
Anonymous Author(s)

Affiliation

Address

email

1 Introduction

In the federated learning and decentralized learning, n participants collaborate to train a global model \mathbf{x} over their joint objectives $\min_{\mathbf{x}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$. Compared to models trained on individual data silos, this model achieves overall better performance on dataset. There is no guarantee that it performs better than standalone training on some workers. Personalized federated learning is one way to address this problem.

2 Shared setting

(A1) *L-smoothness.* For $i \in [n]$, f_i is L -smooth.

(A2) *Lower bound.* For $i \in [n]$, f_i is lower bounded by f_i^* .

3 Dynamic graph and full gradient

3.1 Problem formulation

In this section, we assume not all workers share same stationary points or minimizers.

(A3) **Strong growth condition.** Let $c \subset [n]$ be the a subset of workers that share same stationary point. Then for $\mathbf{x} \in \mathbb{R}^d$ and $i \in c$, we have

$$\|\nabla f_i(\mathbf{x}) - \nabla \bar{f}_c(\mathbf{x})\|_2 \leq M \|\nabla \bar{f}_c(\mathbf{x})\|_2.$$

The (A3) indicates that when an iterate \mathbf{x} reaches the stationary point of \bar{f}_c , then it also reaches the stationary point of f for all $i \in c$. The optimization objective is that

$$\min_{\mathbf{X} \in \mathbb{R}^{d \times n}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_i) + \frac{\rho}{2} \sum_{i < j} w_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2.$$

We optimize the objective with gradient descent with initialization $\mathbf{x}_i^0 = \bar{\mathbf{x}}^0$

$$\mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \eta \left(\nabla f_i(\mathbf{x}_i^t) + \rho \sum_{k=1}^n w_{ik}^t (\mathbf{x}_i^t - \mathbf{x}_k^t) \right). \quad (1)$$

We update w_{ij}^t with the following term

$$w_{ij}^{t+1} = \text{sign}(\alpha - \|\mathbf{x}_i^t - \mathbf{x}_j^t\|_2^2)$$

3.2 Proof Sketch

Notations. We define the following notations

- Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ be the compact form of iterates.

- 22 • Let $\nabla F(\mathbf{X}) := [\nabla f_1(\mathbf{x}_1), \nabla f_2(\mathbf{x}_2), \dots, \nabla f_n(\mathbf{x}_n)]^\top \in \mathbb{R}^{n \times d}$.
- 23 • Let \mathbf{W}^t be the mixing matrix at time t and $\mathbf{D}^t := \text{Diag}(\mathbf{W}^t \mathbf{1})$.
- 24 • Let \mathbf{W}^* be the groundtruth mixing matrix.
- 25 • Let $d_{\max} = \max\{D^*\}$ be the size of largest cluster.

26 Then we update the iterates as follows

$$\mathbf{X}^{t+1} = (\mathbf{I} - \eta \rho (\mathbf{D}^t - \mathbf{W}^t)) \mathbf{X}^t - \eta \nabla F(\mathbf{X}^t). \quad (2)$$

27 **Lemma 1** (Basic properties). *The following equality and inequalities hold true*

- 28 • $(\mathbf{D}^* - \mathbf{W}^*)(\mathbf{D}^* - \mathbf{W}^*) = \mathbf{D}^*(\mathbf{D}^* - \mathbf{W}^*) = (\mathbf{D}^* - \mathbf{W}^*)\mathbf{D}^*$
- 29 • $\|AB\|_F \leq \|A\|_2 \|B\|_F$

30 Note that (1) can be re-written as

$$\mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \eta \left(\nabla f_i(\mathbf{x}_i^t) + \rho \sum_{k=1}^n w_{ik}^* (\mathbf{x}_i^t - \mathbf{x}_k^t) + \rho \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right).$$

31 After averaging i over the same cluster, we have

$$\bar{\mathbf{x}}_c^{t+1} = \bar{\mathbf{x}}_c^t - \eta \left(\frac{1}{c} \sum_{i \in c} \nabla f_i(\mathbf{x}_i^t) + \frac{\rho}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right). \quad (3)$$

32 **Lemma 2** (Sufficient decrease). *Suppose f_i are L -smooth, then by taking $\eta \leq \frac{1}{L}$, we have*

$$\bar{f}_c(\bar{\mathbf{x}}_c^{t+1}) \leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 + \frac{\eta L^2}{c} \sum_{i \in c} \|\mathbf{x}_i^t - \bar{\mathbf{x}}_c^t\|_2^2 + \eta \left\| \frac{\rho}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2.$$

33 Here are related equality

- 34 • $\sum_c \sum_{i \in c} \|\mathbf{x}_i^t - \bar{\mathbf{x}}_c^t\|_2^2 = \|(\mathbf{I} - \mathbf{D}^{-1} \mathbf{W}^*) \mathbf{X}^t\|_F^2$
- 35 • $\sum_c c \left\| \frac{1}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 = \|\mathbf{D}^{-1} \mathbf{W}^* (\mathbf{D}^* - \mathbf{W}^* - \mathbf{D}^t + \mathbf{W}^t) \mathbf{X}^t\|_F^2$
- 36 • $\sum_c \sum_{i \in c} \left\| \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 = \|(\mathbf{D}^* - \mathbf{W}^* - \mathbf{D}^t + \mathbf{W}^t) \mathbf{X}^t\|_F^2$

37 Let $\mathcal{E}_i := \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t)$ be the Misclassification term of worker i .

38 **Lemma 3** (Misclassification error). *Let $\mathcal{E} := \|(\mathbf{D}^* - \mathbf{W}^* - \mathbf{D}^t + \mathbf{W}^t) \mathbf{X}^t\|_F^2$ be the error incurred*
 39 *by misclassification*

$$\mathcal{E} \leq 2\mathcal{E}_{\text{ex}} + 2\mathcal{E}_{\text{in}}.$$

40 where $\mathcal{E}_{\text{ex}} = \sum_c \sum_{i \in c} c \sum_{k \in c, w_{ik}^t = 0} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2$ and $\mathcal{E}_{\text{in}} = \alpha^t n(n-c)^2$.

41 Note that as we initialize all models to be the same, the \mathcal{E}_{ex} is small in the beginning; as we
 42 choose $\alpha^t = \mathcal{O}(\frac{1}{t})$, the \mathcal{E}_{in} is gradually decreasing. The inclusion error eventually vanishes due
 43 to heterogeneity across clusters, in which case we can stop decreasing α^t . The exclusion error is
 44 bounded by generalized strong-growth condition.

45 **Lemma 4.**

Proof.

$$\begin{aligned} \|\mathbf{x}_i^{t+1} - \mathbf{x}_k^{t+1}\|_2^2 &= \|\mathbf{x}_i^t - \mathbf{x}_k^t - \eta(\nabla f_i(\mathbf{x}_i^t) - \nabla f_k(\mathbf{x}_k^t)) - \eta \rho c (\mathbf{x}_i^t - \mathbf{x}_k^t) - \eta \rho (\mathcal{E}_i - \mathcal{E}_k)\|_2^2 \\ &\leq (1 - \eta \rho c) \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 + \eta \rho c \left\| \frac{1}{\rho c} (\nabla f_i(\mathbf{x}_i^t) - \nabla f_k(\mathbf{x}_k^t)) + \frac{1}{c} (\mathcal{E}_i - \mathcal{E}_k) \right\|_2^2. \end{aligned}$$

46 Note that

$$\begin{aligned} \|\nabla f_i(\mathbf{x}_i^t) - \nabla f_k(\mathbf{x}_k^t)\|_2^2 &\leq \|\nabla f_i(\mathbf{x}_i^t) \pm \nabla f_i(\bar{\mathbf{x}}_c^t) \pm \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t) \pm \nabla f_k(\bar{\mathbf{x}}_c^t) - \nabla f_k(\mathbf{x}_k^t)\|_2^2 \\ &\leq 4L^2 \|\mathbf{x}_i^t - \bar{\mathbf{x}}_c^t\|_2^2 + 4L^2 \|\mathbf{x}_k^t - \bar{\mathbf{x}}_c^t\|_2^2 + 8M^2 \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2. \end{aligned}$$

47 Then

$$\frac{1}{c^2} \sum_{i \in c} \sum_{k \in c} \|\nabla f_i(\mathbf{x}_i^t) - \nabla f_k(\mathbf{x}_k^t)\|_2^2 \leq 8L^2 \frac{1}{c} \sum_{k \in c} \|\mathbf{x}_k^t - \bar{\mathbf{x}}_c^t\|_2^2 + 8M^2 \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2.$$

48 Then

$$\begin{aligned} \frac{1}{c^2} \sum_{i \in c} \sum_{k \in c} \|\mathbf{x}_i^{t+1} - \mathbf{x}_k^{t+1}\|_2^2 &\leq (1 - \eta\rho c) \frac{1}{c^2} \sum_{i \in c} \sum_{k \in c} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 \\ &\quad + 2 \frac{\eta}{\rho c} \left(8L^2 \frac{1}{c} \sum_{k \in c} \|\mathbf{x}_k^t - \bar{\mathbf{x}}_c^t\|_2^2 + 8M^2 \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 \right) \\ &\quad + 2 \frac{\eta\rho}{c} \frac{1}{c^2} \sum_{i \in c} \sum_{k \in c} \|\mathcal{E}_i - \mathcal{E}_k\|_2^2. \end{aligned}$$

49 □

50 **Lemma 5** (Consensus distance). *The consensus distance measures the worker model's distance to*
 51 *the center of its own cluster, i.e. $\|(I - D^{-1}W^*)X\|_F^2$ in the compact form.*

52 *Proof.* The distance to their own center is

$$\begin{aligned} (D^* - W^*)X^{t+1} &= (D^* - W^*)(I - \eta\rho(D^t - W^t))X^t - \eta(D^* - W^*)\nabla F(X^t) \\ &= (I - \eta\rho(D^* - W^*))(D^* - W^*)X^t \\ &\quad + \eta\rho(D^* - W^*)(D^* - W^* - D^t + W^t)X^t \\ &\quad - \eta(D^* - W^*)\nabla F(X^t) \\ &= (I - \eta\rho D^*)(D^* - W^*)X^t \\ &\quad + \eta\rho(D^* - W^*)(D^* - W^* - D^t + W^t)X^t \\ &\quad - \eta(D^* - W^*)\nabla F(X^t). \end{aligned}$$

53 Multiply with D^{-1} to both sides yield

$$\begin{aligned} (I - D^{-1}W^*)X^{t+1} &= (I - \eta\rho D^*)(I - D^{-1}W^*)X^t \\ &\quad + \eta\rho(I - D^{-1}W^*)(D^* - W^* - D^t + W^t)X^t \\ &\quad - \eta(I - D^{-1}W^*)\nabla F(X^t). \end{aligned}$$

54 Applying Frobenius norm to the above equation, we have

$$\begin{aligned} &\|(I - D^{-1}W^*)X^{t+1}\|_F^2 \\ &= 3\|(I - \eta\rho D^*)(I - D^{-1}W^*)X^t\|_F^2 + 3\|\eta\rho(I - D^{-1}W^*)(D^* - W^* - D^t + W^t)X^t\|_F^2 \\ &\quad + 3\|\eta(I - D^{-1}W^*)\nabla F(X^t)\|_F^2 \\ &= 3\|I - \eta\rho D^*\|_2^2 \|(I - D^{-1}W^*)X^t\|_F^2 + 3\|\eta\rho(I - D^{-1}W^*)(D^* - W^* - D^t + W^t)X^t\|_F^2 \\ &\quad + 3\|\eta(I - D^{-1}W^*)\nabla F(X^t)\|_F^2 \\ &\leq (1 - \eta\rho d_{\max}) \|(I - D^{-1}W^*)X^t\|_F^2 + 3\|\eta\rho(I - D^{-1}W^*)(D^* - W^* - D^t + W^t)X^t\|_F^2 \\ &\quad + 3\|\eta(I - D^{-1}W^*)\nabla F(X^t)\|_F^2 \end{aligned}$$

55 where we use $\rho \geq \frac{2}{3\eta d_{\max}}$. The second term can be bounded as follows

$$\|(I - D^{-1}W^*)(D^* - W^* - D^t + W^t)X^t\|_F^2 \leq \|I - D^{-1}W^*\|_2^2 \mathcal{E}.$$

56 The last term can be bounded as follows

$$\|(I - D^{-1}W^*)\nabla F(X^t)\|_F^2$$

57 □

58 A Proofs

59 *Proof.* Apply L-smoothness to each function f_i , we have

$$f_i(\bar{\mathbf{x}}_c^{t+1}) \leq f_i(\bar{\mathbf{x}}_c^t) + \langle \nabla f_i(\bar{\mathbf{x}}_c^t), \bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t \rangle + \frac{L}{2} \|\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t\|_2^2.$$

60 Average the above inequality over $i \in c$, we have

$$\bar{f}_c(\bar{\mathbf{x}}_c^{t+1}) \leq \bar{f}_c(\bar{\mathbf{x}}_c^t) + \langle \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t), \bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t \rangle + \frac{L}{2} \|\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t\|_2^2.$$

61 Expand the intermediate term, we have

$$\begin{aligned} \langle \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t), \bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t \rangle &= -\frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 - \frac{1}{2\eta} \|\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t\|_2^2 \\ &\quad + \frac{\eta}{2} \left\| \frac{\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t}{\eta} - \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t) \right\|_2^2. \end{aligned}$$

62 Then by having $\eta \leq \frac{1}{L}$

$$\begin{aligned} \bar{f}_c(\bar{\mathbf{x}}_c^{t+1}) &\leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 + \frac{\eta}{2} \left\| \frac{\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t}{\eta} - \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t) \right\|_2^2 \\ &\quad - \frac{1-L\eta}{2\eta} \|\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t\|_2^2 \\ &\leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 + \frac{\eta}{2} \left\| \frac{\bar{\mathbf{x}}_c^{t+1} - \bar{\mathbf{x}}_c^t}{\eta} - \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t) \right\|_2^2 \end{aligned}$$

63 Using (3) we have that

$$\begin{aligned} \bar{f}_c(\bar{\mathbf{x}}_c^{t+1}) &\leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 \\ &\quad + \frac{\eta}{2} \left\| \frac{1}{c} \sum_{i \in c} \nabla f_i(\mathbf{x}_i^t) - \nabla \bar{f}_c(\bar{\mathbf{x}}_c^t) + \frac{\rho}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*)(\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 \\ &\leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 + \eta \left\| \frac{1}{c} \sum_{i \in c} (\nabla f_i(\mathbf{x}_i^t) - \nabla f_i(\bar{\mathbf{x}}_c^t)) \right\|_2^2 \\ &\quad + \eta \left\| \frac{\rho}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*)(\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2. \end{aligned}$$

64 Using the L -smoothness of f_i

$$\begin{aligned} \bar{f}_c(\bar{\mathbf{x}}_c^{t+1}) &\leq \bar{f}_c(\bar{\mathbf{x}}_c^t) - \frac{\eta}{2} \|\nabla \bar{f}_c(\bar{\mathbf{x}}_c^t)\|_2^2 + \frac{\eta L^2}{c} \sum_{i \in c} \|\mathbf{x}_i^t - \bar{\mathbf{x}}_c^t\|_2^2 \\ &\quad + \eta \left\| \frac{\rho}{c} \sum_{i \in c} \sum_{k=1}^n (w_{ik}^t - w_{ik}^*)(\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2. \end{aligned}$$

65

□

66 **A.1 Misclassification Error**

Proof.

$$\begin{aligned}
\mathcal{E} &= \sum_c \sum_{i \in c} \left\| \sum_{k=1}^n (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 \\
&= \sum_c \sum_{i \in c} \left\| \sum_{k \in c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) + \sum_{k \notin c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 \\
&\leq 2 \underbrace{\sum_c \sum_{i \in c} \left\| \sum_{k \in c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2}_{\text{Exclusion Error } \mathcal{E}_{\text{ex}}} + 2 \underbrace{\sum_c \sum_{i \in c} \left\| \sum_{k \notin c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2}_{\text{Inclusion Error } \mathcal{E}_{\text{in}}}
\end{aligned}$$

67 The exclusion error

$$\begin{aligned}
\mathcal{E}_{\text{ex}} &= \sum_c \sum_{i \in c} \left\| \sum_{k \in c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 \\
&\leq \sum_c \sum_{i \in c} \left(\sum_{k \in c} (w_{ik}^* - w_{ik}^t) \right) \left(\sum_{k \in c} (w_{ik}^* - w_{ik}^t) \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 \right) \\
&\leq \sum_c \sum_{i \in c} c \sum_{k \in c \& w_{ik}^t = 0} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2.
\end{aligned}$$

68 This may be further bounded as a function of α_t (but we may not use it)

$$\mathcal{E}_{\text{ex}} = \sum_c \sum_{i \in c} c \frac{\left(\sum_{k \in c \& w_{ik}^t = 0} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 \right)^2}{\alpha^t} \text{ Or } \frac{\left(\sum_c \sum_{i \in c} c \sum_{k \in c \& w_{ik}^t = 0} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 \right)^2}{\alpha^t}.$$

69 On the other hand, the inclusion error can be bounded as follows

$$\begin{aligned}
\mathcal{E}_{\text{in}} &= \sum_c \sum_{i \in c} \left\| \sum_{k \notin c} (w_{ik}^t - w_{ik}^*) (\mathbf{x}_i^t - \mathbf{x}_k^t) \right\|_2^2 \\
&\leq \sum_c \sum_{i \in c} (n - c) \sum_{k \notin c \& w_{ik}^t = 1} \|\mathbf{x}_i^t - \mathbf{x}_k^t\|_2^2 \\
&\leq \alpha^t n (n - c)^2.
\end{aligned}$$

70

□