

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Proof of Proposition 1

Fix any $\theta', \theta \in \mathbb{R}^d$. The optimality condition to (11) implies that

$$\sum_{i=1}^n \nabla f_i(\mathcal{M}(\theta); \theta) = \mathbf{0}, \quad \sum_{i=1}^n \nabla f_i(\mathcal{M}(\theta'); \theta') = \mathbf{0}. \quad (29)$$

Note that the gradients are taken w.r.t. the first argument in the function f_i . Observe the chain

$$0 = \langle \mathbf{0} | \mathcal{M}(\theta) - \mathcal{M}(\theta') \rangle = \left\langle \sum_{i=1}^n [\nabla f_i(\mathcal{M}(\theta); \theta) - \nabla f_i(\mathcal{M}(\theta'); \theta')] | \mathcal{M}(\theta) - \mathcal{M}(\theta') \right\rangle.$$

Adding and subtracting $\sum_{i=1}^n \nabla f_i(\mathcal{M}(\theta); \theta')$ implies the equality:

$$\begin{aligned} & \sum_{i=1}^n \langle \nabla f_i(\mathcal{M}(\theta); \theta') - \nabla f_i(\mathcal{M}(\theta); \theta) | \mathcal{M}(\theta) - \mathcal{M}(\theta') \rangle \\ & = \sum_{i=1}^n \langle (\nabla f_i(\mathcal{M}(\theta); \theta') - \nabla f_i(\mathcal{M}(\theta'); \theta')) | \mathcal{M}(\theta) - \mathcal{M}(\theta') \rangle. \end{aligned} \quad (30)$$

Applying A1 to the right hand side of (30) lead to:

$$\sum_{i=1}^n \langle (\nabla f_i(\mathcal{M}(\theta); \theta') - \nabla f_i(\mathcal{M}(\theta'); \theta')) | \mathcal{M}(\theta) - \mathcal{M}(\theta') \rangle \geq n\mu \|\mathcal{M}(\theta) - \mathcal{M}(\theta')\|^2.$$

Meanwhile, applying Lemma 2 to the left hand side of (30) gives

$$\begin{aligned} & \sum_{i=1}^n \langle \nabla f_i(\mathcal{M}(\theta); \theta') - \nabla f_i(\mathcal{M}(\theta); \theta) | \mathcal{M}(\theta) - \mathcal{M}(\theta') \rangle \\ & \leq \sum_{i=1}^n \epsilon_i L \|\theta' - \theta\| \|\mathcal{M}(\theta) - \mathcal{M}(\theta')\|. \end{aligned}$$

Substituting back into (30) implies that

$$\|\mathcal{M}(\theta) - \mathcal{M}(\theta')\| \leq \frac{\sum_{i=1}^n \epsilon_i L}{n\mu} \|\theta - \theta'\| = \frac{\epsilon_{\text{avg}} L}{\mu} \|\theta - \theta'\|. \quad (31)$$

Therefore, the map $\mathcal{M} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a contraction if $\epsilon_{\text{avg}} < \mu/L$. Subsequently, by the Banach fixed point theorem [Granas and Dugundji, 2003], the map $\mathcal{M}(\theta)$ admits a unique fixed point which is denoted as θ^{PS} .

To prove the converse, we consider the following instantiation of (11) with

$$\ell(\theta; Z) = \frac{1}{2}(\theta - Z)^2, \quad Z \sim \mathcal{D}_i(\theta) \iff Z \sim \mathcal{N}(\mu_i + \epsilon_i \theta, 1) \quad (32)$$

Note that the above satisfies A1 with $\mu = 1$, A2 with $L = 1$, A3 with ϵ_i for $i = 1, \dots, n$. We consider a case where it holds $\epsilon_{\text{avg}} \geq \mu/L = 1$. We also let $\mu_{\text{avg}} := (1/n) \sum_{i=1}^n \mu_i \neq 0$.

We observe

$$\begin{aligned} f_i(\theta'; \theta) & = \mathbb{E}_{Z \sim \mathcal{D}_i(\theta)} \left[\frac{1}{2}(\theta' - Z)^2 \right] = \mathbb{E}_{\tilde{Z} \sim \mathcal{N}(0,1)} \left[\frac{1}{2}(\theta' - \mu_i - \epsilon_i \theta - \tilde{Z})^2 \right] \\ & = \frac{1}{2}(\theta' - \mu_i - \epsilon_i \theta)^2 + \frac{1}{2}. \end{aligned} \quad (33)$$

For any $\theta \in \mathbb{R}$, it can be shown that

$$\mathcal{M}(\theta) = \arg \min_{\theta' \in \mathbb{R}} \frac{1}{2n} \sum_{i=1}^n (\theta' - \mu_i - \epsilon_i \theta)^2 = \epsilon_{\text{avg}} \theta + \mu_{\text{avg}} \quad (34)$$

Thus, applying the map for T times leads to

$$\mathcal{M}^T(\theta) = \epsilon_{\text{avg}}^T \theta + (1 + \epsilon_{\text{avg}} + \dots + \epsilon_{\text{avg}}^{T-1}) \mu_{\text{avg}} \quad (35)$$

Since $\epsilon_{\text{avg}} > 1$ and $\mu_{\text{avg}} \neq 0$, we have $\lim_{T \rightarrow \infty} |\mathcal{M}^T(\theta)| = \infty$ and the map is not a contraction.

B Proof of Lemma 3

Recall that $\tilde{\theta}^t := \bar{\theta}^t - \theta^{PS}$ is the error of averaged decision at the t th iteration. Using (21), we have

$$\left\| \tilde{\theta}^{t+1} \right\|^2 = \left\| \tilde{\theta}^t \right\|^2 - \frac{2\gamma_{t+1}}{n} \left\langle \tilde{\theta}^t \mid \sum_{i=1}^n \nabla \ell(\theta_i^t; Z_i^{t+1}) \right\rangle + \frac{\gamma_{t+1}^2}{n^2} \left\| \sum_{i=1}^n \nabla \ell(\theta_i^t; Z_i^{t+1}) \right\|^2. \quad (36)$$

We consider taking the conditional expectation $\mathbb{E}_t[\cdot]$ on the both sides. Using the fixed point condition $\sum_{i=1}^n \nabla f_i(\theta^{PS}; \theta^{PS}) = \mathbf{0}$, we observe the following equivalent expression for the last term

$$\left\| \sum_{i=1}^n \nabla \ell(\theta_i^t; Z_i^{t+1}) \right\|^2 = \left\| \sum_{i=1}^n [\nabla \ell(\theta_i^t; Z_i^{t+1}) - \nabla f_i(\theta_i^t; \theta_i^t) + \nabla f_i(\theta_i^t; \theta_i^t) - \nabla f_i(\theta^{PS}; \theta^{PS})] \right\|^2$$

Observe that $Z_i^{t+1}, i = 1, \dots, n$ are independent r.v.s, taking the conditional expectation $\mathbb{E}_t[\cdot]$ yields the upper bound to the above term

$$\begin{aligned} & \mathbb{E}_t \left\| \sum_{i=1}^n \nabla \ell(\theta_i^t; Z_i^{t+1}) \right\|^2 \\ & \leq 2 \sum_{i=1}^n \mathbb{E}_t \left\| \nabla \ell(\theta_i^t; Z_i^{t+1}) - \nabla f_i(\theta_i^t; \theta_i^t) \right\|^2 + 2n \sum_{i=1}^n \left\| \nabla f_i(\theta_i^t; \theta_i^t) - \nabla f_i(\theta^{PS}; \theta^{PS}) \right\|^2 \\ & \leq 2 \sum_{i=1}^n \sigma^2 (1 + \|\theta_i^t - \theta^{PS}\|^2) + 2n \sum_{i=1}^n L^2 (1 + \epsilon_i)^2 \|\theta_i^t - \theta^{PS}\|^2 \\ & \leq 2\sigma^2 n + 4n[\sigma^2 + nL^2(1 + \epsilon_{\max})^2] \left\| \tilde{\theta}^t \right\|^2 + 4[\sigma^2 + nL^2(1 + \epsilon_{\max})^2] \left\| \Theta_o^t \right\|_F^2 \end{aligned} \quad (37)$$

where the first inequality is due to A5 and Lemma 2. We conclude that

$$\frac{1}{n^2} \mathbb{E}_t \left\| \sum_{i=1}^n \nabla \ell(\theta_i^t; Z_i^{t+1}) \right\|^2 \leq \frac{2\sigma^2}{n} + c_2 \left\| \tilde{\theta}^t \right\|^2 + c_2 \frac{1}{n} \left\| \Theta_o^t \right\|_F^2 \quad (38)$$

where we recall the definition that $c_2 = 4 \left(\frac{\sigma^2}{n} + L^2(1 + \epsilon_{\max})^2 \right)$.

Next, we focus on the inner product term in (36), we have

$$\begin{aligned} \left\langle \tilde{\theta}^t \mid \sum_{i=1}^n \nabla f_i(\theta_i^t; \theta_i^t) \right\rangle &= \sum_{i=1}^n \left\langle \tilde{\theta}^t \mid \nabla f_i(\theta_i^t; \theta_i^t) - \nabla f_i(\bar{\theta}^t; \theta^{PS}) \right\rangle \\ &+ \sum_{i=1}^n \left\langle \tilde{\theta}^t \mid \nabla f_i(\bar{\theta}^t; \theta^{PS}) - \nabla f_i(\theta^{PS}; \theta^{PS}) \right\rangle \end{aligned} \quad (39)$$

Applying the Cauchy-Schwarz inequality and A2, A3, we obtain

$$\begin{aligned} \sum_{i=1}^n \left\langle \tilde{\theta}^t \mid \nabla f_i(\theta_i^t; \theta_i^t) - \nabla f_i(\bar{\theta}^t; \theta^{PS}) \right\rangle &\geq - \left\| \tilde{\theta}^t \right\| \sum_{i=1}^n \left(L \left\| \theta_i^t - \bar{\theta}^t \right\| + L\epsilon_i \left\| \theta_i^t - \theta^{PS} \right\| \right) \\ &\geq - \left\| \tilde{\theta}^t \right\| \sum_{i=1}^n \left(L(1 + \epsilon_i) \left\| \theta_i^t - \bar{\theta}^t \right\| + L\epsilon_i \left\| \tilde{\theta}^t \right\| \right). \end{aligned} \quad (40)$$

Meanwhile, using the strong convexity property of $\ell(\cdot; \cdot)$ [cf. A1], we have

$$\sum_{i=1}^n \left\langle \tilde{\theta}^t \mid \nabla f_i(\bar{\theta}^t; \theta^{PS}) - \nabla f_i(\theta^{PS}; \theta^{PS}) \right\rangle \geq n\mu \left\| \tilde{\theta}^t \right\|^2. \quad (41)$$

Summing up the two lower bounds and rearranging terms give

$$\frac{1}{n} \mathbb{E}_t \left\langle \tilde{\theta}^t \mid \sum_{i=1}^n \nabla f_i(\theta_i^t; \theta_i^t) \right\rangle \geq (\mu - L\epsilon_{\text{avg}}) \left\| \tilde{\theta}^t \right\|^2 - \frac{L}{n} (1 + \epsilon_{\max}) \sum_{i=1}^n \left\| \tilde{\theta}^t \right\| \left\| \theta_i^t - \bar{\theta}^t \right\|. \quad (42)$$

For any $\alpha > 0$, using the Young's inequality shows that the above can be further lower bounded by

$$\begin{aligned}
& \left[\mu - L\epsilon_{\text{avg}} - \frac{\alpha}{2n}L(1 + \epsilon_{\text{max}}) \right] \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 - \frac{L(1 + \epsilon_{\text{max}})}{2n\alpha} \sum_{i=1}^n \left\| \boldsymbol{\theta}_i^t - \bar{\boldsymbol{\theta}}^t \right\|^2 \\
& \geq \left[\mu - L\epsilon_{\text{avg}} - \frac{\alpha}{2n}L(1 + \epsilon_{\text{max}}) \right] \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 - \frac{L(1 + \epsilon_{\text{max}})}{2n\alpha} \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 \\
& \geq [\mu - (1 + \delta)L\epsilon_{\text{avg}}] \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 - \frac{L(1 + \epsilon_{\text{max}})^2}{4n^2\delta\epsilon_{\text{avg}}} \left\| \boldsymbol{\Theta}_o^t \right\|_F^2,
\end{aligned} \tag{43}$$

where we have set $\alpha = \frac{2n\delta\epsilon_{\text{avg}}}{1 + \epsilon_{\text{max}}}$ to yield the last inequality.

Substituting (38), (43) back to the inequality (36) gives us the desired result. In particular,

$$\begin{aligned}
\mathbb{E}_t \left\| \tilde{\boldsymbol{\theta}}^{t+1} \right\|^2 & \leq \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 - 2\gamma_{t+1} \left[[\mu - (1 + \delta)L\epsilon_{\text{avg}}] \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 - \frac{L(1 + \epsilon_{\text{max}})^2}{4n^2\delta\epsilon_{\text{avg}}} \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 \right] \\
& \quad + \gamma_{t+1}^2 \left[\frac{2\sigma^2}{n} + c_2 \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + c_2 \frac{1}{n} \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 \right] \\
& = (1 - 2\tilde{\mu}\gamma_{t+1} + c_2\gamma_{t+1}^2) \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + \left[c_1 \frac{\gamma_{t+1}}{n} + c_2 \frac{\gamma_{t+1}^2}{n} \right] \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + \frac{2\sigma^2}{n} \gamma_{t+1}^2 \\
& \leq (1 - \tilde{\mu}\gamma_{t+1}) \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + \left[c_1 \frac{\gamma_{t+1}}{n} + c_2 \frac{\gamma_{t+1}^2}{n} \right] \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + \frac{2\sigma^2}{n} \gamma_{t+1}^2
\end{aligned} \tag{44}$$

where we recall the constants $c_1 := \frac{L(1 + \epsilon_{\text{max}})^2}{2n\delta\epsilon_{\text{avg}}}$, $c_2 := 4 \left(\frac{\sigma^2}{n} + L^2(1 + \epsilon_{\text{max}})^2 \right)$ and $\tilde{\mu} := \mu - (1 + \delta)\epsilon_{\text{avg}}L$ and the last inequality is obtained by observing the condition $\gamma_{t+1} \leq \tilde{\mu}/c_2$.

C Proof of Lemma 4

To simplify notations, we denote

$$\begin{aligned}
\tilde{\nabla}F^t & := (\nabla\ell(\boldsymbol{\theta}_1^t; Z_1^{t+1}), \dots, \nabla\ell(\boldsymbol{\theta}_n^t; Z_n^{t+1}))^\top \in \mathbb{R}^{n \times d}, \\
\boldsymbol{\Theta}^t & := (\boldsymbol{\theta}_1^t, \dots, \boldsymbol{\theta}_n^t)^\top \in \mathbb{R}^{n \times d}, \quad \bar{\boldsymbol{\Theta}}^t := (1/n)\mathbf{1}\mathbf{1}^\top \boldsymbol{\Theta}^t \in \mathbb{R}^n.
\end{aligned} \tag{45}$$

Notice that $\boldsymbol{\Theta}_o^t = \boldsymbol{\Theta}^t - \bar{\boldsymbol{\Theta}}^t = (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top)\boldsymbol{\Theta}^t$. We first observe the following relation:

$$\begin{aligned}
\boldsymbol{\Theta}_o^{t+1} & = \boldsymbol{\Theta}^{t+1} - \bar{\boldsymbol{\Theta}}^{t+1} = \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) \boldsymbol{\Theta}^{t+1} = \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) (\mathbf{W}\boldsymbol{\Theta}^t - \gamma_{t+1}\tilde{\nabla}F^t) \\
& = \left(\mathbf{W} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) \boldsymbol{\Theta}_o^t - \gamma_{t+1} \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla}F^t,
\end{aligned}$$

where the last equality is due to $(\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top)\mathbf{W} = (\mathbf{W} - (1/n)\mathbf{1}\mathbf{1}^\top)(\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top)$ as \mathbf{W} is a doubly stochastic matrix.

Computing the squared norm of the consensus error leads to: for any $\alpha > 0$,

$$\begin{aligned}
\mathbb{E}_t \left\| \boldsymbol{\Theta}_o^{t+1} \right\|_F^2 & \leq (1 + \alpha)(1 - \rho)^2 \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + (1 + \frac{1}{\alpha})\gamma_{t+1}^2 \mathbb{E}_t \left\| \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla}F^t \right\|_F^2 \\
& \leq (1 - \rho) \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} \mathbb{E}_t \left\| \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla}F^t \right\|_F^2,
\end{aligned} \tag{46}$$

where we have applied A4 in the first inequality and set $\alpha = \frac{\rho}{1-\rho}$ in the second inequality. The last term in the above inequality can be bounded as

$$\begin{aligned}
\mathbb{E}_t \left\| \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla} F^t \right\|_F^2 &= \mathbb{E}_t \left[\sum_{i=1}^n \left\| \nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - \frac{1}{n} \sum_{j=1}^n \nabla \ell(\boldsymbol{\theta}_j^t; Z_j^{t+1}) \right\|^2 \right] \\
&\leq 3 \sum_{i=1}^n \mathbb{E}_t \left\| \nabla \ell(\boldsymbol{\theta}_i^t; Z_i^{t+1}) - \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) \right\|^2 + \frac{3}{n} \sum_{j=1}^n \mathbb{E}_t \left\| \nabla \ell(\boldsymbol{\theta}_j^t; Z_j^{t+1}) - \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) \right\|^2 \\
&\quad + 3 \sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) \right\|^2 \\
&\leq 6\sigma^2 \left(n + \sum_{i=1}^n \left\| \boldsymbol{\theta}_i^t - \boldsymbol{\theta}^{PS} \right\|^2 \right) + 3 \sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) \right\|^2 \\
&\leq 6\sigma^2 \left(n + 2n \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + 2 \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 \right) + 3 \sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) \right\|^2
\end{aligned} \tag{47}$$

where the second last inequality is due to A5. For each $i = 1, \dots, n$, we observe

$$\begin{aligned}
&\left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) + \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) - \frac{1}{n} \sum_{j=1}^n [\nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) - \nabla f_j(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t)] \right\|^2 \\
&\leq 3 \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 + 3 \left\| \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 \\
&\quad + \frac{3}{n} \sum_{j=1}^n \left\| \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) - \nabla f_j(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 \\
&\leq 3 \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 + \frac{3}{n} \sum_{j=1}^n \left\| \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) - \nabla f_j(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 + 3\varsigma^2 \left(1 + \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 \right)
\end{aligned} \tag{48}$$

where the last inequality is due to A6. Now, we observe

$$\begin{aligned}
\sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\boldsymbol{\theta}_j^t, \boldsymbol{\theta}_j^t) \right\|^2 &\leq 6 \sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{\theta}_i^t, \boldsymbol{\theta}_i^t) - \nabla f_i(\bar{\boldsymbol{\theta}}^t, \bar{\boldsymbol{\theta}}^t) \right\|^2 + 3n\varsigma^2 \left(1 + \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 \right) \\
&\leq 6L^2(1 + \epsilon_{\max})^2 \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + 3n\varsigma^2 \left(1 + \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 \right)
\end{aligned}$$

where the second inequality is due to Lemma 2 and the definition of $\boldsymbol{\Theta}_o^t$.

Substituting the above bounds into (47) leads to

$$\begin{aligned}
\mathbb{E}_t \left\| \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla} F^t \right\|_F^2 &\leq 6\sigma^2 \left(n + 2n \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + 2 \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 \right) + 18L^2(1 + \epsilon_{\max})^2 \left\| \boldsymbol{\Theta}_o^t \right\|_F^2 + 9n\varsigma^2 \left(1 + \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 \right) \\
&\leq 9n[\sigma^2 + \varsigma^2] + 12n[\sigma^2 + \varsigma^2] \left\| \tilde{\boldsymbol{\theta}}^t \right\|^2 + [12\sigma^2 + 18L^2(1 + \epsilon_{\max})^2] \left\| \boldsymbol{\Theta}_o^t \right\|_F^2
\end{aligned} \tag{49}$$

Let $c_3 := 12\sigma^2 + 18L^2(1 + \epsilon_{\max})^2$. Substituting the above inequality into (46) gives us

$$\begin{aligned} \mathbb{E}_t \|\Theta_o^{t+1}\|_F^2 &\leq (1 - \rho) \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} \left(12n[\sigma^2 + \varsigma^2] \|\tilde{\theta}^t\|^2 + c_3 \|\Theta_o^t\|_F^2 \right) + 9n(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho} \\ &\leq (1 - \rho/2) \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} 12n[\sigma^2 + \varsigma^2] \|\tilde{\theta}^t\|^2 + 9n(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho}, \end{aligned}$$

where the last inequality is due to the step size condition $\gamma_{t+1}^2 \leq \rho^2/2c_3$. The proof is concluded.

Alternative Bound without A6 We consider bounding (48) without using A6. Instead, we only assume that $\max_{i=1, \dots, n} \|\nabla f_i(\theta^{PS}; \theta^{PS})\|^2 \leq \varsigma^2$. We observe

$$\begin{aligned} \left\| \nabla f_i(\bar{\theta}^t, \bar{\theta}^t) - \nabla f(\bar{\theta}^t, \bar{\theta}^t) \right\|^2 &\leq 2 \left\| \nabla f_i(\theta^{PS}; \theta^{PS}) \right\|^2 \\ &\quad + 2 \left\| \nabla f_i(\bar{\theta}^t, \bar{\theta}^t) - \nabla f_i(\theta^{PS}; \theta^{PS}) + \nabla f(\theta^{PS}; \theta^{PS}) - \nabla f(\bar{\theta}^t, \bar{\theta}^t) \right\|^2 \\ &\leq 2 \left\| \nabla f_i(\theta^{PS}; \theta^{PS}) \right\|^2 + 8L^2(1 + \epsilon_{\max})^2 \|\tilde{\theta}^t\|^2 \leq 2\varsigma^2 + 8L^2(1 + \epsilon_{\max})^2 \|\tilde{\theta}^t\|^2, \end{aligned} \quad (50)$$

for all $i = 1, \dots, n$. This leads to

$$\begin{aligned} &\sum_{i=1}^n \left\| \nabla f_i(\theta_i^t, \theta_i^t) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(\theta_j^t, \theta_j^t) \right\|^2 \\ &\leq 6L^2(1 + \epsilon_{\max})^2 \|\Theta_o^t\|_F^2 + 2n\varsigma^2 + 8nL^2(1 + \epsilon_{\max})^2 \|\tilde{\theta}^t\|^2. \end{aligned}$$

Subsequently,

$$\begin{aligned} &\mathbb{E}_t \left\| \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla} F^t \right\|_F^2 \\ &\leq 6\sigma^2 \left(n + 2n \|\tilde{\theta}^t\|^2 + 2 \|\Theta_o^t\|_F^2 \right) + 6n\varsigma^2 + 6L^2(1 + \epsilon_{\max})^2 \left(3 \|\Theta_o^t\|_F^2 + 4n \|\tilde{\theta}^t\|^2 \right) \\ &= 6n[\sigma^2 + \varsigma^2] + 12n \left[\sigma^2 + 2L^2(1 + \epsilon_{\max})^2 \right] \|\tilde{\theta}^t\|^2 + [12\sigma^2 + 18L^2(1 + \epsilon_{\max})^2] \|\Theta_o^t\|_F^2 \end{aligned} \quad (51)$$

Taking $c_3 := 12\sigma^2 + 18L^2(1 + \epsilon_{\max})^2$ as before and substituting the inequality into (46) yields

$$\begin{aligned} &\mathbb{E}_t \|\Theta_o^{t+1}\|_F^2 \\ &\leq (1 - \rho) \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} \left(12n \left[\sigma^2 + 2L^2(1 + \epsilon_{\max})^2 \right] \|\tilde{\theta}^t\|^2 + c_3 \|\Theta_o^t\|_F^2 \right) + 6n(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho} \\ &\leq (1 - \rho/2) \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} 12n \left[\sigma^2 + 2L^2(1 + \epsilon_{\max})^2 \right] \|\tilde{\theta}^t\|^2 + 6n(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho}, \end{aligned}$$

where the last inequality is due to $\sup_{t \geq 1} \gamma_t \leq \rho/\sqrt{2c_3}$. The above can be simplified into

$$\frac{1}{n} \mathbb{E}_t \|\Theta_o^{t+1}\|_F^2 \leq \left(1 - \frac{\rho}{2} \right) \frac{1}{n} \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} 12 \left[\sigma^2 + 2L^2(1 + \epsilon_{\max})^2 \right] \|\tilde{\theta}^t\|^2 + 6(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho}. \quad (52)$$

Compared to (23), we observe that the above bound entails a larger coefficient for $\|\tilde{\theta}^t\|^2$ which lead to a (slightly) worse convergence bound for the DSGD-GD scheme.

Lastly, we should mention that as in the original Lemma 4, (52) can also be combined with Lemma 3 to develop an alternate version of Lemma 5. Subsequently, we can achieve a similar result as Theorem 1 without assuming A6.

D Proof of Lemma 5

Combining Lemmas 3 and 4 leads to

$$\mathcal{L}_{t+1} \leq (1 - \tilde{\mu}\gamma_{t+1}) \mathbb{E} \|\tilde{\theta}^t\|^2 + [c_1\gamma_{t+1} + c_2\gamma_{t+1}^2] \frac{1}{n} \mathbb{E} \|\Theta_o^t\|_F^2 + \frac{2\sigma^2}{n} \gamma_{t+1}^2$$

$$\begin{aligned}
& + \gamma_{t+1} \frac{8c_1}{\rho} \left(\left(1 - \frac{\rho}{2}\right) \frac{1}{n} \mathbb{E} \|\Theta_o^t\|_F^2 + \frac{\gamma_{t+1}^2}{\rho} 12[\sigma^2 + \varsigma^2] \mathbb{E} \|\tilde{\theta}^t\|^2 + 9(\sigma^2 + \varsigma^2) \frac{\gamma_{t+1}^2}{\rho} \right) \\
& = \left(1 - \tilde{\mu}\gamma_{t+1} + \frac{96c_1}{\rho^2} [\sigma^2 + \varsigma^2] \gamma_{t+1}^3\right) \mathbb{E} \|\tilde{\theta}^t\|^2 + \frac{2\sigma^2}{n} \gamma_{t+1}^2 + \frac{72c_1}{\rho^2} (\sigma^2 + \varsigma^2) \gamma_{t+1}^3 \\
& + \gamma_t \frac{8c_1}{\rho} \left(\frac{\gamma_{t+1}}{\gamma_t} \left(1 - \frac{\rho}{2}\right) + \frac{\rho}{8} + \frac{c_2\rho}{8c_1} \gamma_{t+1} \right) \frac{1}{n} \mathbb{E} \|\Theta_o^t\|_F^2
\end{aligned}$$

Note that by the step size conditions specified in the lemma, we have

$$\begin{aligned}
1 - \tilde{\mu}\gamma_{t+1} + \frac{96c_1}{\rho^2} [\sigma^2 + \varsigma^2] \gamma_{t+1}^3 &\leq 1 - \tilde{\mu}\gamma_{t+1}/2 \\
\frac{\gamma_{t+1}}{\gamma_t} \left(1 - \frac{\rho}{2}\right) + \frac{\rho}{8} + \frac{c_2\rho}{8c_1} \gamma_{t+1} &\leq 1 - \tilde{\mu}\gamma_{t+1}/2.
\end{aligned} \tag{53}$$

Thus, we obtain

$$\mathcal{L}_{t+1} \leq (1 - \tilde{\mu}\gamma_{t+1}/2) \mathcal{L}_t + \frac{2\sigma^2}{n} \gamma_{t+1}^2 + \frac{72c_1}{\rho^2} (\sigma^2 + \varsigma^2) \gamma_{t+1}^3. \tag{54}$$

This concludes the first part of the lemma, i.e., (25). For the second part, we further expand (54) to obtain

$$\begin{aligned}
\mathcal{L}_{t+1} &\leq \prod_{i=1}^{t+1} \left(1 - \frac{\tilde{\mu}\gamma_i}{2}\right) \mathsf{D} + \sum_{s=1}^{t+1} \prod_{i=s+1}^{t+1} (1 - \tilde{\mu}\gamma_i/2) \left(\frac{2\sigma^2}{n} \gamma_s^2 + \frac{72c_1}{\rho^2} (\sigma^2 + \varsigma^2) \gamma_s^3 \right) \\
&\leq \prod_{i=1}^{t+1} \left(1 - \frac{\tilde{\mu}\gamma_i}{2}\right) \mathsf{D} + \frac{288c_1(\sigma^2 + \varsigma^2)}{\rho^2 \tilde{\mu}} \gamma_{t+1}^2 + \frac{8\sigma^2}{\tilde{\mu}n} \gamma_{t+1}.
\end{aligned} \tag{55}$$

where we recall that $\mathsf{D} := \|\tilde{\theta}^0\|^2 + \frac{8\gamma_1 c_1}{\rho n} \|\Theta_o^0\|_F^2$ and the last inequality is due to Lemma 6 together with the specified step size conditions. The proof is thus concluded.

E Auxilliary Results

Lemma 6. Consider a sequence of non-negative, non-increasing step sizes $\{\gamma_t\}_{t \geq 1}$. Let $a > 0$, $p \in \mathbb{Z}_+$ and $\gamma_1 < 2/a$. If $\gamma_t^p / \gamma_{t+1}^p \leq 1 + (a/2) \gamma_{t+1}^p$ for any $t \geq 1$, then

$$\sum_{j=1}^t \gamma_j^{p+1} \prod_{\ell=j+1}^t (1 - \gamma_\ell a) \leq \frac{2}{a} \gamma_t^p, \quad \forall t \geq 1. \tag{56}$$

Proof. Observe that:

$$\begin{aligned}
\sum_{j=1}^t \gamma_j^{p+1} \prod_{\ell=j+1}^t (1 - \gamma_\ell a) &= \gamma_t^p \sum_{j=1}^t \gamma_j \prod_{\ell=j+1}^t \frac{\gamma_{\ell-1}^p}{\gamma_\ell^p} (1 - \gamma_\ell a) \\
&\stackrel{(a)}{\leq} \gamma_t^p \sum_{j=1}^t \gamma_j \prod_{\ell=j+1}^t \left(1 - \gamma_\ell \frac{a}{2}\right) \\
&= \frac{2\gamma_t^p}{a} \sum_{j=1}^t \left(\prod_{\ell=j+1}^t (1 - \gamma_\ell a/2) - \prod_{\ell'=j}^t (1 - \gamma_{\ell'} a/2) \right) \\
&= \frac{2\gamma_t^p}{a} \left(1 - \prod_{\ell'=1}^t (1 - \gamma_{\ell'} a/2)\right) \leq \frac{2\gamma_t^p}{a},
\end{aligned}$$

where (a) is due to the following observation

$$\frac{\gamma_{\ell-1}^p}{\gamma_\ell^p} (1 - \gamma_\ell a) \leq \left(1 + \frac{a}{2} \gamma_\ell^p\right) (1 - \gamma_\ell a) \leq 1 - \frac{a}{2} \gamma_\ell.$$

The proof is concluded. \square

Tasks	$\bar{\epsilon} = \epsilon_{\text{avg}}$	a_0	a_1	batch
Gaussian Mean Estimation	see §4	50	10000	1
Spam Email Classification	see §4	50	100000	32
LEAF Synthetic Data (Hetero & Homo)	0.1	200	1000	32
LEAF Synthetic Data (Hetero & Homo)	10.0	1	1000	32

Table 1: Parameters for the numerical experiments.

Lemma 7. Consider a sequence of non-negative, non-increasing step sizes $\{\gamma_t\}_{t \geq 1}$. Let $p \in \mathbb{Z}^+$. If $\sup_{t \geq 1} \gamma_t^p / \gamma_{t+1}^p \leq 1 + \frac{\rho}{4-2\rho}$, then for any $t \geq 0$, it holds that

$$\sum_{i=1}^{t+1} \left(1 - \frac{\rho}{2}\right)^{t+1-i} \gamma_i^p \leq \frac{4}{\rho} \gamma_{t+1}^p. \quad (57)$$

Proof. We observe the following chain:

$$\begin{aligned} \sum_{i=1}^{t+1} \left(1 - \frac{\rho}{2}\right)^{t+1-i} \gamma_i^p &= \gamma_{t+1}^p \sum_{i=1}^{t+1} \left(1 - \frac{\rho}{2}\right)^{t+1-i} \left(\frac{\gamma_i}{\gamma_{i+1}}\right)^p \left(\frac{\gamma_{i+1}}{\gamma_{i+2}}\right)^p \cdots \left(\frac{\gamma_t}{\gamma_{t+1}}\right)^p \\ &\leq \gamma_{t+1}^p \sum_{i=1}^{t+1} \left(1 - \frac{\rho}{4}\right)^{t+1-i} \leq \frac{4}{\rho} \gamma_{t+1}^p \end{aligned}$$

where the second last inequality is due to:

$$\left(1 - \frac{\rho}{2}\right) \left(\frac{\gamma_{i+1}}{\gamma_{i+2}}\right)^p \leq 1 - \frac{\rho}{4}$$

since $\sup_{k \geq 1} \gamma_{k-1}^p / \gamma_k^p \leq 1 + \frac{\rho}{4-2\rho}$. This completes the proof. \square

F Details of Numerical Experiments and Additional Results

This section provides details for the numerical experiments conducted in §4. We also describe an additional numerical experiment based on the logistic regression Example 1 on synthetic data. The latter examines the effects of heterogeneous data on the convergence rate of DSGD-GD.

For all our experiments, we have performed DSGD-GD with the step size $\gamma_t = a_0 / (a_1 + t)$. Moreover, at each iteration of DSGD-GD, the i th agent draws $\text{batch} \geq 1$ samples from $\mathcal{D}_i(\theta_i^t)$. The parameters $a_0 > 0, a_1 \geq 0, \text{batch} \geq 1$ used for different tasks are specified in Table 1. For both Gaussian mean estimation and spambase logistic regression, we use the same parameters for all settings of $\bar{\epsilon} = \epsilon_{\text{avg}}$. We consider using $n = 25$ agents in all experiments, connected on a ring graph. We set the mixing matrix weights as $W_{ij} = 1/3$ for all $(i, j) \in E$, and $W_{ij} = 0$ if $(i, j) \notin E$.

Spam Email Classification. In Fig. 4, we provide additional results for the experiment in the main paper [cf. Fig. 3]. In particular, we compare the training loss $f(\bar{\theta}^t; \bar{\theta}^t)$ and training accuracy against the iteration number t . We also plot the gap to an *approximate* Multi-PS solution in Fig. 4 (right) for $\|\bar{\theta}^t - \theta^{\hat{P}S}\|^2$. Note that the Multi-PS solution compared here is only an approximation obtained by applying a similar method to repeated gradient descent in [Perdomo et al., 2020] on $\min_{\theta} \sum_{i=1}^n f_i(\theta; \theta)$, where we used 1000 gradient descent iterations together with an outer loop of 10^4 deployments. Note that this process is only guaranteed to find a near-optimal solution, denoted as $\theta^{\hat{P}S}$. Nevertheless, we observe that when the decision dependent distributions becomes more sensitive ($\epsilon_{\text{avg}} = 1$), the DSGD-GD scheme seems unable to reach $\theta^{\hat{P}S}$.

Logistic Regression on LEAF Synthetic Data. To study the effect of homogeneity of data distribution [cf. A6] on the convergence of DSGD-GD, we conduct an additional experiment based on Example 1 but on the LEAF synthetic data [Caldas et al., 2019]. Here, we set the sensitivity parameter at $\epsilon_i = \bar{\epsilon}$ for $i = 1, \dots, 25$ and generate synthetic data using the framework in [Caldas et al., 2019] with the

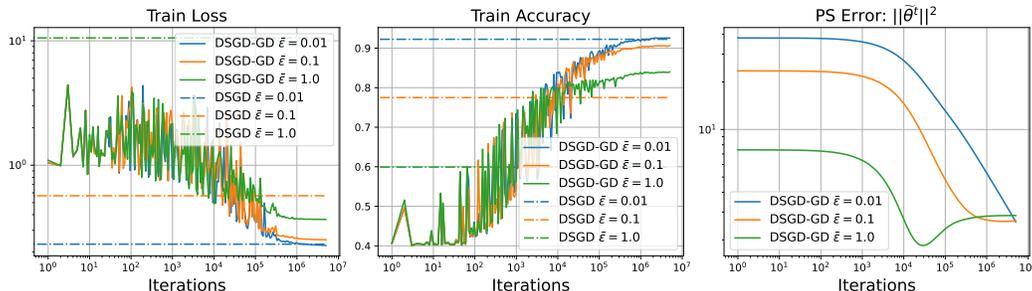


Figure 4: **Additional Results for Spam Email Classification.** (Left) Training Loss. (Middle) Training Accuracy. (Right) Approximate Gap to Multi-PS solution (see below). We also compare the non-performative optimal solution (dashed lines) on the shifted dataset.

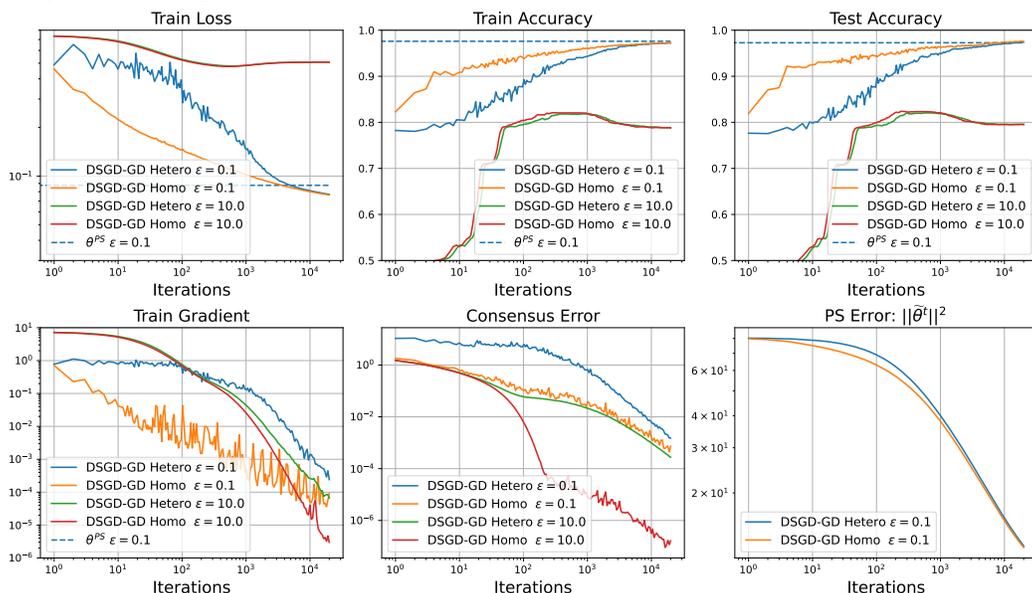


Figure 5: **Logistic Regression on LEAF Synthetic Data.** DSGD-GD in homogeneous and heterogeneous data distribution converge to the same Multi-PS solution.

standard deviation $\sigma = 1$ that represents the degree of heterogeneity of the dataset. Note that the framework produces $m_i = 100$ training samples with $d = 100$ features for each agent, denoted as $(\mathbf{X}_k^i, Y_k^i)_{k=1}^{100}$, for $i = 1, \dots, 25$ agents.

We consider two settings and describe them using the notations as in Example 1. In the heterogeneous data setting, the base data distribution \mathcal{D}_i^0 for agent i is taken to be $(\mathbf{X}_k^i, Y_k^i)_{k=1}^{100}$ such that $\mathcal{D}_i(\theta) \neq \mathcal{D}_j(\theta)$. In the homogeneous data setting, the base data distribution \mathcal{D}_i^0 for agent i is taken to be $((\mathbf{X}_k^i, Y_k^i)_{k=1}^{100})_{i=1}^{25}$, i.e., the entire dataset generated from LEAF. Note that in this case, $\mathcal{D}_i^0 \equiv \mathcal{D}_j^0$ and thus $\mathcal{D}_i(\theta) \equiv \mathcal{D}_j(\theta)$ for any $\theta \in \mathbb{R}^d$ and $i, j = 1, \dots, n$ since $\epsilon_i = \epsilon_{\text{avg}}$. Note that the Multi-PS solution θ^{PS} (if exists) in both settings are unique and identical. Meanwhile, the homogeneous case satisfies A6 with $\zeta = 0$, thus the DSGD-GD scheme applied to it is expected to converge at a faster rate than in the heterogeneous case.

Our numerical results are presented in Fig. 5, and we show in Table 1 the simulation parameters. Observe that with $\epsilon_{\text{avg}} = 10$, the local data distributions are too sensitive and the Multi-PS solution θ^{PS} may not exist. With $\epsilon_{\text{avg}} = 0.1$, we observe that the convergence of test accuracy, training loss, etc. are faster with the homogeneous case initially. However, as the iteration number t grows, the gap between the homogeneous and heterogeneous cases fade. This corroborates with our finite-time analysis in (18), where the fluctuation term $\sigma^2 \gamma_t / (n \bar{\mu})$ becomes dominant as $t \gg 1$ in all cases, yet the transient time can be shorter when $\zeta = 0$ as predicted by (19).

G Extension to Time-varying Graph

This section shows how to extend our analysis for DSGD-GD to the setting with time-varying communication graph. Let $G^{(t)} = (V, E^{(t)})$ be a simple, undirected graph which is possibly not connected and the graph is associated with a weighted adjacency matrix $\mathbf{W}^{(t)}$. Note that the graph $G^{(t)}$ consists of a fixed set of agents V and a set of time-varying edges $E^{(t)}$.

In lieu of A4, we assume that:

A7. *The time-varying undirected graph sequence $\{G^{(t)}\}_{t \geq 1} = \{(V, E^{(t)})\}_{t \geq 1}$ is B -connected. Specifically, for any $t \geq 1$, there exists a positive integer B such that the undirected graph $(V, E^{(t)} \cup \dots \cup E^{(t+B-1)})$ is connected. For any $t \geq 1$, the mixing matrix $\mathbf{W}^{(t)} \in \mathbb{R}^{n \times n}$ satisfies:*

1. (Topology) $\mathbf{W}_{ij}^{(t)} = 0$ if $(i, j) \notin E^{(t)}$.
2. (Doubly stochastic) $\mathbf{W}^{(t)}\mathbf{1} = (\mathbf{W}^{(t)})^\top \mathbf{1} = \mathbf{1}$.
3. (Fast mixing) Let $\mathbf{A}^{(t)} := \mathbf{W}^{(t)} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$, there exists $\bar{\rho} \in (0, 1]$ such that $\|\mathbf{A}^{(t+B-1)} \dots \mathbf{A}^{(t)}\|_2 \leq 1 - \bar{\rho}$.

The last condition can be guaranteed under the bounded communication setting, i.e., when the combined graph $(V, E^{(t)} \cup \dots \cup E^{(t+B)})$ is connected for any $t \geq 0$.

Notations. Throughout, we denote $\Theta(m, n) := \mathbb{E}[\|\Theta_o^m\|_F^2 + \dots + \|\Theta_o^n\|_F^2]$ and $\tilde{\theta}(m, n) := \mathbb{E}[\|\tilde{\theta}^m\|^2 + \dots + \|\tilde{\theta}^n\|^2]$, which is the aggregation of consensus error and performative stable gap in one time block whose length is B , respectively.

Proof Sketch. Below we provide a proof sketch for the convergence of DSGD-GD scheme when the latter is applied on a time varying graph satisfying A7. We begin by considering the extensions of Lemmas 3 and 4. As follows,

Lemma 8 (Extension of Lemma 3). *Fix any $\delta > 0$ and let $\epsilon_{\text{avg}} \leq \frac{\mu}{(1+\delta)L}$. Under A1, A2, A3, A5 and let the step sizes satisfy $\sup_{t \geq 0} \gamma_{t+1} \leq \frac{\tilde{\mu}}{c_2}$, the following bound holds for any $t \geq 0$,*

$$\begin{aligned} \tilde{\theta}(t+1, t+B) &\leq (1 - \tilde{\mu}\gamma_{t+B})^B \tilde{\theta}(t-B+1, t) + \frac{2B\sigma^2}{n} \gamma_{t+1}^2 \\ &\quad + B \left(c_1 \frac{\gamma_{t+1}}{n} + c_2 \frac{\gamma_{t+1}^2}{n} \right) [\Theta(t-B+1, t) + \Theta(t+1, t+B)]. \end{aligned}$$

Proof. Recall the inequality (22) in Lemma 3,

$$\mathbb{E}_t \left\| \tilde{\theta}^{t+1} \right\|^2 \leq (1 - \tilde{\mu}\gamma_{t+1}) \left\| \tilde{\theta}^t \right\|^2 + [c_1\gamma_{t+1} + c_2\gamma_{t+1}^2] \frac{1}{n} \|\Theta_o^t\|_F^2 + \frac{2\sigma^2}{n} \gamma_{t+1}^2. \quad (58)$$

This implies

$$\tilde{\theta}(t+1, t+B) \leq (1 - \tilde{\mu}\gamma_{t+B}) \tilde{\theta}(t, t+B-1) + \left(\frac{c_1\gamma_{t+1}}{n} + \frac{c_2\gamma_{t+1}^2}{n} \right) \Theta(t, t+B-1) + \frac{2B\sigma^2}{n} \gamma_{t+1}^2,$$

where we have summed (58) from $t+1$ th to $t+B$ th iteration and noted that the step size γ_t is non-increasing. Applying the above inequality for B times, we can link two consecutive B performative stable gap $\tilde{\theta}(t+1, t+B)$ and $\tilde{\theta}(t-B+1, t)$ by

$$\begin{aligned} \tilde{\theta}(t+1, t+B) &\leq (1 - \tilde{\mu}\gamma_{t+B})^B \tilde{\theta}(t-B+1, t) + \frac{2B\sigma^2}{n} \gamma_{t+1}^2 \\ &\quad + \left(\frac{c_1\gamma_{t+1}}{n} + \frac{c_2\gamma_{t+1}^2}{n} \right) [\Theta(t, t+B-1) + \Theta(t-1, t+B-2) + \dots + \Theta(t-B+1, t)]. \end{aligned} \quad (59)$$

For the first term $\Theta(t, t+B-1)$ in the last quantity, we observe the crude bound

$$\Theta(t, t+B-1) \leq \Theta(t+1, t+B) + \mathbb{E}\|\Theta_o^t\|_F^2 \leq \Theta(t-B+1, t) + \Theta(t+1, t+B).$$

Following the same trick, we get another crude bound as

$$\begin{aligned} & [\Theta(t, t+B-1) + \Theta(t-1, t+B-2) + \dots + \Theta(t-B+1, t)] \\ & \leq B[\Theta(t-B+1, t) + \Theta(t+1, t+B)]. \end{aligned}$$

Substituting back to inequality (59) derives the final bound

$$\begin{aligned} \tilde{\theta}(t+1, t+B) & \leq (1 - \tilde{\mu}\gamma_{t+B})^B \tilde{\theta}(t-B+1, t) + \frac{2B\sigma^2}{n} \gamma_{t+1}^2 \\ & \quad + B \left(c_1 \frac{\gamma_{t+1}}{n} + c_2 \frac{\gamma_{t+1}^2}{n} \right) [\Theta(t-B+1, t) + \Theta(t+1, t+B)]. \end{aligned}$$

□

Lemma 9 (Extension of Lemma 4). *Under A2–A5 and A7 and let the step sizes satisfy*

$$\sup_{t \geq 0} \gamma_{t+1} \leq \rho / \sqrt{2Bc_3},$$

then it holds for any $t \geq 0$ that

$$\begin{aligned} \Theta(t+1, t+B) & \leq \frac{1 - \bar{\rho}/2}{1 - Bc_3\gamma_{t-B+1}^2/\bar{\rho}} \Theta(t-B+1, t) \\ & \quad + \frac{\gamma_{t-B+1}^2}{\rho - Bc_3\gamma_{t-B+1}^2} \left\{ B^2 d_1 + d_2 B [\tilde{\theta}(t-B+1, t) + \tilde{\theta}(t+1, t+B)] \right\}, \end{aligned} \quad (60)$$

where $d_1 := 9n(\sigma^2 + \zeta^2)$, $d_2 := 12n(\sigma^2 + \zeta^2)$.

Proof. Recall the notations (45) and observe that

$$\Theta_o^{t+1} = \Theta^{t+1} - \bar{\Theta}^{t+1} = \underbrace{\left(\mathbf{W}^{t+1} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right)}_{=\mathbf{A}^{t+1}} \Theta_o^t - \gamma_{t+1} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \tilde{\nabla} F^t.$$

Therefore, we can obtain the following consensus error recursion

$$\Theta_o^{t+1} = \mathbf{A}^{t+1} \Theta_o^t - \gamma_{t+1} (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^t.$$

Then, we aim to link Θ_o^{t+1} to Θ_o^{t-B+1} .

$$\begin{aligned} \Theta_o^{t+1} & = \mathbf{A}^{t+1} \Theta_o^t - \gamma_t (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^{t-1} \\ & = \mathbf{A}^{t+1} \mathbf{A}^t \Theta_o^{t-1} - \gamma_t \mathbf{A}^{t+1} (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^{t-1} - \gamma_{t+1} (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^t \\ & \quad \vdots \\ & = \mathbf{A}^{t+1} \mathbf{A}^t \mathbf{A}^{t-1} \dots \mathbf{A}^{t-B+1} \Theta_o^{t-B+1} - \sum_{s=t-B+1}^t \gamma_{s+1} \mathbf{A}^{s+2} (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s. \end{aligned}$$

Taking Frobenius norm on both sides and applying the Young's inequality give

$$\begin{aligned} \|\Theta_o^{t+1}\|_F^2 & \leq (1 + \alpha) \|\mathbf{A}^{t+1} \mathbf{A}^t \mathbf{A}^{t-1} \dots \mathbf{A}^{t-B+1}\|^2 \|\Theta_o^{t-B+1}\|_F^2 \\ & \quad + (1 + \alpha^{-1}) \sum_{s=t-B+1}^t \gamma_{s+1}^2 \|\mathbf{A}^{s+2}\|^2 \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2, \end{aligned}$$

which holds for any $\alpha > 0$. Using A7 and setting $\alpha = \frac{\rho}{1-\rho}$, we have

$$\|\Theta_o^{t+1}\|_F^2 \leq (1 - \bar{\rho}) \|\Theta_o^{t-B+1}\|_F^2 + \sum_{s=t-B+1}^t \gamma_{s+1}^2 \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2.$$

Similarly, we get

$$\begin{aligned} \|\Theta_o^{t+2}\|_F^2 &\leq (1-\bar{\rho})\|\Theta_o^{t-B+2}\|_F^2 + \sum_{s=t-B+2}^{t+1} \gamma_{s+1}^2 \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2 \\ &\vdots \\ \|\Theta_o^{t+B}\|_F^2 &\leq (1-\bar{\rho})\|\Theta_o^t\|_F^2 + \sum_{s=t}^{t+B-1} \gamma_{s+1}^2 \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2. \end{aligned}$$

Adding these B consensus errors together leads to

$$\begin{aligned} \Theta(t+1, t+B) &\leq (1-\bar{\rho})\Theta(t-B+1, t) \\ &+ \frac{\gamma_{t-B+1}^2}{\rho} \left\{ \sum_{s=t-B+1}^t \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2 + \cdots + \sum_{s=t}^{t+B} \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2 \right\}. \end{aligned} \quad (61)$$

Using the inequality (49) in the proof of Lemma 4, we get

$$\mathbb{E}_s \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2 \leq d_1 + d_2 \|\tilde{\theta}^s\|^2 + c_3 \|\Theta_o^s\|_F^2,$$

where $d_1 := 9n(\sigma^2 + \zeta^2)$, $d_2 := 12n(\sigma^2 + \zeta^2)$ and $c_3 = 12\sigma^2 + 18L^2(1 + \epsilon_{\max})^2$. Then, we have

$$\begin{aligned} \sum_{s=t-B+1}^t \mathbb{E} \left\| (\mathbf{I} - (1/n)\mathbf{1}\mathbf{1}^\top) \tilde{\nabla} F^s \right\|_F^2 &\leq \sum_{s=t-B+1}^t \mathbb{E} \left[d_1 + d_2 \|\tilde{\theta}^s\|^2 + c_3 \|\Theta_o^s\|_F^2 \right] \\ &= Bd_1 + d_2 \sum_{s=t-B+1}^t \mathbb{E} \|\tilde{\theta}^s\|^2 + c_3 \sum_{s=t-B+1}^t \mathbb{E} \|\Theta_o^s\|_F^2. \end{aligned}$$

Substituting back to (61) give us

$$\begin{aligned} \Theta(t+1, t+B) &\leq (1-\bar{\rho})\Theta(t-B+1, t) + \frac{\gamma_{t-B+1}^2}{\rho} \left\{ B^2 d_1 + d_2 [\tilde{\theta}(t-B+1, t) + \cdots + \tilde{\theta}(t, t+B)] \right. \\ &\quad \left. + c_3 [\Theta(t-B+1, t) + \cdots + \Theta(t, t+B)] \right\}. \end{aligned}$$

The above can be simplified to

$$\begin{aligned} \Theta(t+1, t+B) &\leq (1-\bar{\rho})\Theta(t-B+1, t) + \frac{\gamma_{t-B+1}^2}{\rho} \left\{ B^2 d_1 + d_2 B [\tilde{\theta}(t+1, t+B) + \tilde{\theta}(t, t+B)] \right. \\ &\quad \left. + c_3 B [\Theta(t-B+1, t) + \Theta(t+1, t+B)] \right\}. \end{aligned}$$

Setting $\sup_{k \geq 1} \gamma_k \leq \frac{\bar{\rho}}{\sqrt{2c_3 B}}$ and rearranging terms give us

$$\begin{aligned} \Theta(t+1, t+B) &\leq \frac{1-\bar{\rho}/2}{1-Bc_3\gamma_{t-B+1}^2/\bar{\rho}} \Theta(t-B+1, t) \\ &\quad + \frac{\gamma_{t-B+1}^2}{\rho - Bc_3\gamma_{t-B+1}^2} \left\{ B^2 d_1 + d_2 B [\tilde{\theta}(t-B+1, t) + \tilde{\theta}(t+1, t+B)] \right\}, \end{aligned}$$

which gives us desired upper bound for $\Theta(t+1, t+B)$. \square

Convergence of $\tilde{\theta}^t$ to θ^{PS} with Time varying graph. We conclude our proof sketch through analyzing the following Lyapunov function. For any $t \geq 0$, we define:

$$\mathcal{L}_{t+1}^{t+B} := \tilde{\theta}(t+1, t+B) + \Theta(t+1, t+B) \geq 0.$$

Combing Lemma 8 and 9 leads to

$$\begin{aligned}
& \left(1 - Bc_1 \frac{\gamma_{t+1}}{n} - Bc_2 \frac{\gamma_{t+1}^2}{n}\right) \Theta(t+1, t+B) + \left(1 - \frac{d_2 B \gamma_{t-B+1}^2}{\rho - Bc_3 \gamma_{t-B+1}^2}\right) \tilde{\theta}(t+1, t+B) \\
& \leq \left(\frac{1 - \bar{\rho}/2}{1 - Bc_3 \gamma_{t-B+1}^2 / \bar{\rho}} + Bc_1 \frac{\gamma_{t+1}}{n} + Bc_2 \frac{\gamma_{t+1}^2}{n}\right) \Theta(t-B+1, t) \\
& \quad + \left((1 - \tilde{\mu} \gamma_{t+B})^B + \frac{d_2 B \gamma_{t-B+1}^2}{\rho - Bc_3 \gamma_{t-B+1}^2}\right) \tilde{\theta}(t-B+1, t) + \frac{B^2 d_1 \gamma_{t-B+1}^2}{\rho - Bc_3 \gamma_{t-B+1}^2} + \frac{2B\sigma^2}{n} \gamma_{t+1}^2.
\end{aligned} \tag{62}$$

We focus on the l.h.s. of above inequality. If the step size satisfies

$$\sup_{k \geq 1} \gamma_k \leq \min \left\{ \frac{c_1}{c_2}, \sqrt{\frac{\bar{\rho}}{2Bc_3}}, \frac{\bar{\rho}c_1}{n} \right\}$$

then, the l.h.s. of (62) can be lower bounded by

$$\text{l.h.s. of (62)} \geq \left(1 - 2Bc_1 \frac{\gamma_{t+1}}{n}\right) \left[\Theta(t+1, t+B) + \tilde{\theta}(t+1, t+B)\right].$$

Next, we consider the r.h.s. of (62). Suppose that $\sup_{k \geq 1} \gamma_k \leq \sqrt{\frac{\rho}{(4-\bar{\rho})Bc_3}}$, it holds

$$\frac{1 - \bar{\rho}/2}{1 - Bc_3 \gamma_{t-B+1}^2 / \bar{\rho}} \leq 1 - \bar{\rho}/4, \quad \frac{B^2 d_1 \gamma_{t-B+1}^2}{\rho - Bc_3 \gamma_{t-B+1}^2} \leq \frac{2B^2 d_1}{\rho} \gamma_{t-B+1}^2.$$

If step size also satisfies

$$\sup_{k \geq 1} \gamma_k \leq \min \left\{ \frac{1}{\sqrt{b}}, \frac{\tilde{\mu} \bar{\rho}}{2^{2B+1} d_2 B} \right\},$$

where b such that $\gamma_k^2 / \gamma_{k+1}^2 \leq 1 + b \gamma_{k+1}^2$, then it holds:

$$\begin{aligned}
& \text{r.h.s. of (62)} \\
& \leq \left(1 - \frac{\bar{\rho}}{4} + 2Bc_1 \frac{\gamma_{t+1}}{n}\right) \Theta(t-B+1, t) + \left(1 - \frac{\tilde{\mu} \gamma_{t+B}}{2}\right) \tilde{\theta}(t-B+1, t) \\
& \quad + \frac{2B^2 d_1}{\rho} \gamma_{t-B+1}^2 + \frac{2B\sigma^2}{n} \gamma_{t+1}^2.
\end{aligned}$$

Combining the above inequalities lead to:

$$\begin{aligned}
& \left(1 - 2Bc_1 \frac{\gamma_{t+1}}{n}\right) \left[\Theta(t+1, t+B) + \tilde{\theta}(t+1, t+B)\right] \leq \left(1 - \frac{\bar{\rho}}{4} + 2Bc_1 \frac{\gamma_{t+1}}{n}\right) \Theta(t-B+1, t) \\
& \quad + \left(1 - \frac{\tilde{\mu} \gamma_{t+B}}{2}\right) \tilde{\theta}(t-B+1, t) + \frac{2B^2 d_1}{\rho} \gamma_{t-B+1}^2 + \frac{2B\sigma^2}{n} \gamma_{t+1}^2.
\end{aligned}$$

If the step size satisfying

$$\sup_{k \geq 1} \gamma_k \leq \frac{\bar{\rho}}{8Bc_1/n + 2\tilde{\mu}},$$

then the main recursion can be simplified as

$$\begin{aligned}
& \left(1 - 2Bc_1 \frac{\gamma_{t+1}}{n}\right) \left[\Theta(t+1, t+B) + \tilde{\theta}(t+1, t+B)\right] \\
& \leq \left(1 - \frac{\tilde{\mu} \gamma_{t+B}}{2}\right) \left[\Theta(t-B+1, t) + \tilde{\theta}(t-B+1, t)\right] + \frac{2B^2 d_1}{\rho} \gamma_{t-B+1}^2 + \frac{2B\sigma^2}{n} \gamma_{t+1}^2.
\end{aligned}$$

Dividing $(1 - 2Bc_1 \frac{\gamma_{t+1}}{n})$ for the both sides, we obtain that

$$\mathcal{L}_{t+1}^{t+B} \leq \frac{(1 - \tilde{\mu} \gamma_{t+B}/2)}{1 - 2Bc_1 \gamma_{t+1}/n} \mathcal{L}_{t-B+1}^t + \left(\frac{2B^2 d_1}{\rho} + \frac{2B^2 \sigma^2}{n}\right) \frac{\gamma_{t-B+1}^2}{(1 - 2Bc_1 \gamma_{t+1}/n)}.$$

Observe that with sufficiently small step size, the above recursion can be simplified to give a similar form as (25). Solving the recursion then lead to $\mathcal{L}_{t+1}^{t+B} = \mathcal{O}(\gamma_{t-B+1})$ and the convergence of $\tilde{\theta}(t+1, t+B) \rightarrow 0$.

Lastly, we remark that the above analysis only gives a crude bound to the convergence of DSGD-GD in the time varying graph setting. It is possible to give tighter bounds through further optimizing the constants in the above analysis.

H Extension to Local Distributions Influenced by All Agents

This section outlines how to extend our analysis to the scenario when the local distributions $\mathcal{D}_i(\cdot)$ are simultaneously influenced by other agents in the network similar to the competitive Multi-PfD considered by [Narang et al., 2022, Piliouras and Yu, 2022].

We define the concatenated decision vector $\boldsymbol{\vartheta} := (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n) \in \mathbb{R}^{nd}$ and state the modified consensus Multi-PfD problem (1) as follows

$$\min_{\boldsymbol{\theta}_i \in \mathbb{R}^d, i=1, \dots, n} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\vartheta})} [\ell(\boldsymbol{\theta}_i; Z_i)] \quad \text{s.t. } \boldsymbol{\theta}_i = \boldsymbol{\theta}_j, \forall (i, j) \in E. \quad (63)$$

With a slight abuse of notation, we also define $f_i(\boldsymbol{\theta}; \boldsymbol{\vartheta}) := \mathbb{E}_{Z_i \sim \mathcal{D}_i(\boldsymbol{\vartheta})} [\ell(\boldsymbol{\theta}_i; Z_i)]$.

Following [Narang et al., 2022], we consider the following modification to A3:

A8. For any $i = 1, \dots, n$, there exists a constant $\epsilon_i > 0$ such that

$$\mathcal{W}_1(\mathcal{D}_i(\boldsymbol{\vartheta}), \mathcal{D}_i(\boldsymbol{\vartheta}')) \leq \epsilon_i \|\boldsymbol{\vartheta} - \boldsymbol{\vartheta}'\|, \forall \boldsymbol{\vartheta}', \boldsymbol{\vartheta} \in \mathbb{R}^{nd}, \quad (64)$$

where $\mathcal{W}_1(\mathcal{D}, \mathcal{D}')$ denotes the Wasserstein-1 distance between the distributions $\mathcal{D}, \mathcal{D}'$.

Specifically, we notice that if $\boldsymbol{\vartheta}$ satisfies the consensus constraint, i.e., $\boldsymbol{\vartheta} = \mathbf{1}_n \otimes \boldsymbol{\theta} = (\boldsymbol{\theta}, \dots, \boldsymbol{\theta})$, then A8 is equivalent to A3 with the latter's sensitivity parameter given by $\epsilon'_i = \sqrt{n}\epsilon_i$ since $\|\mathbf{1}_n \otimes \boldsymbol{\theta} - \mathbf{1}_n \otimes \boldsymbol{\theta}'\| = \sqrt{n}\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|$. This observation immediately leads to the following corollary of Proposition 1:

Corollary 1. Under A1, A2, A8. Define the map $\mathcal{M} : \mathbb{R}^d \rightarrow \mathbb{R}^d$

$$\mathcal{M}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}' \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{\theta}'; \mathbf{1}_n \otimes \boldsymbol{\theta}) \quad (65)$$

If $\sqrt{n}\epsilon_{\text{avg}} < \mu/L$, then the map $\mathcal{M}(\boldsymbol{\theta})$ is a contraction with the unique fixed point $\boldsymbol{\theta}^{PS} = \mathcal{M}(\boldsymbol{\theta}^{PS})$. If $\sqrt{n}\epsilon_{\text{avg}} \geq \mu/L$, then there exists an instance of (11) where $\lim_{T \rightarrow \infty} \|\mathcal{M}^T(\boldsymbol{\theta})\| = \infty$.

The proof is attained by simply observing that if $\boldsymbol{\theta}_i = \boldsymbol{\theta}_j$ (as constrained by (65) (and (63))), then A8 is equivalent to A3 with $\epsilon'_i = \sqrt{n}\epsilon_i$.

Comparison to [Narang et al., 2022]. Notice that in [Narang et al., 2022], the existence of a performative stable equilibrium requires $\sqrt{\sum_{i=1}^n \epsilon_i^2} < \mu/L$. Meanwhile, Corollary 1 requires $(1/\sqrt{n}) \sum_{i=1}^n \epsilon_i < \mu/L$. Due to norm equivalence, we have

$$(1/\sqrt{n}) \sum_{i=1}^n \epsilon_i \leq \sum_{i=1}^n \epsilon_i^2.$$

Thus, the consensus constrained performative stable solution in cooperative Multi-PfD will be attainable under a more relaxed condition than the competitive Multi-PfD.

DSGD-GD Algorithm for (63). The extension of Theorem 1 to (63) via the DSGD-GD algorithm is more involved and thus the details are skipped in this brief discussion. However, it remains straightforward to extend the analysis through a careful modification of Lemma 3 with A8. In particular, one only needs to pay attention to the use of A8 in (37) and (40) for the proof.

Remarks. We emphasize that as explained in the main paper, the original scenario considered by (1) and A3 is relevant to the decentralized learning scenario of the current paper. It captures the effects of ‘geographical’ barriers where the population of users are not simultaneously influenced by all agents. Nevertheless, a future direction is to study the Multi-PfD problem (cooperative or competitive) where users can be influenced by the decisions from a *few* neighboring agents.