm $\|\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})\|_{\sigma} = \lambda_{max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})),$

$$\begin{split} \lambda_{max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) &\leq \max\left\{ \left| \lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) \right|, \\ \left| \lambda_{\min}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) \right| \right\} \\ &< \rho^i_{\tau}. \end{split}$$

Thus, $\lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) \leq \lambda_{\max}(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) - \nabla^2_{\tau\tau}U_{\tilde{\theta}}(\hat{\theta},\hat{\tau})) - \alpha \leq \rho^i_{\tau} - \alpha$, where $\rho^i_{\tau} \geq 0$. Hence,

$$\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta},\hat{\tau}) \preceq (\rho^i_{\tau} - \alpha) \mathbf{I}.$$

921 When $\rho_{ au}^i \leq \alpha$, then $\nabla_{ au au}^2 U^i_{\hat{\theta}}(\hat{\theta}, \hat{\tau}) \leq 0$.

Now, since $(\hat{\theta}, \hat{\tau})$ satisfy the equilibrium condition (25), thus $\nabla^2_{\tau\tau} U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \prec 0$ and the Schur complement of $\nabla^2_{\tau\tau} U_{\tilde{\theta}}(\hat{\theta}, \hat{\tau})$ is positive semi-definite. Now when $\rho^i_{\tau} < \alpha$, it follows from above that $\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta}, \hat{\tau}) \prec 0$, hence $\left(\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta}, \hat{\tau})\right)^{-1}$ exists. Now, we need to show that Schur complement of $\nabla^2_{\tau\tau} U^i_{\tilde{\theta}}(\hat{\theta}, \hat{\tau})$ is positive semi-definite.

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

Define $S^i := \left[\nabla^2_{\theta\theta} U^i_{\tilde{\theta}} - \nabla^2_{\theta\tau} U^i_{\tilde{\theta}} \left(\nabla^2_{\tau\tau} U^i_{\tilde{\theta}} \right)^{-1} \nabla^2_{\tau\theta} U^i_{\tilde{\theta}} \right]$. We aim to prove $\lambda_{\min}(S^i) \geq 0$ to show S^i is positive semidefinite (PSD).

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

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

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

Thus, we have

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

We can write $S^i - S$ as

$$\begin{split} S^i - S &= (\nabla^2_{\theta\theta} U^i_{\tilde{\theta}} - \nabla^2_{\theta\theta} U_{\tilde{\theta}}) - \left[(\nabla^2_{\theta\tau} U^i_{\tilde{\theta}} - \nabla^2_{\theta\tau} U_{\tilde{\theta}}) (\nabla^2_{\tau\tau} U^i_{\tilde{\theta}})^{-1} \nabla^2_{\tau\theta} U^i_{\tilde{\theta}} \right. \\ &+ \nabla^2_{\theta\tau} U_{\tilde{\theta}} (\nabla^2_{\tau\tau} U^i_{\tilde{\theta}})^{-1} (\nabla^2_{\tau\theta} U^i_{\tilde{\theta}} - \nabla^2_{\tau\theta} U_{\tilde{\theta}}) + \nabla^2_{\theta\tau} U_{\tilde{\theta}} \Big((\nabla^2_{\tau\tau} U^i_{\tilde{\theta}})^{-1} - (\nabla^2_{\tau\tau} U_{\tilde{\theta}})^{-1} \Big) \nabla^2_{\tau\theta} U_{\tilde{\theta}} \Big]. \end{split}$$

Hence,

$$\begin{split} \|S^i - S\|_{\sigma} &\leq \|\nabla^2_{\theta\theta} U^i_{\bar{\theta}} - \nabla^2_{\theta\theta} U_{\bar{\theta}}\|_{\sigma} + \underbrace{\|(\nabla^2_{\theta\tau} U^i_{\bar{\theta}} - \nabla^2_{\theta\tau} U_{\bar{\theta}})(\nabla^2_{\tau\tau} U^i_{\bar{\theta}})^{-1} \nabla^2_{\tau\theta} U^i_{\bar{\theta}}\|_{\sigma}}_{T_1} \\ &+ \underbrace{\|\nabla^2_{\theta\tau} U_{\bar{\theta}} (\nabla^2_{\tau\tau} U^i_{\bar{\theta}})^{-1} (\nabla^2_{\tau\theta} U^i_{\bar{\theta}} - \nabla^2_{\tau\theta} U_{\bar{\theta}})\|_{\sigma}}_{T_2} \\ &+ \underbrace{\|\nabla^2_{\theta\tau} U_{\bar{\theta}} \left((\nabla^2_{\tau\tau} U^i_{\bar{\theta}})^{-1} - (\nabla^2_{\tau\tau} U_{\bar{\theta}})^{-1}\right) \nabla^2_{\tau\theta} U_{\bar{\theta}}\|_{\sigma}}_{T_3}. \end{split}$$

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

$$\nabla^2 U^i_{\tilde{\theta}} = \begin{bmatrix} \nabla^2_{\theta\theta} U^i_{\tilde{\theta}} & \nabla^2_{\theta\tau} U^i_{\tilde{\theta}} \\ \nabla^2_{\tau\theta} U^i_{\tilde{\theta}} & \nabla^2_{\tau\tau} U^i_{\tilde{\theta}} \end{bmatrix},$$

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

$$\|\nabla^2_{\theta\theta}U^i_{\tilde{\theta}}\|_{\sigma} \leq L, \quad \|\nabla^2_{\theta\tau}U^i_{\tilde{\theta}}\|_{\sigma} \leq L, \quad \|\nabla^2_{\tau\theta}U^i_{\tilde{\theta}}\|_{\sigma} \leq L, \quad \|\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}\|_{\sigma} \leq L.$$

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

$$T_{1} = \|(\nabla_{\theta\tau}^{2} U_{\tilde{\theta}}^{i} - \nabla_{\theta\tau}^{2} U_{\tilde{\theta}})(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} \nabla_{\tau\theta}^{2} U_{\tilde{\theta}}^{i} \|_{\sigma}$$

$$\leq \|(\nabla_{\theta\tau}^{2} U_{\tilde{\theta}}^{i} - \nabla_{\theta\tau}^{2} U_{\tilde{\theta}})\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1}\|_{\sigma} \cdot \|\nabla_{\tau\theta}^{2} U_{\tilde{\theta}}^{i}\|_{\sigma}$$

$$\leq L\rho_{\theta\tau}^{i} \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})|}$$

Similarly, bounding T_2

$$T_{2} = \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}}(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} (\nabla_{\tau\theta}^{2} U_{\tilde{\theta}}^{i} - \nabla_{\tau\theta}^{2} U_{\tilde{\theta}})\|_{\sigma}$$

$$\leq \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}}\|_{\sigma} \cdot \|(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1}\|_{\sigma} \cdot \|(\nabla_{\tau\theta}^{2} U_{\tilde{\theta}}^{i} - \nabla_{\tau\theta}^{2} U_{\tilde{\theta}})\|_{\sigma}$$

$$\leq L \rho_{\tau\theta}^{i} \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})|}$$

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

$$T_{3} = \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}} \Big((\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} - (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}})^{-1} \Big) \nabla_{\tau\theta}^{2} U_{\tilde{\theta}} \|_{\sigma}$$

$$\leq \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}} \|_{\sigma} \cdot \| (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} - (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}})^{-1} \|_{\sigma} \cdot \|\nabla_{\tau\theta}^{2} U_{\tilde{\theta}} \|_{\sigma}$$

$$= \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}} \|_{\sigma} \cdot \| (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}} - \nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i}) (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}})^{-1} \|_{\sigma} \cdot \|\nabla_{\tau\theta}^{2} U_{\tilde{\theta}} \|_{\sigma}$$

$$\leq \|\nabla_{\theta\tau}^{2} U_{\tilde{\theta}} \|_{\sigma} \cdot \| (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i})^{-1} \|_{\sigma} \cdot \|\nabla_{\tau\tau}^{2} U_{\tilde{\theta}} - \nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i} \|_{\sigma} \cdot \| (\nabla_{\tau\tau}^{2} U_{\tilde{\theta}})^{-1} \|_{\sigma} \cdot \|\nabla_{\tau\theta}^{2} U_{\tilde{\theta}} \|_{\sigma}$$

$$\leq L^{2} \rho_{\tau}^{i} \frac{1}{|\lambda_{\max}(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}}^{i}) \cdot \lambda_{\max}(\nabla_{\tau\tau}^{2} U_{\tilde{\theta}})|}$$

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

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

where $\lambda_{\max}(S) = \beta$ such that $\beta \geq 0$. Hence, we obtain $\left[\nabla^2_{\theta\theta}U^i_{\tilde{\theta}} - \nabla^2_{\theta\tau}U^i_{\tilde{\theta}}\left(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}\right)^{-1}\nabla^2_{\tau\theta}U^i_{\tilde{\theta}}\right](\hat{\theta},\hat{\tau}) \succeq (\beta - B^i)I.$ When $\beta \geq B^i$, then S^i is

positive semi-definite. When
$$B^i=0$$
, hence $\left[\nabla^2_{\theta\theta}U^i_{\tilde{\theta}}-\nabla^2_{\theta\tau}U^i_{\tilde{\theta}}\left(\nabla^2_{\tau\tau}U^i_{\tilde{\theta}}\right)^{-1}\nabla^2_{\tau\theta}U^i_{\tilde{\theta}}\right](\hat{\theta},\hat{\tau})\succeq\beta I$,

thus it will be positive semidefinite. When $\rho_{\tau}^{i} < \alpha$ and $\beta > B^{i}$, then the suuficient condition for ε^{i} -approximate equilibrium is satisfied. And we obtain the result.

Thus, for each client i, any approximation error ε^i that satisfies:

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

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

B.3 Consistency

B.3.1 ASSUMPTIONS

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

1014
1015 **Assumption 6** (Identification). θ_0 is the unique $\theta \in \Theta$ such that $\psi(f^i; \theta) = 0$ for all $f^i \in \mathcal{F}$, where $i \in [n]$.

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

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

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 variables for all $i \in [N]$.

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

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

B.3.2 Proof of Theorem 2

Theorem 6 (Restatement of 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 $m^i(\theta, \tau, \tilde{\theta}) := f^i(Z^i; \tau)(Y^i - g(X^i; \theta)) - \frac{1}{4}f^i(Z^i; \tau)^2(Y^i - g(X^i; \tilde{\theta}))^2$, $M^i(\theta) = \sup_{\tau \in \mathcal{T}} \mathbb{E}[m^i(\theta, \tau, \tilde{\theta})]$ and $\eta^i(\epsilon) := \inf_{d(\theta, \theta_0) \ge \epsilon} M^i(\theta) - M^i(\theta_0)$ for every $\epsilon > 0$. Let $(\hat{\theta}_n, \hat{\tau}_n)$ be a solution that satisfies the approximate equilibrium for each of the client $i \in [N]$ as

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

for some δ_0 , such that for any $\delta \in (0, \delta_0]$, and any θ, τ such that $\|\theta - \hat{\theta}\| \le \delta$ and $\|\tau - \hat{\tau}\| \le \delta$ and a function $h(\delta) \to 0$ as $\delta \to 0$. 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 θ_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 $i \in [N]$.

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

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

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

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

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

Note that $\tilde{\theta}_n$ is a data-dependent sequence for the global model. Practically, the previous global 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 $\tilde{\theta}_n \stackrel{p}{\to} \tilde{\theta}$.

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

$$\begin{split} \sup_{\theta} |M_{n_{i}}(\theta) - M^{i}(\theta)| &= \sup_{\theta} \left| \sup_{\tau \in \mathcal{T}} \mathbb{E}_{n_{i}}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] - \sup_{\tau \in \mathcal{T}} \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta})] \right| \\ &\leq \sup_{\theta, \tau} \left| \mathbb{E}_{n_{i}}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] - \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta})] \right| \\ &\leq \sup_{\theta, \tau} \left| \mathbb{E}_{n_{i}}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] - \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] \right| + \sup_{\theta, \tau} \left| \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] - \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta})] \right| \\ &\leq \sup_{\theta, \iota, \theta > \tau} \left| \mathbb{E}_{n_{i}}[m^{i}(\theta_{1}, \tau, \theta_{2})] - \mathbb{E}[m^{i}(\theta_{1}, \tau, \theta_{2})] \right| + \sup_{\theta, \tau} \left| \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta}_{n})] - \mathbb{E}[m^{i}(\theta, \tau, \tilde{\theta})] \right| \end{split}$$

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 $m^i_k(\theta,\tau,\tilde{\theta}_n)=f^i(Z^i_k;\tau)(Y^i_k-g^i(X^i_k;\theta))-\frac{1}{4}f^i(Z^i_k;\tau)^2(Y^i_k-g^i(X^i_k;\tilde{\theta}))^2$ and using $m^{i'}_k(\theta,\tau,\tilde{\theta}'_n)$ with ghost variables $\tilde{\theta}'_n$ for symmetrization and ϵ_k as k i.i.d. Rademacher random variables, we

obtain

$$\begin{split} \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}\left[m_{k}^{i}{}'(\theta_{1},\tau,\theta_{2}')\right] \right| \right] \\ &\leq \mathbb{E}\left[\sup_{\theta_{1},\theta_{2},\tau} \left| \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} \left(m_{k}^{i}(\theta_{1},\tau,\theta_{2}) - m_{k}^{i}{}'(\theta_{1},\tau,\theta_{2}')\right) \right| \right] \\ &\leq \mathbb{E}\left[\sup_{\theta_{1},\theta_{2},\tau} \left| \frac{1}{n_{i}} \sum_{k=1}^{n_{i}} \epsilon_{k} \left(m_{k}^{i}(\theta_{1},\tau,\theta_{2}) - m_{k}^{i}{}'(\theta_{1},\tau,\theta_{2}')\right) \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] \\ &+ \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] \\ &+ \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};\tilde{\theta}))^{4} \right| \right] \\ &+ \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 ϕ , $\mathcal{R}_{n_i}(\phi \circ f^i) \leq L\mathcal{R}_{n_i}(\mathcal{F}^i)$, 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, $\mathbb{E}[B_1] \to 0$. 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:

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

$$\le \frac{b}{n_i},$$

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

Now, we will handle B_2 . For that

$$B_{2} = \sup_{\theta,\tau} \left| \mathbb{E} \left[m^{i}(\theta,\tau,\tilde{\theta}_{n}) \right] - \mathbb{E} \left[m^{i}(\theta,\tau,\tilde{\theta}) \right] \right|$$

$$= \sup_{\theta,\tau} \left| \mathbb{E} \left[f^{i}(Z^{i};\tau)(Y^{i} - g(X^{i};\theta)) - \frac{1}{4} f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}_{n}))^{2} \right] \right|$$

$$- \mathbb{E} \left[f^{i}(Z^{i};\tau)(Y^{i} - g(X^{i};\theta)) - \frac{1}{4} f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}))^{2} \right] \right|$$

$$= \sup_{\theta,\tau} \frac{1}{4} \left| \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}_{n}))^{2} \right] - \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}))^{2} \right] \right|$$

$$= \sup_{\theta,\tau} \frac{1}{4} \left| \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}_{n}))^{2} \right] + \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}))^{2} \right] \right|$$

$$- \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}))^{2} \right] - \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}(Y^{i} - g(X^{i};\tilde{\theta}))^{2} \right] \right|$$

$$\leq \frac{1}{4} \sup_{\tau} \left| \mathbb{E} \left[f^{i}(Z^{i};\tau)^{2}\omega_{n} \right] \right|$$

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, $\tilde{\theta}_n \stackrel{p}{\to} \tilde{\theta}$, thus $\omega_n \stackrel{p}{\to} 0$ due to Slutsky's and continuous mapping theorem. Since, $f^i(Z; \tau)$ is uniformly bounded, thus for some constant b' > 0, we have

$$B_2 \le \frac{b'}{4} \sup_{\tau} \frac{1}{N} \sum_{i=1}^{N} |\mathbb{E} [\omega_n]|$$
$$\le \frac{b'}{4} \sup_{\tau} \frac{1}{N} \sum_{i=1}^{N} \mathbb{E} [|\omega_n|]$$

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

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 $\inf_{d(\theta,\theta_0)>\epsilon} M^i(\theta) > M^i(\theta_0)$.

 $M^i(\theta_0)$ is the unique minimizer of $M^i(\theta)$. By assumption (6) and (7), θ_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 θ will have at least one τ 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 $cf^{i'}$ for any $c\in(0,1)$. Using assumption (7), $cf^{i'}\in\mathcal{F}^i$, thus

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

$$\leq c \mathbb{E} \left[f^{i'}(Z^{i})(Y^{i} - g(X^{i}; \theta')) \right] - \frac{c^{2}}{4} \mathbb{E} \left[f^{i'}(Z^{i})^{2} (Y^{i} - g(X^{i}; \tilde{\theta}))^{2} \right]$$

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 θ_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 $M^i(\theta_n)\to M^i(\theta_0)$. By construction, $B(\theta_0,\epsilon)^c$ is closed and the corresponding limit parameters $\theta^*=\lim_{n\to\infty}\theta_n\in B(\theta_0,\epsilon)^c$ must satisfy $M^i(\theta^*)=M^i(\theta_0)$ using assumption (8).

But $d(\theta^*, \theta_0) \ge \epsilon > 0$, thus $\theta^* \ne \theta_0$. This contradicts that θ_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 ε^i - approximate equilibrium condition, given as:

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

for a function $h(\delta) \to 0$ as $\delta \to 0$ and some δ_0 , such that for any $\delta \in (0, \delta_0]$, and any θ, τ such that $\|\theta - \hat{\theta}\| \le \delta$ and $\|\tau - \hat{\tau}\| \le \delta$. Assume that this is true with $o_p(1)$, hence

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

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

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

Thus, we have

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

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

$$M^{i}(\hat{\theta}_{n}) - M^{i}(\theta_{0}) \leq M^{i}(\hat{\theta}_{n}) - M_{n_{i}}(\hat{\theta}_{n}) + 2\varepsilon^{i} + o_{p}(1)$$

$$\leq \sup_{\theta} |M^{i}(\hat{\theta}) - M_{n_{i}}(\hat{\theta})| + 2\varepsilon^{i} + o_{p}(1)$$

$$\leq 2\varepsilon^{i} + o_{p}(1)$$

Hence, we obtain

$$M^{i}(\hat{\theta}_{n}) - M^{i}(\theta_{0}) - 2\varepsilon^{i} \leq M^{i}(\hat{\theta}_{n}) - M_{n_{i}}(\hat{\theta}_{n}) + o_{p}(1)$$

$$\leq \sup_{\theta} |M^{i}(\hat{\theta}) - M_{n_{i}}(\hat{\theta})| + o_{p}(1)$$

$$\leq o_{p}(1)$$

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, $\mathbb{P}[d(\hat{\theta}_n,\theta_0)\geq \epsilon]\leq \mathbb{P}[M^i(\hat{\theta}_n)-M^i(\theta_0)\geq \eta^i(\epsilon)]=\mathbb{P}[M^i(\hat{\theta}_n)-M^i(\theta_0)-2\epsilon^i\geq \eta^i(\epsilon)-2\epsilon^i]$. For every $\epsilon>0$, we have $\eta^i(\epsilon)>0$ from claim 2, and $M^i(\hat{\theta}_n)-M^i(\theta_0)-2\epsilon^i=o_p(1)$. Thus, $\eta^i(\epsilon)-2\epsilon^i>0$ when $\epsilon^i<\frac{\eta^i(\epsilon)}{2}$. We have that for every $\epsilon>0$ and $\epsilon^i<\frac{\eta^i(\epsilon)}{2}$, the RHS probability converges to 0, thus $d(\hat{\theta}_n,\theta_0)=o_p(1)$, hence $\hat{\theta}_n$ converges in probability to θ_0 for each client $i\in[N]$.

C LIMIT POINTS OF FEDGDA

We first discuss the γ - FEDGDA flow.

C.1 FEDGDA FLOW

The FEDGDA updates can be written as

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

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

$$\lim_{\eta \to 0} \frac{\theta_{t+1} - \theta_t}{\eta} = \lim_{\eta \to 0} -\frac{1}{\gamma} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^{R} \left(\nabla_{\theta} U_{\tilde{\theta}}(\theta_t, \tau_t) + \left(\nabla_{\theta} U_{\tilde{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\theta} U_{\tilde{\theta}}^i(\theta_t, \tau_t) \right) + \left(\nabla_{\theta} U_{\tilde{\theta}}^i(\theta_t, \tau_t) - \nabla_{\theta} U_{\tilde{\theta}}(\theta_t, \tau_t) \right) \right)$$

$$\lim_{\eta \to 0} \frac{\tau_{t+1} - \tau_t}{\eta} = \lim_{\eta \to 0} \frac{1}{N} \sum_{i \in [N]} \sum_{r=1}^{R} \left(\nabla_{\tau} U_{\tilde{\theta}}(\theta_t, \tau_t) + \left(\nabla_{\tau} U_{\tilde{\theta}}^i(\theta_{t,r}^i, \tau_{t,r}^i) - \nabla_{\tau} U_{\tilde{\theta}}^i(\theta_t, \tau_t) \right) + \left(\nabla_{\tau} U_{\tilde{\theta}}^i(\theta_t, \tau_t) - \nabla_{\tau} U_{\tilde{\theta}}^i(\theta_t, \tau_t) \right) \right)$$

We obtain the gradient flow equations as

$$\frac{d\theta}{dt} = -\frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\theta} U_{\tilde{\theta}}(\theta(t), \tau(t)) \right) - \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\theta} U_{\tilde{\theta}}^{i}(\theta^{i}(t), \tau^{i}(t)) - \nabla_{\theta} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) \right)
- \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\theta} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) - \nabla_{\theta} U_{\tilde{\theta}}(\theta(t), \tau(t)) \right), \tag{26}$$

$$\frac{d\tau}{dt} = R \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\tau} U_{\tilde{\theta}}(\theta(t), \tau(t)) \right) + R \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\tau} U_{\tilde{\theta}}^{i}(\theta^{i}(t), \tau^{i}(t)) - \nabla_{\tau} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) \right)
+ R \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\tau} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) - \nabla_{\tau} U_{\tilde{\theta}}(\theta(t), \tau(t)) \right). \tag{27}$$

Using Assumption 3

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

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

Thus,

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

Since $U^i_{\tilde{\mathbf{a}}}$ is Lipschitz smooth by assumption 2, we have

$$\left\| \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\theta} U_{\tilde{\theta}}^{i}(\theta^{i}(t), \tau^{i}(t)) - \nabla_{\theta} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) \right) \right\| \leq L \frac{R}{\gamma} \frac{1}{N} \sum_{i \in [N]} \|(\theta^{i}(t), \tau^{i}(t)) - (\theta(t), \tau(t))\|,$$

$$\left\| R \frac{1}{N} \sum_{i \in [N]} \left(\nabla_{\tau} U_{\tilde{\theta}}^{i}(\theta^{i}(t), \tau^{i}(t)) - \nabla_{\tau} U_{\tilde{\theta}}^{i}(\theta(t), \tau(t)) \right) \right\| \leq L R \frac{1}{N} \sum_{i \in [N]} \|(\theta^{i}(t), \tau^{i}(t)) - (\theta(t), \tau(t))\|.$$

Substituting these bounds into Equations (26) and (27), we obtain

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

Since the local update follows

$$\theta^{i}(t) = \theta(t) - \frac{\eta}{\gamma} \sum_{j=1}^{R} \nabla_{\theta} U_{\bar{\theta}}^{i}(\theta_{j}^{i}(t), \tau_{j}^{i}(t)),$$

$$\tau^{i}(t) = \tau(t) + \eta \sum_{j=1}^{R} \nabla_{\tau} U_{\bar{\theta}}^{i}(\theta_{j}^{i}(t), \tau_{j}^{i}(t)),$$

Using bounded gradient assumption, i.e. $\|\nabla_{\theta}U_{\bar{\theta}}^{i}(\theta,\tau)\|^{2} \leq G_{\theta}$ and $\|\nabla_{\tau}U_{\bar{\theta}}^{i}(\theta,\tau)\|^{2} \leq G_{\tau}$ for all i, as $\eta \to 0$ and R is fixed and finite, the deviation $\|(\theta^{i}(t),\tau^{i}(t)) - (\theta(t),\tau(t))\|$ vanish, leading to

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

C.2 PROOF OF THEOREM 3

Proof. Let $A = \nabla^2_{\theta\theta} U_{\tilde{\theta}}(\theta, \tau)$, $B = \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta, \tau)$ and $C = \nabla^2_{\theta\tau} U_{\tilde{\theta}}(\theta, \tau)$. Consider $\epsilon = \frac{1}{\gamma}$, thus for sufficiently small ϵ (hence a large γ), the Jacobian J of FEDGDA for a point (θ, τ) is given as:

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

Using Lemma 5, J_{ϵ} has d_1+d_2 complex eigenvalues $\{\Lambda_j\}_{j=1}^{d_1+d_2}$ such that

$$|\Lambda_j + \epsilon \mu_j| = o(\epsilon) \qquad 1 \le j \le d_1$$

$$|\Lambda_{j+d_1} - \nu_j| = o(1), \qquad 1 \le j \le d_2,$$
(28)

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

We now prove the theorem statement:

By definition of \limsup and \liminf , we know that $\infty - \mathcal{FGDA} \subset \overline{\infty - \mathcal{FGDA}}$.

Now we show \mathcal{L} oc \mathcal{M} inimax $\subset \underline{\infty - \mathcal{F}\mathcal{G}\mathcal{D}\mathcal{A}}$. Consider a strict local minimax point (θ, τ) , then by sufficient condition it follows that:

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

Thus, $R\mathbf{B} < 0$, and $R(\mathbf{A} - C\mathbf{B}^{-1}C^{\top}) > 0$, where R is always positive. Hence, $\{\nu_j\}_{j=1}^{d_1} < 0$ and $\{\mu_j\}_{j=1}^{d_2} < 0$. Using equations 28, for some small $\epsilon_0 < \epsilon$, $\operatorname{Re}(\Lambda_j) < 0$ for all j. Thus, (θ, τ) is a strict linearly stable point of $\frac{1}{\epsilon}$ -FEDGDA.

Now, we show $\overline{\infty - \mathcal{FGDA}} \subset \mathcal{L}oc\mathcal{M}inimax \cup \{(\theta, \tau) | (\theta, \tau) \text{ is stationary and } \nabla^2_{\tau\tau} U_{\tilde{\theta}}(\theta, \tau) \text{ is degenerate}\}$. Consider (θ, τ) a strict linearly stable point of $\frac{1}{\epsilon}$ -FEDGDA, such that for some small ϵ , $\operatorname{Re}(\Lambda_i) < 0$ for all j. By equation 28, assuming B^{-1} exists

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

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

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

Thus, J_{ϵ} has 0 as its eigenvalue, which is a contradiction because for strict linearly stable point $\operatorname{Re}(\Lambda_j) < 0$ for all j. Thus, $A - CB^{-1}C^{\top} \succ 0$. Hence, (θ, τ) is a strict local minimax point.

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

We aim to show that the set

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

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: $\nabla G(\theta,\tau) = \operatorname{Tr}\left(\operatorname{Adj}(\nabla^2_{\tau\tau}U_{\tilde{\theta}})\cdot\nabla(\nabla^2_{\tau\tau}U_{\tilde{\theta}})\right)$. If $\det(\nabla^2_{\tau\tau}U_{\tilde{\theta}}(\theta,\tau))=0$, then the Hessian $\nabla^2_{\tau\tau}U_{\tilde{\theta}}$ 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 (θ, τ) 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: $G^{-1}(0) = \{(\theta,\tau) \mid \det(\nabla^2_{\tau\tau}U_{\tilde{\theta}}(\theta,\tau)) = 0\}$ has measure zero in $\mathbb R^{d+k}$.

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

Lemma 4. (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

$$|b_k - a_k| < \delta$$
 for $0 \le k \le n$, and $|b_k| < \delta$ for $n + 1 \le k \le m$,

then the roots β_k , $1 \le k \le m$, 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 α_k , $1 \le k \le n$, of $p_n(z)$, the inequalities

$$|\beta_k - \alpha_k| < \epsilon$$
 for $1 \le k \le n$, and $|\beta_k| > 1/\epsilon$ for $n + 1 \le k \le m$.

Lemma 5. For any symmetric matrix $A \in \mathbb{R}^{d_1 \times d_1}$, $B \in \mathbb{R}^{d_2 \times d_2}$, any rectangular matrix $C \in \mathbb{R}^{d_1 \times d_2}$ $\mathbb{R}^{d_1 imes d_2}$ and a scalar R, assume that B is non-degenerate. Then, matrix

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

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

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

$$\begin{split} |\Lambda_j + \epsilon \mu_j| &= o(\epsilon) \qquad 1 \leq j \leq d_1 \\ |\Lambda_{j+d_1} - \nu_j| &= o(1), \qquad 1 \leq j \leq d_2, \end{split}$$
 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 4 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 $R(\boldsymbol{A}-\boldsymbol{C}\boldsymbol{B}^{-1}\boldsymbol{C}^{\top})$ and $R\boldsymbol{B}$, respectively, and is thus

D RELATED WORK

1404

1405

1406

1407 1408

1409

1415

1416 1417 1418

1419 1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447 1448

1449

1450

1451 1452 1453

1454 1455

1456

1457

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

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

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

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

BENCHMARK CONSIDERATIONS AND ADDITIONAL EXPERIMENTS

THE EXPERIMENTAL BENCHMARK DESIGN

As stated, our experiments take the Bennett et al. (2019)'s experiments as a centralized-setting baseline. Therefore, we have used the same synthetic dataset as DEEPGMM, which they use in their experiments to benchmarks against the baselines therein such as DEEPIV (Hartford et al., 2017). It is standard to perform experimental analysis on synthetic datasets for unavailability of ground truth for causal inference; for example see Section 4.1.1 of Vo et al. (2022b). As the learning process essentially involves estimating the true parameter θ_0 by $\hat{\theta}$, to measure the performance of the learning procedure, we use the MSE of the estimate $\hat{g} := g(., \hat{\theta})$ against the true g_0 averaged over the clients. Nonetheless, an experimental comparison of our work with recent works on federated Bayesian methods for causal effect estimations does not apply directly. We discuss that below.

The two works in the domain of federated Bayesian methods for causal effect estimations are CAUSALRFF (Vo et al., 2022a) and FEDCI (Vo et al., 2022b). The aim of CAUSALRFF (Vo et al., 2022a) is to estimate the conditional average treatment effect (CATE) and average treatment effect (ATE), whereas FEDCI (Vo et al., 2022b) aims to estimate individual treatment effect (ITE) and ATE. For this, (Vo et al., 2022a) consider a setting of Y, W, and X to be random variables denoting the outcome, treatment, and proxy variable, respectively. Along with that, they also consider a confounding variable Z. However, their causal dependency builds on the dependence of each of Y, W, and X on Z besides dependency of Y on W. Consequently, to compute CATE and ATE, they need to estimate the conditional probabilities $p(w^i|x^i)$, $p(y^i|x^i,w^i)$, $p(z^i|x^i,y^i,w^i)$, $p(y^i|w^i,z^i)$, where the superscript i represents a client. Their experiments compare the estimates of CATE and ATE with the Bayesian baselines (Hill, 2011), (Shalit et al., 2017), (Louizos et al., 2017), etc. in a centralized setting without any consideration of data decentralization or heterogeneity native to federated learning. Further, they compare against the same baselines in a *one-shot federated* setting, where at the end of training on separate data sources independently, the predicted treatment effects are averaged. Similar is the experimental evaluation of (Vo et al., 2022b).

By contrast, the setting of IV analysis as in our work does not consider dependency of the outcome variable Y on the confounder Z, though the treatment variable X could be endogenous and depend on Z. For us, computing the treatment effects and thereby comparing it against these works is not direct. Furthermore, it is unclear, if the approach of (Vo et al., 2022a) and (Vo et al., 2022b), where the predicted inference over a number of datasets is averaged as the final result, would be comparable to our approach where the problem is solved 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

	$Dir_S(\alpha) = 0.1$		$Dir_S(\alpha) = 1.0$	
Estimations	FDEEPGMM-	FDEEPGMM-	FDEEPGMM-	FDEEPGMM-
	GDA	SGDA	GDA	SGDA
$FEMNIST_{x}$	0.27 ± 0.04	0.23 ± 0.02	0.17 ± 0.01	0.19 ± 0.03
$FEMNIST_{x,z}$	0.21 ± 0.01	0.24 ± 0.04	0.16 ± 0.03	0.18 ± 0.02
$\mathbf{FEMNIST_{z}}$	0.29 ± 0.02	0.25 ± 0.03	0.20 ± 0.04	0.23 ± 0.01
${ m CIFAR10_x}$	0.26 ± 0.01	0.27 ± 0.01	0.18 ± 0.01	0.15 ± 0.02
${ m CIFAR10_{x,z}}$	0.29 ± 0.02	0.30 ± 0.01	0.21 ± 0.02	0.13 ± 0.01
${ m CIFAR10_z}$	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.