## Supplemental Material

We give a table to summarize the content of the supplemental material.

| Section | Content |
|---|---|
| Appendix A | some useful lemmas as technical tools |
| Appendix B | proof of Theorem 1 for $K = 1$ |
| Appendix C | proof of Theorem 2 for $K < \infty$ |
| Appendix D | proof of Theorem 3 for $K = \infty$ |
| Appendix E | a table for some important notations |

**Table 1: Outline of the supplemental material.**

## A   USEFUL LEMMAS

In this section, we provide some useful lemmas. Specifically, Lemma 1 is used to support the claim of the convergence speed in Insight 4. Lemmas 2 to 4 are some results about the Gaussian random matrices that can be found in the literature. We want to highlight Lemma 5 as part of our technical novelty, which gives the exact values of terms related to the projection formed by each agent's training inputs. Lemma 6 is used to justify the definition of model error.

LEMMA 1. *Recalling the definition of $C$ in Eq. (32), we have*

$$\lim_{t=p \ln p, \ p \to \infty} C^t = 0.$$

PROOF. We have $C^t \geq 0$ and

$$C^t \leq \left(1 - \frac{n}{p}\right)^t \quad \text{(since } C \leq \left(1 - \frac{n}{p}\right) \text{ because } \left(1 - \frac{n}{p}\right)^2 \leq \left(1 - \frac{n}{p}\right))$$

$$= \left(1 + \frac{1}{\frac{p}{n} - 1}\right)^{-t} \quad \text{(since } 1 - \frac{n}{p} = \frac{1}{1 + \frac{1}{\frac{p}{n} - 1}})$$

$$= \left(1 + \frac{1}{\frac{p}{n} - 1}\right)^{-p \ln p} \quad \text{(since } t = p \ln p)$$

$$= \left(1 + \frac{1}{\frac{p}{n} - 1}\right)^{-\frac{p}{n} \cdot n \cdot \ln p}$$

$$\leq \left(1 + \frac{1}{\frac{p}{n} - 1}\right)^{-\left(\frac{p}{n} - 1\right) \cdot n \cdot \ln p}.$$

Notice that

$$\lim_{p \to \infty} \left(1 + \frac{1}{\frac{p}{n} - 1}\right)^{-\left(\frac{p}{n} - 1\right) \cdot n \cdot \ln p} = \lim_{p \to \infty} e^{-n \ln p} = 0,$$

where we use the fact that $\lim_{x \to \infty}(1 + x^{-1})^x = e$. The result of this lemma thus follows by the squeeze theorem. □

The result of the following lemma can be found in the literature (e.g., [19, 60]).

LEMMA 2. *Consider a random matrix $\mathbf{K} \in \mathbb{R}^{p \times n}$ where $p$ and $n$ are two positive integers and $p > n + 1$. Each element of $\mathbf{K}$ is i.i.d. according to standard Gaussian distribution. For any fixed vector $\boldsymbol{a} \in \mathbb{R}^p$, we must have*

$$\mathbb{E} \left\| \left(\mathbf{I}_p - \mathbf{K} \left(\mathbf{K}^\top \mathbf{K}\right)^{-1} \mathbf{K}^\top \right) \boldsymbol{a} \right\|^2 = \left(1 - \frac{n}{p}\right) \|\boldsymbol{a}\|^2,$$

$$\mathbb{E} \left\| \mathbf{K} \left(\mathbf{K}^\top \mathbf{K}\right)^{-1} \mathbf{K}^\top \boldsymbol{a} \right\|^2 = \frac{n}{p} \|\boldsymbol{a}\|^2.$$

The following lemma can be found in Lemma 8 of [61].

LEMMA 3. *Consider a random matrix* $\mathbf{K} \in \mathbb{R}^{a \times b}$ *where* $a > b + 1$. *Each element of* $\mathbf{K}$ *is i.i.d. following standard Gaussian distribution* $\mathcal{N}(0,1)$. *Consider three Gaussian random vectors* $\boldsymbol{\alpha}, \boldsymbol{\gamma} \in \mathbb{R}^a$ *and* $\boldsymbol{\beta} \in \mathbb{R}^b$ *such that* $\boldsymbol{\alpha} \sim \mathcal{N}(0, \sigma_\alpha^2 \mathbf{I}_a)$, $\boldsymbol{\gamma} \sim \mathcal{N}(0, \mathrm{diag}(d_1^2, d_2^2, \cdots, d_a^2))$, *and* $\boldsymbol{\beta} \sim \mathcal{N}(0, \sigma_\beta^2 \mathbf{I}_b)$. *Here* $\mathbf{K}, \boldsymbol{\alpha}, \boldsymbol{\gamma}$, *and* $\boldsymbol{\beta}$ *are independent of each other. We then must have*

$$\mathbb{E}\left[(\mathbf{K}^\top \mathbf{K})^{-1}\right] = \frac{\mathbf{I}_b}{a - b - 1}, \tag{44}$$

$$\mathbb{E}\left\|\mathbf{K}(\mathbf{K}^\top \mathbf{K})^{-1}\boldsymbol{\beta}\right\|^2 = \frac{b\sigma_\beta^2}{a - b - 1}, \tag{45}$$

$$\mathbb{E}\left\|(\mathbf{K}^\top \mathbf{K})^{-1}\mathbf{K}^\top \boldsymbol{\alpha}\right\|^2 = \frac{b\sigma_\alpha^2}{a - b - 1}, \tag{46}$$

$$\mathbb{E}\left\|(\mathbf{K}^\top \mathbf{K})^{-1}\mathbf{K}^\top \boldsymbol{\gamma}\right\|^2 = \frac{b \sum_{i=1}^a d_i^2}{a(a - b - 1)}. \tag{47}$$

The following lemma can be found in [62] and Lemma 13 of [60].

LEMMA 4. *Consider a random matrix* $\mathbf{K} \in \mathbb{R}^{a \times b}$ *whose each element follows* i.i.d. *standard Gaussian distribution (i.e.,* i.i.d. $\mathcal{N}(0,1)$). *We mush have*

$$\mathbb{E}[\mathbf{K}^\top \mathbf{K}] = a\mathbf{I}_b,$$
$$\mathbb{E}[\mathbf{K}\mathbf{K}^\top] = b\mathbf{I}_a,$$
$$\mathbb{E}[\mathbf{K}\mathbf{K}^\top \mathbf{K}\mathbf{K}^\top] = b(b + a + 1)\mathbf{I}_a.$$

LEMMA 5. *For any* $i \in [m]$ *and* $t$, *we must have*

$$\mathop{\mathbb{E}}_{\mathbf{P}_{(i),t}}\left[\mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty}\right] = \frac{n_{(i),t}}{p}\Delta_{t-1}^{K=\infty}. \tag{48}$$

*Consequently, when* $i \neq j$, *we have*

$$\mathop{\mathbb{E}}_{\mathbf{P}_{(i),t}, \mathbf{P}_{(j),t}}\left[\Delta_{t-1}^{K=\infty \top}\mathbf{P}_{(i),t}\mathbf{P}_{(j),t}\Delta_{t-1}^{K=\infty}\right] = \frac{n_{(j),t}n_{(i),t}}{p^2}\left\|\Delta_{t-1}^{K=\infty}\right\|^2.$$

PROOF. Let $C := \left\|\Delta_{t-1}^{K=\infty}\right\|$. Since we are calculating expected projection of $\Delta_{t-1}^{K=\infty}$ onto the column space of $\mathbf{X}_{(i),t}$, by the symmetry of $\mathbf{X}_{(i),t}$, without loss of generality we let

$$\Delta_{t-1}^{K=\infty} = C \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}. \tag{49}$$

Define

$$\tilde{\mathbf{X}}_{(i),t} := \begin{bmatrix} -1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}\mathbf{X}_{(i),t}. \tag{50}$$

Since each element of $\mathbf{X}_{(i),t}$ follows *i.i.d.* standard Gaussian distribution, we know that $\tilde{\mathbf{X}}_{(i),t}$ and $\mathbf{X}_{(i),t}$ has identical distributioin. Thus, we have

$$\int \mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty}d\mu(\mathbf{X}_{(i),t}) = \int \tilde{\mathbf{X}}_{(i),t}(\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t})\tilde{\mathbf{X}}_{(i),t}\Delta_{t-1}^{K=\infty}d\mu(\mathbf{X}_{(i),t}), \tag{51}$$

where $\mu(\mathbf{X}_{(i),t})$ denotes the joint probability distribution of $\mathbf{X}_{(i),t}$.

By Eq. (50), we have

$$\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t} = \mathbf{X}_{(i),t}^\top \begin{bmatrix} -1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}\begin{bmatrix} -1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}\mathbf{X}_{(i),t} = \mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t},$$

$$\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} = [\mathbf{X}_{(i),t}]_{1,:}, \quad \tilde{\mathbf{X}}_{(i),t}^\top \Delta_{t-1}^{K=\infty} = -[\mathbf{X}_{(i),t}]_{1,:} \text{ (here } [\cdot]_{1,:} \text{ denotes the first row of a matrix)}.$$

Thus, we have

$$\tilde{\mathbf{X}}_{(i),t}(\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t})^{-1}\tilde{\mathbf{X}}_{(i),t}^\top \Delta_{t-1}^{K=\infty} = -\tilde{\mathbf{X}}_{(i),t}(\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty}. \tag{52}$$

Therefore, we have

$$\mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} + \tilde{\mathbf{X}}_{(i),t}(\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t})^{-1}\tilde{\mathbf{X}}_{(i),t}^\top \Delta_{t-1}^{K=\infty}$$

$$=(\mathbf{X}_{(i),t} - \tilde{\mathbf{X}}_{(i),t})(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} \quad \text{(by Eq. (52))}$$

$$=\begin{bmatrix} 2 & & & \\ & 0 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix} \mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} \quad \text{(by Eq. (50))}$$

$$=\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 2 & 0 & \cdots & 0 \end{bmatrix} \mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty}$$

$$=2\frac{\Delta_{t-1}^{K=\infty}}{C^2}\Delta_{t-1}^{K=\infty\top}\mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} \quad \text{(by Eq. (49))}$$

$$=2\frac{\Delta_{t-1}^{K=\infty}}{C^2}\Delta_{t-1}^{K=\infty\top}\mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty} \quad \text{(by Eq. (83))}$$

$$=2\frac{\Delta_{t-1}^{K=\infty}}{C^2}\left\| \mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty}\right\|^2 \quad (\text{since } \mathbf{P}_{(i),t}^\top \mathbf{P}_{(i),t} = \mathbf{P}_{(i),t} \text{ as } \mathbf{P}_{(i),t} \text{ is an orthogonal projection}). \tag{53}$$

Thus, we have

$$\underset{\mathbf{X}_{(i),t}}{\mathbb{E}}[\mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty}]$$

$$= \int \mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} d\mu(\mathbf{X}_{(i),t})$$

$$=\frac{1}{2}\int \left(\mathbf{X}_{(i),t}(\mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t})^{-1}\mathbf{X}_{(i),t}^\top \Delta_{t-1}^{K=\infty} + \tilde{\mathbf{X}}_{(i),t}(\tilde{\mathbf{X}}_{(i),t}^\top \tilde{\mathbf{X}}_{(i),t})\tilde{\mathbf{X}}_{(i),t}^\top \Delta_{t-1}^{K=\infty}\right) d\mu(\mathbf{X}_{(i),t}) \quad \text{(by Eq. (51))}$$

$$= \int \frac{\Delta_{t-1}^{K=\infty}}{C^2}\left\| \mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty}\right\|^2 d\mu(\mathbf{X}_{(i),t})$$

$$=\frac{\Delta_{t-1}^{K=\infty}}{C^2}\underset{\mathbf{X}_{(i),t}}{\mathbb{E}}\left\| \mathbf{P}_{(i),t}\Delta_{t-1}^{K=\infty}\right\|^2$$

$$=\frac{n_{(i),t}}{p}\Delta_{t-1}^{K=\infty} \quad \text{(by Lemma 2)}.$$

The result of this lemma thus follows.                                                                                                                                     □

LEMMA 6. *Let the noise in every test sample have zero mean and variance $\sigma^2$. For any learning result $\hat{\mathbf{w}}$, the mean square test error must equal to $\|\hat{\mathbf{w}} - \mathbf{w}^*\|^2 + \sigma^2$. Therefore, the mean squared test error for noise-free test samples equals to the model error $L^{model}(\hat{\mathbf{w}}) = \|\hat{\mathbf{w}} - \mathbf{w}^*\|^2$.*

PROOF. Considering $(\mathbf{x}, y)$ as a randomly generated test sample by the ground truth $y = \mathbf{x}^\top \mathbf{w}^* + \epsilon$, the mean squared error is equal to

$$\underset{\mathbf{x},y}{\mathbb{E}}\left\| \mathbf{x}^\top \hat{\mathbf{w}} - y \right\| = \underset{\mathbf{x},\epsilon}{\mathbb{E}}\left\| \mathbf{x}^\top \hat{\mathbf{w}} - (\mathbf{x}^\top \mathbf{w}^* + \epsilon) \right\|^2$$

$$= \underset{\mathbf{x},\epsilon}{\mathbb{E}}\left\| \mathbf{x}^\top (\hat{\mathbf{w}} - \mathbf{w}^*) + \epsilon \right\|^2$$

$$= \underset{\mathbf{x}}{\mathbb{E}}\left\| \mathbf{x}^\top (\hat{\mathbf{w}} - \mathbf{w}^*) \right\|^2 + \underset{\epsilon}{\mathbb{E}}\left\| \epsilon \right\|^2$$

(since the noise $\epsilon$ has zero mean and is independent of other random variables)

$$= \left\| \hat{\mathbf{w}} - \mathbf{w}^* \right\|^2 + \sigma^2$$

(notice that $\mathbf{x}$ follows standard Gaussian distribution and is independent of $\hat{\mathbf{w}}$).

□

# B PROOF OF THEOREM 1

Calculating the gradient of the training loss defined at Eq. (4), we have

$$
\begin{aligned}
\frac{\partial L(\hat{\boldsymbol{w}})}{\partial \hat{\boldsymbol{w}}} &= \frac{\partial (\boldsymbol{y} - \mathbf{X}^\top \hat{\boldsymbol{w}})}{\partial \hat{\boldsymbol{w}}} \cdot \frac{\partial \frac{1}{2n} \left\| \boldsymbol{y} - \mathbf{X}^\top \hat{\boldsymbol{w}} \right\|^2}{\partial (\boldsymbol{y} - \mathbf{X}^\top \hat{\boldsymbol{w}})} \quad \text{(by the chain rule)} \\
&= -\mathbf{X} \cdot \frac{1}{n} (\boldsymbol{y} - \mathbf{X}^\top \hat{\boldsymbol{w}}) \\
&= \frac{1}{n} (\mathbf{X}\mathbf{X}^\top \hat{\boldsymbol{w}} - \mathbf{X}\boldsymbol{y}).
\end{aligned}
$$

When $K = 1$, with step size $\alpha_{(i),t} > 0$, we thus have

$$
\begin{aligned}
\hat{\boldsymbol{w}}_{(i),t}^{K=1} &= \left( \mathbf{I}_p - \frac{\alpha_{(i),t}}{n_{(i),t}} \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \right) \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=1} + \frac{\alpha_{(i),t}}{n_{(i),t}} \mathbf{X}_{(i),t} \boldsymbol{y}_{(i),t} \\
&= \left( \mathbf{I}_p - \frac{\alpha_{(i),t}}{n_{(i),t}} \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \right) \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=1} + \frac{\alpha_{(i),t}}{n_{(i),t}} \mathbf{X}_{(i),t} \left( \mathbf{X}_{(i),t}^\top \boldsymbol{w}_{(i),t} + \boldsymbol{\epsilon}_{(i),t} \right) \quad \text{(by Eq. (2)).}
\end{aligned}
$$

Thus, we have

$$
\begin{aligned}
\hat{\boldsymbol{w}}_{\text{avg},t}^{K=1} &= \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} n_{(i),t} \hat{\boldsymbol{w}}_{(i),t}^{K=1} \\
&= \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=1} + \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} \alpha_{(i),t} \left( -\mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=1} + \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \boldsymbol{w}_{(i),t} + \mathbf{X}_{(i),t} \boldsymbol{\epsilon}_{(i),t} \right).
\end{aligned}
\tag{54}
$$

By Eqs. (3) and (8), we have

$$
\begin{aligned}
&\Delta_t^{K=1} \\
&= \Delta_{t-1}^{K=1} + \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} \alpha_{(i),t} \left( \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top (\boldsymbol{\gamma}_{(i),t} - \Delta_{t-1}^{K=1}) - \mathbf{X}_{(i),t} \boldsymbol{\epsilon}_{(i),t} \right) \\
&= \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} \left( \underbrace{\left( n_{(i),t} \mathbf{I}_p - \alpha_{(i),t} \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \right) \Delta_{t-1}^{K=1}}_{\boldsymbol{q}_{1i}} + \underbrace{\alpha_{(i),t} \mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top \boldsymbol{\gamma}_{(i),t}}_{\boldsymbol{q}_{2i}} - \underbrace{\alpha_{(i),t} \mathbf{X}_{(i),t} \boldsymbol{\epsilon}_{(i),t}}_{\boldsymbol{q}_{3i}} \right) \\
&\quad \text{(since } \Delta_{t-1}^{K=1} = \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} n_{(i),t} \Delta_{t-1}^{K=1} \text{).}
\end{aligned}
\tag{55}
$$

Considering the three types of terms $\boldsymbol{q}_{1i}, \boldsymbol{q}_{2i}, \boldsymbol{q}_{3i}$ defined in Eq. (55), by Assumption 1, we have

$$
\begin{aligned}
\mathbb{E}_t \boldsymbol{q}_{1i} &= n_{(i),t} \left( 1 - \alpha_{(i),t} \right) \Delta_{t-1}^{K=1}, \\
\mathbb{E}_t \boldsymbol{q}_{2i} &= \alpha_{(i),t} n_{(i),t} \boldsymbol{\gamma}_{(i),t}, \\
\mathbb{E}_t \boldsymbol{q}_{3i} &= \mathbf{0}.
\end{aligned}
\tag{56}
$$

Notice that we use $\mathbb{E}$ to denote the expectation on all randomness and use $\mathbb{E}_t$ to denote the expectation on the randomness at the $t$-th round, i.e., on the randomness of $\mathbf{X}_{(i),t}$ and $\boldsymbol{\epsilon}_{(i),t}$ for all $i \in [m]$. By Eqs. (55) and (56), we thus have

$$
\mathbb{E}_t \Delta_t^{K=1} = \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} \left( n_{(i),t} \left( 1 - \alpha_{(i),t} \right) \Delta_{t-1}^{K=1} + \alpha_{(i),t} n_{(i),t} \boldsymbol{\gamma}_{(i),t} \right).
\tag{57}
$$

Applying Eq. (57) recursively and recalling Eq. (15), we thus have

$$
\mathbb{E}[\Delta_t^{K=1}] = \boldsymbol{g}_t^{K=1}.
\tag{58}
$$

By Assumption 1, we know that $\boldsymbol{\epsilon}_{(i),t}$ is independent of $\mathbf{X}_{(j),t}$ for all $i, j \in [m]$ and $\mathbb{E}\boldsymbol{\epsilon}_{(i),t} = \mathbf{0}$. Thus, we have

$$
\mathbb{E}_t[\boldsymbol{q}_{1i}^\top \boldsymbol{q}_{3j}] = \mathbb{E}_t[\boldsymbol{q}_{2i}^\top \boldsymbol{q}_{3j}] = 0.
$$

Thus, we have

$$
\begin{aligned}
\mathop{\mathbb{E}}_t \left\| \Delta_t^{K=1} \right\|^2 =& \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \left( \sum_{i \in [m]} \left( \mathop{\mathbb{E}}_t \|q_{1i}\|^2 + \mathop{\mathbb{E}}_t \|q_{2i}\|^2 + \mathop{\mathbb{E}}_t \|q_{3i}\|^2 + 2 \mathop{\mathbb{E}}_t [q_{1i}^\top q_{2i}] \right) \right. \\
& \left. + \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} \left( \mathop{\mathbb{E}}_t [q_{1i}^\top q_{1j}] + \mathop{\mathbb{E}}_t [q_{1i}^\top q_{2j}] + \mathop{\mathbb{E}}_t [q_{1j}^\top q_{2i}] + \mathop{\mathbb{E}}_t [q_{2i}^\top q_{2j}] \right) \right) \\
=& \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \left( \sum_{i \in [m]} \left( \mathop{\mathbb{E}}_t \|q_{1i}\|^2 + \mathop{\mathbb{E}}_t \|q_{2i}\|^2 + \mathop{\mathbb{E}}_t \|q_{3i}\|^2 + 2 \mathop{\mathbb{E}}_t [q_{1i}^\top q_{2i}] \right) \right. \\
& \left. + \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} \left( \mathop{\mathbb{E}}_t [q_{1i}^\top q_{1j}] + 2 \mathop{\mathbb{E}}_t [q_{1i}^\top q_{2j}] + \mathop{\mathbb{E}}_t [q_{2i}^\top q_{2j}] \right) \right) \\
& (\text{since } \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} q_{1i}^\top q_{2j} + q_{1j}^\top q_{2i} = 2 \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} q_{1i}^\top q_{2j}).
\end{aligned}
\tag{59}
$$

By Lemma 4, for any $i \in [m]$, we have

$$
\begin{aligned}
\mathop{\mathbb{E}}_t \|q_{1i}\|^2 =& \left( n_{(i),t}^2 - 2\alpha_{(i),t} n_{(i),t}^2 + \alpha_{(i),t}^2 n_{(i),t} (n_{(i),t} + p + 1) \right) \left\| \Delta_{t-1}^{K=1} \right\|^2 \\
=& \left( \left( 1 - \alpha_{(i),t} \right)^2 n_{(i),t}^2 + \alpha_{(i),t}^2 n_{(i),t} (p + 1) \right) \left\| \Delta_{t-1}^{K=1} \right\|^2, \\
\mathop{\mathbb{E}}_t \|q_{2i}\|^2 =& \alpha_{(i),t}^2 n_{(i),t} (n_{(i),t} + p + 1) \left\| \gamma_{(i),t} \right\|^2, \\
\mathop{\mathbb{E}}_t \|q_{3i}\|^2 =& \alpha_{(i),t}^2 p n_{(i),t} \sigma_{(i),t}^2, \\
\mathop{\mathbb{E}}_t [q_{1i}^\top q_{2i}] =& \left( \alpha_{(i),t} n_{(i),t}^2 - \alpha_{(i),t}^2 n_{(i),t} (n_{(i),t} + p + 1) \right) \Delta_{t-1}^{K=1\top} \gamma_{(i),t}.
\end{aligned}
\tag{60}
$$

Similarly, by Lemma 4, for any $i, j \in [m]$ where $i \neq j$, we have

$$
\begin{aligned}
\mathbb{E}[q_{1i}^\top q_{1j}] =& n_{(i),t} n_{(j),t} \left( 1 - \alpha_{(i),t} \right) \left( 1 - \alpha_{(j),t} \right) \left\| \Delta_{t-1}^{K=1} \right\|^2, \\
\mathbb{E}[q_{1i}^\top q_{2j}] =& \left( \alpha_{(j),t} n_{(i),t} n_{(j),t} - \alpha_{(i),t} \alpha_{(j),t} n_{(i),t} n_{(j),t} \right) \Delta_{t-1}^{K=1\top} \gamma_{(j),t} \\
=& n_{(i),t} n_{(j),t} \alpha_{(j),t} \left( 1 - \alpha_{(i),t} \right) \Delta_{t-1}^{K=1\top} \gamma_{(j),t}, \\
\mathbb{E}[q_{2i}^\top q_{2j}] =& \alpha_{(i),t} \alpha_{(j),t} n_{(i),t} n_{(j),t} \gamma_{(i),t}^\top \gamma_{(j),t}.
\end{aligned}
\tag{61}
$$

Plugging Eqs. (60) and (61) into Eq. (59), we thus have

$$
\begin{aligned}
& \mathop{\mathbb{E}}_t [ \left\| \Delta_t^{K=1} \right\|^2 ] \\
=& \frac{\left\| \Delta_{t-1}^{K=1} \right\|^2}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \left( \sum_{i \in [m]} \left( (1 - \alpha_{(i),t})^2 n_{(i),t}^2 + \alpha_{(i),t}^2 n_{(i),t} (p + 1) \right) + \sum_{i \in [m]} \sum_{j \in [m] \setminus \{j\}} n_{(i),t} n_{(j),t} (1 - \alpha_{(i),t})(1 - \alpha_{(j),t}) \right) \\
& + \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \sum_{i \in [m]} \alpha_{(i),t}^2 \left( p n_{(i),t} \sigma_{(i),t}^2 + n_{(i),t} (n_{(i),t} + p + 1) \left\| \gamma_{(i),t} \right\|^2 \right) \\
& + 2 \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \sum_{i \in [m]} \left( \alpha_{(i),t} n_{(i),t}^2 - \alpha_{(i),t}^2 n_{(i),t} (n_{(i),t} + p + 1) \right) \Delta_{t-1}^{K=1\top} \gamma_{(i),t} \\
& + \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} \left( 2 n_{(i),t} n_{(j),t} \alpha_{(j),t} \left( 1 - \alpha_{(i),t} \right) \Delta_{t-1}^{K=1\top} \gamma_{(j),t} + \alpha_{(i),t} \alpha_{(j),t} n_{(i),t} n_{(j),t} \gamma_{(i),t}^\top \gamma_{(j),t} \right).
\end{aligned}
\tag{62}
$$

Notice that

$$
\left( \sum_{i \in [m]} \left( (1 - \alpha_{(i),t})^2 n_{(i),t}^2 + \alpha_{(i),t}^2 n_{(i),t}(p+1) \right) + \sum_{i \in [m]} \sum_{j \in [m] \setminus \{j\}} n_{(i),t} n_{(j),t} (1 - \alpha_{(i),t})(1 - \alpha_{(j),t}) \right)
$$

$$
= \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \left( \sum_{i \in [m]} n_{(i),t}(1 - \alpha_{(i),t})^2 \right)^2 + \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \alpha_{(i),t}^2 n_{(i),t}(p+1)
$$

$= H_t$ (recalling Eq. (16)),

and

$$
\frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \alpha_{(i),t}^2 n_{(i),t}(n_{(i),t} + p + 1) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2
$$

$$
+ \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} \alpha_{(i),t} \alpha_{(j),t} n_{(i),t} n_{(j),t} \boldsymbol{\gamma}_{(i),t}^\top \boldsymbol{\gamma}_{(j),t}
$$

$$
= \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \left\| \sum_{i \in [m]} \alpha_{(i),t} n_{(i),t} \boldsymbol{\gamma}_{(i),t} \right\|^2 + \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \alpha_{(i),t}^2 n_{(i),t}(p+1) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2,
$$

and

$$
2 \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \left( \alpha_{(i),t} n_{(i),t}^2 - \alpha_{(i),t}^2 n_{(i),t}(n_{(i),t} + p + 1) \right) \Delta_{t-1}^{K=1 \top} \boldsymbol{\gamma}_{(i),t}
$$

$$
+ \frac{1}{(\sum_{i \in [m]} n_{(i),t})^2} \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} \left( 2 n_{(i),t} n_{(j),t} \alpha_{(j),t} \left( 1 - \alpha_{(i),t} \right) \Delta_{t-1}^{K=1 \top} \boldsymbol{\gamma}_{(j),t} \right)
$$

$$
= \frac{2}{(\sum_{i \in [m]} n_{(i),t})^2} \left( \sum_{i \in [m]} n_{(i),t}(1 - \alpha_{(i),t}) \right) \cdot \left( \sum_{i \in [m]} n_{(i),t} \alpha_{(i),t} \Delta_{t-1}^{K=1 \top} \boldsymbol{\gamma}_{(i),t} \right)
$$

$$
- \frac{2 \sum_{i \in [m]} \alpha_{(i),t}^2 n_{(i),t}(p+1) \Delta_{t-1}^{K=1 \top} \boldsymbol{\gamma}_{(i),t}}{(\sum_{i \in [m]} n_{(i),t})^2}.
$$

Further, by Eq. (58) and recalling Eq. (17), we thus can rewrite Eq. (62) as

$$
\mathbb{E} \left\| \Delta_t^{K=1} \right\|^2 = H_t \, \mathbb{E} \left\| \Delta_{t-1}^{K=1} \right\|^2 + G_t. \tag{63}
$$

Applying Eq. (63) recursively, we thus have Eq. (18).

## C PROOF OF THEOREM 2

Define

$$
g_l^{K < \infty} := \mathcal{F} \left( l, \Delta_0, \text{seq}_t \left( \frac{\sum_{i \in [m]} n_{(i),t}(1 - \alpha_{(i),t})^K}{\sum_{i \in [m]} n_{(i),t}} \right), \text{seq}_t \left( \frac{\sum_{i \in [m]} n_{(i),t} \left( 1 - (1 - \alpha_{(i),t})^K \right) \boldsymbol{\gamma}_{(i),t}}{\sum_{i \in [m]} n_{(i),t}} \right) \right) \tag{64}
$$

$$
\mathcal{A}_{(i),t} := (1 - \alpha_{(i),t})^2 + \frac{\alpha_{(i),t}^2 (p+1)}{\tilde{n}_{(i),t}}, \tag{65}
$$

$$
\mathcal{B}_{(i),t,k} := \frac{\alpha_{(i),t}^2 p \sigma_{(i),t}^2}{\tilde{n}_{(i),t}}
$$

$$
+ \left( \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) + 2 \alpha_{(i),t} \left( 1 - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) \left( 1 - (1 - \alpha_{(i),t})^{k-1} \right) \right) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2
$$

$$
+ 2 \left( \alpha_{(i),t} - \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) (1 - \alpha_{(i),t})^{k-1} \boldsymbol{\gamma}_{(i),t}^\top g_{t-1}^{K < \infty}, \tag{66}
$$

$$\mathcal{J}_t := \frac{\sum_{i\in[m]} n_{(i),t}^2 \mathcal{A}_{(i),t}^K}{(\sum_{i\in[m]} n_{(i),t})^2} + \frac{\sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} (1-\alpha_{(i),t})^K (1-\alpha_{(j),t})^K}{(\sum_{i\in[m]} n_{(i),t})^2},$$ (67)

$$Q_t := \frac{\sum_{i\in[m]} n_{(i),t}^2 \sum_{k=1}^K \mathcal{B}_{(i),t,k} \mathcal{A}_{(i),t}^{K-k}}{(\sum_{i\in[m]} n_{(i),t})^2}$$

$$+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} \left( 2(1-\alpha_{(i),t})^K (1-(1-\alpha_{(j),t})^K) \boldsymbol{\gamma}_{(j),t}^\top \boldsymbol{g}_{t-1}^{K<\infty} \right.$$

$$\left. +(1-(1-\alpha_{(i),t})^K)(1-(1-\alpha_{(j),t})^K) \boldsymbol{\gamma}_{(i),t}^\top \boldsymbol{\gamma}_{(j),t} \right).$$ (68)

In the following, we use $\mathbb{E}_k$ to denote the expectation with respect to the randomness in the $k$-th batch. We have

$$\Delta_t^{K<\infty} = \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{\text{avg},t}^{K<\infty}$$

$$= \boldsymbol{w}^* - \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} \hat{\boldsymbol{w}}_{(i),t}$$

$$= \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t}) \quad (\text{since } \boldsymbol{w}^* = \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} \boldsymbol{w}^*).$$

Thus, we have

$$\left\| \Delta_t^{K<\infty} \right\|^2 = \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} n_{(i),t}^2 \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t} \right\|^2$$

$$+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t})^\top (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(j),t}).$$ (69)

By Assumption 1, we know that at round $t$, different agents' data are independent with each other. Thus, we have

$$\mathbb{E}_t (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t})^\top (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(j),t}) = \mathbb{E}_t (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t})^\top \mathbb{E}_t (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(j),t}).$$

Thus, by Eq. (69), to calculate $\mathbb{E}_t \left\| \Delta_t^{K<\infty} \right\|^2$, it remains to calculate $\mathbb{E}_t \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t} \right\|^2$ and $\mathbb{E}_t (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t})$ for all $i \in [m]$. To that end, we have

$$\hat{\boldsymbol{w}}_{(i),t,k} = \left( \mathbf{I}_p - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \mathbf{X}_{(i),t,k} \mathbf{X}_{(i),t,k}^\top \right) \hat{\boldsymbol{w}}_{(i),t,k-1} + \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \mathbf{X}_{(i),t,k} (\mathbf{X}_{(i),t,k}^\top \boldsymbol{w}_{(i),t} + \boldsymbol{\epsilon}_{(i),t,k}).$$

We thus have

$$\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k} = \left( \mathbf{I}_p - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \mathbf{X}_{(i),t,k} \mathbf{X}_{(i),t,k}^\top \right) (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k-1}) + \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \mathbf{X}_{(i),t,k} \mathbf{X}_{(i),t,k}^\top (\boldsymbol{w}^* - \boldsymbol{w}_{(i),t})$$

$$+ \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \mathbf{X}_{(i),t,k} \boldsymbol{\epsilon}_{(i),t,k}.$$ (70)

By Lemma 4 and recalling Eq. (3), we thus have

$$\mathbb{E}_k (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k}) = (1-\alpha_{(i),t})(\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k-1}) + \alpha_{(i),t} \boldsymbol{\gamma}_{(i),t}.$$ (71)

Applying Eq. (71) recursively and recalling that $\hat{\boldsymbol{w}}_{(i),t,0} = \Delta_{t-1}^{K<\infty}$, we thus have

$$\mathbb{E}_{1,2,\cdots,k} (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k}) = (1-\alpha_{(i),t})^k \Delta_{t-1}^{K<\infty} + \left( 1 - (1-\alpha_{(i),t})^k \right) \boldsymbol{\gamma}_{(i),t}.$$ (72)

By letting $k = K$ in Eq. (72) and $\hat{\boldsymbol{w}}_{(i),t,K} = \hat{\boldsymbol{w}}_{(i),t}$, we thus have

$$\mathbb{E}_t (\boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t}) = (1-\alpha_{(i),t})^K \Delta_{t-1}^{K<\infty} + \left( 1 - (1-\alpha_{(i),t})^K \right) \boldsymbol{\gamma}_{(i),t}.$$ (73)

Plugging Eq. (73) into Eq. (69), we thus have

$$
\begin{aligned}
\mathbb{E}_t \left\| \Delta_t^{K<\infty} \right\|^2 =& \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} n_{(i),t}^2 \, \mathbb{E}_t \left\| w^* - \hat{w}_{(i),t} \right\|^2 \\
&+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} \, \mathbb{E}_t (w^* - \hat{w}_{(i),t})^\top \, \mathbb{E}_t (w^* - \hat{w}_{(j),t}) \\
=& \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} n_{(i),t}^2 \, \mathbb{E}_t \left\| w^* - \hat{w}_{(i),t} \right\|^2 \\
&+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} \left( (1-\alpha_{(i),t})^K (1-\alpha_{(j),t})^K \left\| \Delta_{t-1}^{K<\infty} \right\|^2 \right. \\
&+ (1-\alpha_{(i),t})^K (1-(1-\alpha_{(j),t})^K) \gamma_{(j),t}^\top \Delta_{t-1}^{K<\infty} + (1-\alpha_{(j),t})^K (1-(1-\alpha_{(i),t})^K) \gamma_{(i),t}^\top \Delta_{t-1}^{K<\infty} \\
&+ (1-(1-\alpha_{(i),t})^K)(1-(1-\alpha_{(j),t})^K) \gamma_{(i),t}^\top \gamma_{(j),t} \Big) \\
=& \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} n_{(i),t}^2 \, \mathbb{E}_t \left\| w^* - \hat{w}_{(i),t} \right\|^2 \\
&+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} \left( (1-\alpha_{(i),t})^K (1-\alpha_{(j),t})^K \left\| \Delta_{t-1}^{K<\infty} \right\|^2 \right. \\
&+ 2(1-\alpha_{(i),t})^K (1-(1-\alpha_{(j),t})^K) \gamma_{(j),t}^\top \Delta_{t-1}^{K<\infty} \\
&+ (1-(1-\alpha_{(i),t})^K)(1-(1-\alpha_{(j),t})^K) \gamma_{(i),t}^\top \gamma_{(j),t} \Big).
\end{aligned}
$$

(74)

(75)

Notice that in Eq. (74) we use $\mathbb{E}_t (w^* - \hat{w}_{(i),t})^\top (w^* - \hat{w}_{(j),t}) = \mathbb{E}_t (w^* - \hat{w}_{(i),t})^\top \mathbb{E}_t (w^* - \hat{w}_{(j),t})$ for $i \neq j$, since $\hat{w}_{(i),t}$ and $\hat{w}_{(j),t}$ are independent with respect to the randomness during the local updates at round $t$.

By Eqs. (5) and (73), we thus have

$$
\mathbb{E}\, \Delta_t^{K<\infty} = \frac{\sum_{i\in[m]} n_{(i),t}(1-\alpha_{(i),t})^K}{\sum_{i\in[m]} n_{(i),t}} \mathbb{E}\, \Delta_{t-1}^{K<\infty} + \frac{\sum_{i\in[m]} n_{(i),t}\left(1-(1-\alpha_{(i),t})^K\right)\gamma_{(i),t}}{\sum_{i\in[m]} n_{(i),t}}.
$$

(76)

Applying Eq. (76) recursively and recalling Eq. (9), we thus have

$$
\mathbb{E}[\Delta_l^{K<\infty}] = g_l^{K<\infty},
$$

(77)

where $g_l^{K<\infty}$ is defined in Eq. (64).

By Eq. (70), we have

$$
\begin{aligned}
& \mathbb{E}_k \left\| w^* - \hat{w}_{(i),t,k} \right\|^2 \\
=& (w^* - \hat{w}_{(i),t,k-1})^\top \left( I_p - 2\frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} X_{(i),t,k} X_{(i),t,k}^\top + \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}^2} X_{(i),t,k} X_{(i),t,k}^\top X_{(i),t,k} X_{(i),t,k}^\top \right) (w^* - \hat{w}_{(i),t,k-1}) \\
&+ \gamma_{(i),t}^\top \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}^2} X_{(i),t,k} X_{(i),t,k}^\top X_{(i),t,k} X_{(i),t,k}^\top \gamma_{(i),t} + \epsilon_{(i),t,k}^\top \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}^2} X_{(i),t,k}^\top X_{(i),t,k} \epsilon_{(i),t,k} \\
&+ 2\frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} \gamma_{(i),t}^\top X_{(i),t,k} X_{(i),t,k}^\top \left( I_p - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} X_{(i),t,k} X_{(i),t,k}^\top \right) (w^* - \hat{w}_{(i),t,k-1}) \\
=& \left( 1 - 2\alpha_{(i),t} + \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) \left\| w^* - \hat{w}_{(i),t,k-1} \right\|^2 + \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \left\| \gamma_{(i),t} \right\|^2 \\
&+ \alpha_{(i),t}^2 \frac{p}{\tilde{n}_{(i),t}} \sigma_{(i),t}^2 + 2\alpha_{(i),t} \left( 1 - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) \gamma_{(i),t}^\top (w^* - \hat{w}_{(i),t,k-1}) \quad \text{(by Lemma 4).}
\end{aligned}
$$

(78)

Plugging Eq. (72) into Eq. (78), we have

$$
\begin{aligned}
&\mathop{\mathbb{E}}_{1,2,\cdots,k} \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k} \right\|^2 \\
&= \left( (1 - \alpha_{(i),t})^2 + \frac{\alpha_{(i),t}^2 (p+1)}{\tilde{n}_{(i),t}} \right) \mathop{\mathbb{E}}_{1,2,\cdots,k-1} \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k-1} \right\|^2 + \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2 \\
&\quad + \alpha_{(i),t}^2 \frac{p}{\tilde{n}_{(i),t}} \sigma_{(i),t}^2 + 2\alpha_{(i),t} \left( 1 - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) (1 - \alpha_{(i),t})^{k-1} \boldsymbol{\gamma}_{(i),t}^\top \Delta_{t-1}^{K<\infty} \\
&\quad + 2\alpha_{(i),t} \left( 1 - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) \left( 1 - (1 - \alpha_{(i),t})^{k-1} \right) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2 \\
&= \mathcal{A}_{(i),t} \, \mathbb{E} \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t,k-1} \right\|^2 + \mathcal{B}'_{(i),t,k},
\end{aligned}
\tag{79}
$$

where $\mathcal{A}_{(i),t}$ is defined in Eq. (65) and

$$
\begin{aligned}
&\mathcal{B}'_{(i),t,k} \\
&:= \frac{\alpha_{(i),t}^2 p \sigma_{(i),t}^2}{\tilde{n}_{(i),t}} \\
&\quad + \left( \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) + 2\alpha_{(i),t} \left( 1 - \frac{\alpha_{(i),t}}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) \left( 1 - (1 - \alpha_{(i),t})^{k-1} \right) \right) \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2 \\
&\quad + 2 \left( \alpha_{(i),t} - \frac{\alpha_{(i),t}^2}{\tilde{n}_{(i),t}} (\tilde{n}_{(i),t} + p + 1) \right) (1 - \alpha_{(i),t})^{k-1} \boldsymbol{\gamma}_{(i),t}^\top \Delta_{t-1}^{K<\infty}.
\end{aligned}
$$

We also define $\mathcal{B}_{(i),t,k}$ by replacing $\Delta_{t-1}^{K<\infty}$ in $\mathcal{B}'_{(i),t,k}$ with $\mathcal{F}_{t-1}$, i.e., Eq. (66).

Applying Eq. (79) recursively over $k = 1, 2, \cdots, K$, we thus have

$$
\mathop{\mathbb{E}}_{t} \left\| \boldsymbol{w}^* - \hat{\boldsymbol{w}}_{(i),t} \right\|^2 = \mathcal{A}_{(i),t}^K \left\| \Delta_{t-1}^{K<\infty} \right\|^2 + \sum_{k=1}^{K} \mathcal{B}_{(i),t,k} \mathcal{A}_{(i),t}^{K-k}.
\tag{80}
$$

Plugging Eqs. (77) and (80) into Eq. (75), we thus have

$$
\mathbb{E} \left\| \Delta_t^{K<\infty} \right\|^2 = \mathcal{J}_t \, \mathbb{E} \left\| \Delta_{t-1}^{K<\infty} \right\|^2 + Q_t,
\tag{81}
$$

where $\mathcal{J}_t$ is defined in Eq. (67) and $Q_t$ is defined in Eq. (68).

Applying Eq. (81) recursively, we thus have Eq. (20). □

## D PROOF OF THEOREM 3

PROOF. In the overparameterized situation, after each agent trains to converge, we have

$$
\hat{\boldsymbol{w}}_{(i),t}^{K=\infty} = \mathbf{X}_{(i),t} \left( \mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t} \right)^{-1} \left( \boldsymbol{y}_{(i),t} - \mathbf{X}_{(i),t}^\top \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=\infty} \right) + \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=\infty}.
\tag{82}
$$

For any $i \in [m]$, we define $\mathbf{P}_{(i),t} \in \mathbb{R}^{p \times p}$ as

$$
\mathbf{P}_{(i),t} := \mathbf{X}_{(i),t} \left( \mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t} \right)^{-1} \mathbf{X}_{(i),t}^\top.
\tag{83}
$$

(We know $\mathbf{P}_{(i),t}$ is an orthogonal projection since $\mathbf{P}_{(i),t} \mathbf{P}_{(i),t} = \mathbf{P}_{(i),t}$ and $\mathbf{P}_{(i),t}^\top = \mathbf{P}_{(i),t}$.) By Eqs. (2), (82) and (83), we thus have

$$
\hat{\boldsymbol{w}}_{(i),t}^{K=\infty} = \mathbf{P}_{(i),t} \boldsymbol{w}_{(i),t} + (\mathbf{I}_p - \mathbf{P}_{(i),t}) \hat{\boldsymbol{w}}_{\text{avg},t-1}^{K=\infty} + \mathbf{X}_{(i),t} \left( \mathbf{X}_{(i),t}^\top \mathbf{X}_{(i),t} \right)^{-1} \boldsymbol{\epsilon}_{(i),t}.
\tag{84}
$$

We thus have

$$\Delta_t^{K=\infty}$$

$$= w^* - \hat{w}_{\text{avg},t}^{K=\infty} \quad \text{(by Eq. (8))}$$

$$= w^* - \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} \left( P_{(i),t} w_{(i),t} + (I_p - P_{(i),t})\hat{w}_{\text{avg},t-1}^{K=\infty} + X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right)$$

(by Eqs. (5) and (84))

$$= \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} \left( P_{(i),t}(w^* - w_{(i),t}) + (I_p - P_{(i),t})(w^* - \hat{w}_{\text{avg},t-1}^{K=\infty}) - X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right)$$

$$\left( \text{since } w^* = \frac{\sum_{i\in[m]} n_{(i),t}(P_{(i),t} + I_p - P_{(i),t})w^*}{\sum_{i\in[m]} n_{(i),t}} \right)$$

$$= \frac{1}{\sum_{i\in[m]} n_{(i),t}} \sum_{i\in[m]} n_{(i),t} \left( P_{(i),t} \gamma_{(i),t} + (I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} - X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right)$$

(by Eqs. (3) and (8)). $\hspace{8cm}$ (85)

For any $i, j \in [m]$, because $\epsilon_{(j),t}$ is independent of $\Delta_{t-1}^{K=\infty}$ and $X_{(i),t}$, and also because $\epsilon_{(j),t}$ has zero mean (by Assumption 1), we have

$$\mathbb{E}\left[ \left( P_{(i),t} \gamma_{(i),t} \right)^\top X_{(j),t} \left( X_{(j),t}^\top X_{(j),t} \right)^{-1} \epsilon_{(j),t} \right]$$

$$= \mathbb{E}\left[ ((I_p - P_{(i),t})\Delta_{t-1}^{K=\infty})^\top X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right]$$

$$= 0, \hspace{10cm} (86)$$

and

$$\mathbb{E}\left[ X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right] = 0. \hspace{6cm} (87)$$

Since $P_{(i),t}(I_p - P_{(i),t}) = 0$, we have

$$\left( P_{(i),t} \gamma_{(i),t} \right)^\top (I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} = 0. \hspace{6cm} (88)$$

Thus, by Eqs. (85), (86) and (88), we have

$$\mathbb{E}_t \left\| \Delta_t^{K=\infty} \right\|^2$$

$$= \frac{\sum_{i\in[m]} n_{(i),t}^2 \left( \mathbb{E}_t \left\| (I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} \right\|^2 + \mathbb{E}_t \left\| P_{(i),t} \gamma_{(i),t} \right\|^2 + \mathbb{E}_t \left\| X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right\|^2 \right)}{\left( \sum_{i\in[m]} n_{(i),t} \right)^2}$$

$$+ \frac{1}{\left( \sum_{i\in[m]} n_{(i),t} \right)^2} \sum_{i\in[m]} \sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t} \left( \gamma_{(j),t}^\top P_{(j),t} P_{(i),t} \gamma_{(i),t} \right.$$

$$+ \Delta_{t-1}^{K=\infty\top}(I_p - P_{(j),t})(I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} + 2\gamma_{(j),t}^\top P_{(j),t}(I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} \right). \hspace{2cm} (89)$$

For any $i \in [m]$, we have

$$\mathbb{E}_t \left\| P_{(i),t} \gamma_{(i),t} \right\|^2 = \frac{n_{(i),t}}{p} \left\| \gamma_{(i),t} \right\|^2 \quad \text{(by Lemma 2)}, \hspace{4cm} (90)$$

$$\mathbb{E}_t \left\| (I_p - P_{(i),t})\Delta_{t-1}^{K=\infty} \right\|^2 = \left( 1 - \frac{n_{(i),t}}{p} \right) \left\| \Delta_{t-1}^{K=\infty} \right\|^2 \quad \text{(by Lemma 2)}, \hspace{2cm} (91)$$

$$\mathbb{E}_t \left\| X_{(i),t} \left( X_{(i),t}^\top X_{(i),t} \right)^{-1} \epsilon_{(i),t} \right\|^2 = \frac{n_{(i),t}\sigma_i^2}{p - n_{(i),t} - 1} \quad \text{(by Lemma 3)}. \hspace{2cm} (92)$$

For any $i, j \in [m]$ where $i \neq j$, we have

$$\mathbb{E}_t \left[ \Delta_{t-1}^{K=\infty}{}^\top (\mathbf{I}_p - \mathbf{P}_{(j),t})(\mathbf{I}_p - \mathbf{P}_{(i),t}) \Delta_{t-1}^{K=\infty} \right]$$

$$= \mathbb{E}_t \left[ (\mathbf{I}_p - \mathbf{P}_{(i),t}) \Delta_{t-1}^{K=\infty} \right]^\top \mathbb{E}_t \left[ (\mathbf{I}_p - \mathbf{P}_{(j),t}) \Delta_{t-1}^{K=\infty} \right]$$

$$(\text{since } \mathbf{P}_{(i),t} \text{ and } \mathbf{P}_{(j),t} \text{ are independent when } i \neq j)$$

$$= \left( 1 - \frac{n_{(i),t}}{p} \right) \left( 1 - \frac{n_{(j),t}}{p} \right) \left\| \Delta_{t-1}^{K=\infty} \right\|^2 \quad (\text{by Lemma 5}). \tag{93}$$

Similarly, for $i \neq j$, we have

$$\mathbb{E}_t \left[ \boldsymbol{\gamma}_{(j),t}^\top \mathbf{P}_{(j),t} \mathbf{P}_{(i),t} \boldsymbol{\gamma}_{(i),t} \right] = \frac{n_{(i),t} n_{(j),t}}{p^2} \boldsymbol{\gamma}_{(j),t}^\top \boldsymbol{\gamma}_{(i),t} \quad (\text{by Lemma 5}), \tag{94}$$

and

$$\mathbb{E}_t \left[ \boldsymbol{\gamma}_{(j),t}^\top \mathbf{P}_{(j),t} (\mathbf{I}_p - \mathbf{P}_{(i),t}) \Delta_{t-1}^{K=\infty} \right] = \frac{n_{(j),t}}{p} \left( 1 - \frac{n_{(i),t}}{p} \right) \boldsymbol{\gamma}_{(j),t}^\top \Delta_{t-1}^{K=\infty} \quad (\text{by Lemma 5}). \tag{95}$$

Plugging Eqs. (93) to (95) and (90) to (92) into Eq. (89), we thus have

$$\mathbb{E}_t \left\| \Delta_t^{K=\infty} \right\|^2$$

$$= \frac{\sum_{i \in [m]} n_{(i),t}^2 \left( \left( 1 - \frac{n_{(i),t}}{p} \right) \left\| \Delta_{t-1}^{K=\infty} \right\|^2 + \frac{n_{(i),t}}{p} \left\| \boldsymbol{\gamma}_{(i),t} \right\|^2 + \frac{n_{(i),t} \sigma_{(i),t}^2}{p - n_{(i),t} - 1} \right)}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2}$$

$$+ \frac{1}{\left( \sum_{i \in [m]} n_{(i),t} \right)^2} \sum_{i \in [m]} \sum_{j \in [m] \setminus \{i\}} n_{(i),t} n_{(j),t} \left( \frac{n_{(i),t} n_{(j),t}}{p^2} \boldsymbol{\gamma}_{(j),t}^\top \boldsymbol{\gamma}_{(i),t} \right.$$

$$\left. + \left( 1 - \frac{n_{(i),t}}{p} \right) \left( 1 - \frac{n_{(j),t}}{p} \right) \left\| \Delta_{t-1}^{K=\infty} \right\|^2 + 2 \frac{n_{(j),t}}{p} \left( 1 - \frac{n_{(i),t}}{p} \right) \boldsymbol{\gamma}_{(j),t}^\top \Delta_{t-1}^{K=\infty} \right). \tag{96}$$

By Eq. (85), we also have

$$\mathbb{E}_t [\Delta_t^{K=\infty}] = \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} n_{(i),t} \left( \frac{n_{(i),t}}{p} \boldsymbol{\gamma}_{(i),t} + \left( 1 - \frac{n_{(i),t}}{p} \right) \Delta_{t-1}^{K=\infty} \right). \tag{97}$$

Applying Eq. (97) recursively, we thus have

$$\mathbb{E}[\Delta_t^{K=\infty}] = g_t^{K=\infty}, \tag{98}$$

where $g_t^{K=\infty}$ is defined in Eq. (23).

By Eqs. (96) and (98), we thus have

$$\mathbb{E} \left\| \Delta_t^{K=\infty} \right\|^2 = C_t \cdot \mathbb{E} \left\| \Delta_{t-1}^{K=\infty} \right\|^2 + D_t, \tag{99}$$

where $C_t$ denotes the coefficient of $\left\| \Delta_{t-1}^{K=\infty} \right\|^2$ and $D_t$ denotes the remaining parts. The specific expressions of $C_t$ and $D_t$ are in Eqs. (26) and (27). Applying Eq. (99) recursively, we thus have Eq. (28).

**Underparameterized situation**

In the underparameterized situation, the convergence point of local steps in each round corresponds to the solution that minimizes the training loss, i.e.,

$$\hat{\boldsymbol{w}}_{(i),t}^{K=\infty} = (\mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top)^{-1} \mathbf{X}_{(i),t} \boldsymbol{y}_{(i),t}$$

$$= (\mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top)^{-1} \mathbf{X}_{(i),t} (\mathbf{X}_{(i),t}^\top \boldsymbol{w}_{(i),t} + \boldsymbol{\epsilon}_{(i),t}) \quad (\text{by Eq. (2)})$$

$$= \boldsymbol{w}_{(i),t} + (\mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top)^{-1} \mathbf{X}_{(i),t} \boldsymbol{\epsilon}_{(i),t}.$$

Also recalling Eqs. (3) and (8), we thus have

$$\Delta_t^{K=\infty} = \frac{1}{\sum_{i \in [m]} n_{(i),t}} \sum_{i \in [m]} n_{(i),t} (\boldsymbol{\gamma}_{(i),t} - (\mathbf{X}_{(i),t} \mathbf{X}_{(i),t}^\top)^{-1} \mathbf{X}_{(i),t} \boldsymbol{\epsilon}_{(i),t}). \tag{100}$$

For any $i, j \in [m]$, because $\epsilon_{(j),t}$ is independent of $\mathbf{X}_{(i),t}$ and $\epsilon_{(i),t}$, and also because $\epsilon_{(j),t}$ has zero mean (by Assumption 1), we have

$$\mathbb{E}\left[\boldsymbol{\gamma}_{(j),t}^{\top}(\mathbf{X}_{(i),t}\mathbf{X}_{(i),t}^{\top})^{-1}\mathbf{X}_{(i),t}\epsilon_{(i),t}\right] = 0 \ \text{ for all } i, j \in [m],$$

$$\mathbb{E}\left[\left(\mathbf{X}_{(j),t}\mathbf{X}_{(j),t}^{\top}\right)^{-1}\mathbf{X}_{(j),t}\epsilon_{(j),t}\right)^{\top}(\mathbf{X}_{(i),t}\mathbf{X}_{(i),t}^{\top})^{-1}\mathbf{X}_{(i),t}\epsilon_{(i),t}\right] = 0 \ \text{ for all } i \neq j.$$

Thus, by Eq. (100), we have

$$\mathbb{E}\left\|\Delta_t^{K=\infty}\right\|^2 = \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2}\sum_{i\in[m]} n_{(i),t}^2\left(\left\|\boldsymbol{\gamma}_{(i),t}\right\|^2 + \mathbb{E}\left\|(\mathbf{X}_{(i),t}\mathbf{X}_{(i),t}^{\top})^{-1}\mathbf{X}_{(i),t}\epsilon_{(i),t}\right\|^2\right)$$

$$+ \frac{1}{(\sum_{i\in[m]} n_{(i),t})^2}\sum_{i\in[m]}\sum_{j\in[m]\setminus\{i\}} n_{(i),t} n_{(j),t}\boldsymbol{\gamma}_{(i),t}^{\top}\boldsymbol{\gamma}_{(j),t}$$

$$= \left\|\frac{\sum_{i\in[m]} n_{(i),t}\boldsymbol{\gamma}_{(i),t}}{\sum_{i\in[m]} n_{(i),t}}\right\|^2 + \frac{\sum_{i\in[m]}\frac{n_{(i),t}^2 p\sigma_{(i),t}^2}{n_{(i),t}-p-1}}{(\sum_{i\in[m]} n_{(i),t})^2} \ \text{ (by Eq. (46) in Lemma 3).}$$

We thus have proven Eq. (29).

The result of this theorem thus follows.                                                                                                                $\square$

# E   A TABLE FOR NOTATIONS

We provide a table of some important notations used in this paper.

| symbol | meaning |
|---|---|
| $n_{(i),t}$ | number of training samples |
| $\tilde{n}_{(i),t}$ | batch size |
| $p$ | number of parameters |
| $\sigma_{(i),t}$ | noise level |
| $\mathbf{X}_{(i),t}$ | matrix for input of training samples |
| $\boldsymbol{y}_{(i),t}$ | vector for output of training samples |
| $\epsilon_{(i),t}$ | vector for noise of training samples |
| $\hat{\boldsymbol{w}}_0$ | the pre-trained parameters (initialization) |
| $\boldsymbol{w}^*$ | the learning target |
| $\boldsymbol{w}_{(i),t}$ | the ground-truth of agent $i$ at round $t$ |
| $\hat{\boldsymbol{w}}_{(i),t}^{K=1}, \hat{\boldsymbol{w}}_{(i),t}, \hat{\boldsymbol{w}}_{(i),t}^{K=\infty}$ | the local learning result of agent $i$ at round $t$ |
| $\hat{\boldsymbol{w}}_{(i),t,k}$ | learning result after $k$-th batch (for $K < \infty$ case) |
| $\hat{\boldsymbol{w}}_{\mathrm{avg},t}^{K=1}, \hat{\boldsymbol{w}}_{\mathrm{avg},t}^{K<\infty}, \hat{\boldsymbol{w}}_{\mathrm{avg},t}^{K=\infty}$ | the FedAvg result at round $t$ |
| $\left\|\Delta_t^{K=1}\right\|^2, \left\|\Delta_t^{K<\infty}\right\|^2, \left\|\Delta_t^{K=\infty}\right\|^2$ | model error |
| $\left\|\Delta_0\right\|^2$ | initial (pre-trained) model error |
| $\alpha_{(i),t}$ | learning rate (step size) |
| $\boldsymbol{\gamma}_{(i),t}$ | measurement of heterogeneity |

**Table 2: Table for some notations.**