

# FedAlign: Differentially Private Distribution Alignment for Non-IID Federated Learning

## Supplementary Material

### A. Experimental Settings

#### A.1. Experimental Environments

All experiments were conducted on a workstation equipped with an Intel(R) Xeon(R) Bronze 3206R CPU @ 1.90 GHz (16 cores) and an NVIDIA GeForce RTX 3090 GPU with 24 GB of memory. The system runs CentOS with Python 3.8. PyTorch (version 2.0.0) was used to implement *FedAlign* and all baseline models, ensuring a consistent and reproducible experimental setup.

#### A.2. Implementation Details

All experiments are conducted with a fixed set of 10 clients, where all clients participate in every communication round (participation ratio = 1). Each client performs 1 local epoch per round, and the training batch size is set to 10. Models are trained using the stochastic gradient descent (SGD) optimizer with a fixed learning rate of 0.01 and a momentum coefficient of 0.9.

The experiments are designed to evaluate the performance of *FedAlign* under both system and data heterogeneity conditions.

**System Heterogeneity:** To simulate variations in computational capabilities among clients, a hyperparameter delay is introduced. Larger values of delay correspond to higher disparities in client processing speeds. In the reported experiments, the delay is set to 0.1 to represent moderate system heterogeneity.

**Data Heterogeneity:** Two types of Non-IID scenarios are considered. First, label and quantity skew are induced by partitioning each clients data according to a Dirichlet distribution with concentration parameter  $\alpha = 0.5$ . Second, feature skew is modeled by adding Gaussian noise  $\mathcal{N}(0, \beta)$  to the training and testing datasets, with  $\beta \in \{0.05, 0.1\}$ .

All reported results represent the mean and standard deviation over five independent runs, excluding the maximum and minimum performance values. The best performance in each setting is highlighted in bold. Other hyperparameters, including dropout rate and data partitioning strategy, are uniformly configured across all clients to ensure reproducibility.

### B. Detailed Proofs

#### B.1. Proof of Theorem 1

*Proof.* Let  $\tilde{g}_k$  be the stochastic gradient of client  $k$ , and let the global DP gradient be

$$\tilde{g}^{\text{DP}} = \sum_{j=1}^K \frac{n_j}{N} \tilde{g}_j + \eta. \quad (15)$$

Using the identity

$$\|\tilde{g}_k - \tilde{g}^{\text{DP}}\|^2 = \|\tilde{g}_k\|^2 + \|\tilde{g}^{\text{DP}}\|^2 - 2\langle \tilde{g}_k, \tilde{g}^{\text{DP}} \rangle, \quad (16)$$

and taking expectations yields

$$\mathbb{E}\|\tilde{g}_k - \tilde{g}^{\text{DP}}\|^2 = \mathbb{E}\|\tilde{g}_k\|^2 - 2\mathbb{E}\langle \tilde{g}_k, \tilde{g}^{\text{DP}} \rangle + \mathbb{E}\|\tilde{g}^{\text{DP}}\|^2. \quad (17)$$

**Inner product expansion.** Substituting (15),

$$\mathbb{E}\langle \tilde{g}_k, \tilde{g}^{\text{DP}} \rangle = \sum_{j=1}^K \frac{n_j}{N} \mathbb{E}\langle \tilde{g}_k, \tilde{g}_j \rangle + \mathbb{E}\langle \tilde{g}_k, \eta \rangle. \quad (18)$$

Since  $\eta$  is zero-mean and independent of  $\tilde{g}_k$ ,

$$\mathbb{E}\langle \tilde{g}_k, \eta \rangle = 0. \quad (19)$$

For arbitrary random vectors  $A, B$ ,

$$\mathbb{E}[A^\top B] = \mathbb{E}[A]^\top \mathbb{E}[B] + \text{tr}(\text{Cov}(A, B)). \quad (20)$$

Thus,

$$\mathbb{E}\langle \tilde{g}_k, \tilde{g}_j \rangle = \langle g_k, g_j \rangle + \text{tr}(\text{Cov}(\tilde{g}_k, \tilde{g}_j)). \quad (21)$$

**Expected squared norms.** The local gradient satisfies

$$\mathbb{E}\|\tilde{g}_k\|^2 = \|g_k\|^2 + \text{tr}(\text{Var}(\tilde{g}_k)). \quad (22)$$

The global DP gradient yields

$$\begin{aligned} \mathbb{E}\|\tilde{g}^{\text{DP}}\|^2 &= \|g\|^2 + \text{tr}\left(\text{Var}\left(\sum_{j=1}^K \frac{n_j}{N} \tilde{g}_j\right)\right) + \text{tr}(\text{Var}(\eta)) \\ &= \|g\|^2 + \text{tr}\left(\sum_{i=1}^K \sum_{j=1}^K \frac{n_i n_j}{N^2} \text{Cov}(\tilde{g}_i, \tilde{g}_j)\right) \\ &\quad + \text{tr}(\text{Var}(\eta)). \end{aligned} \quad (23)$$

**Collecting terms.** The mean-gradient components simplify to

$$\|g_k\|^2 - 2\langle g_k, g \rangle + \|g\|^2 = \|g_k - g\|^2. \quad (24)$$

The variance–covariance terms combine to

$$\begin{aligned} \text{tr}(\text{Var}(\tilde{g}_k)) - 2 \sum_{j=1}^K \frac{n_j}{N} \text{tr}(\text{Cov}(\tilde{g}_k, \tilde{g}_j)) \\ + \sum_{i=1}^K \sum_{j=1}^K \frac{n_i n_j}{N^2} \text{tr}(\text{Cov}(\tilde{g}_i, \tilde{g}_j)). \end{aligned} \quad (25)$$

Adding  $\text{tr}(\text{Var}(\eta))$  yields the decomposition

$$\begin{aligned} \mathbb{E}\|\tilde{g}_k - \tilde{g}^{\text{DP}}\|^2 &= \|g_k - g\|^2 + \text{tr}(\text{Var}(\tilde{g}_k)) \\ &\quad - 2 \sum_{j=1}^K \frac{n_j}{N} \text{tr}(\text{Cov}(\tilde{g}_k, \tilde{g}_j)) \\ &\quad + \sum_{i=1}^K \sum_{j=1}^K \frac{n_i n_j}{N^2} \text{tr}(\text{Cov}(\tilde{g}_i, \tilde{g}_j)) \\ &\quad + \text{tr}(\text{Var}(\eta)), \end{aligned} \quad (26)$$

thereby proving the theorem.  $\square$

## B.2. Proof of Proposition 1

*Proof.* Let  $g_k = \mathbb{E}_{x \sim P_k}[\nabla \ell(w; x)]$  and  $g = \mathbb{E}_{x \sim P}[\nabla \ell(w; x)]$  denote the local and global expected gradients. We adopt a standard smoothness assumption in which the gradient with respect to the input can be locally approximated by a first-order expansion,

$$\nabla \ell(w; x) \approx H(w)(x - \mu_0) + b(w),$$

where  $H(w)$  approximates the Hessian and  $\mu_0$  is a reference point.

**Expected-gradient expansion.** Taking expectations under the local and global data distributions yields

$$g_k \approx H(w)(\mu_k - \mu_0) + b(w), \quad g \approx H(w)(\mu - \mu_0) + b(w),$$

where  $\mu_k$  and  $\mu$  denote the local and global means. Subtracting the two expressions gives

$$g_k - g \approx H(w)(\mu_k - \mu),$$

and therefore

$$\|g_k - g\|^2 \leq \|H(w)\|_F^2 \|\mu_k - \mu\|^2.$$

**Effect of DP noise on global mean estimation.** Each client reports a perturbed mean

$$\tilde{\mu}_k = \mu_k + \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, \sigma_\mu^2 I).$$

The server aggregates them to form

$$\tilde{\mu} = \sum_{j=1}^K \frac{n_j}{N} \tilde{\mu}_j = \mu + \epsilon, \quad \epsilon = \sum_{j=1}^K \frac{n_j}{N} \epsilon_j.$$

Since the noise across clients is independent and Gaussian,

$$\mathbb{E}\|\epsilon\|^2 = \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 \mathbb{E}\|\epsilon_j\|^2 = \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 d\sigma_\mu^2,$$

where  $d$  is the input dimension.

**Mean estimation error.** We decompose

$$\|\mu_k - \tilde{\mu}\|^2 = \|\mu_k - \mu - \epsilon\|^2 = \|\mu_k - \mu\|^2 + \|\epsilon\|^2 - 2\langle \mu_k - \mu, \epsilon \rangle.$$

The cross-term has zero expectation because  $\epsilon$  is independent of  $\mu_k - \mu$  and has zero mean. Thus,

$$\mathbb{E}\|\mu_k - \tilde{\mu}\|^2 = \|\mu_k - \mu\|^2 + \mathbb{E}\|\epsilon\|^2.$$

Rearranging yields

$$\|\mu_k - \mu\|^2 = \mathbb{E}\|\mu_k - \tilde{\mu}\|^2 - \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 d\sigma_\mu^2.$$

**Combining with gradientmean relationship.** Substituting into the gradient bound leads to

$$\|g_k - g\|^2 \leq \|H(w)\|_F^2 \left( \mathbb{E}\|\mu_k - \tilde{\mu}\|^2 - \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 d\sigma_\mu^2 \right),$$

which completes the proof.  $\square$

## B.3. Proof of Proposition 2

*Proof.* We first relate the variance of the stochastic gradient to input second moments via the Lipschitz property of the sample-wise gradient. Let  $\ell(\cdot)$  be the per-sample loss and assume its gradient with respect to model parameters is  $L$ -Lipschitz in the input, i.e.

$$\|\nabla_w \ell(w; x) - \nabla_w \ell(w; x')\| \leq L\|x - x'\|. \quad (27)$$

For the stochastic gradient  $\tilde{g}_k$  computed on client  $k$  we then have

$$\mathbb{E}[\|\tilde{g}_k - g_k\|^2] \leq L^2 \mathbb{E}[\|x\|^2], \quad (28)$$

where the expectation is taken over the local data sampling on client  $k$  and  $g_k = \mathbb{E}[\tilde{g}_k]$  denotes the corresponding mean gradient.

Decomposing the second moment of the input around the local mean  $\mu_k$  gives

$$\mathbb{E}[\|x\|^2] = \mathbb{E}[\|x - \mu_k\|^2] + \|\mu_k\|^2. \quad (29)$$

Using the identity  $\mathbb{E}[\|x - \mu_k\|^2] = \text{tr}(\Sigma_k)$ , where  $\Sigma_k$  is the covariance matrix of client  $k$ 's data, we obtain

$$\mathbb{E}[\|\tilde{g}_k - g_k\|^2] \leq L^2(\text{tr}(\Sigma_k) + \|\mu_k\|^2). \quad (30)$$

Since  $\text{tr}(\text{Var}(\tilde{g}_k)) = \mathbb{E}\|\tilde{g}_k - g_k\|^2$ , the above inequality yields the desired bound

$$\text{tr}(\text{Var}(\tilde{g}_k)) \leq L^2(\text{tr}(\Sigma_k) + \|\mu_k\|^2), \quad (31)$$

which is Eq. (9).

We now account for DP noise injected into the covariance estimation used for normalization/alignment. Let each client release a perturbed covariance

$$\tilde{\Sigma}_k = \Sigma_k + \xi_k, \quad \xi_k \sim \mathcal{N}(0, \sigma_{\text{DP},k}^2 I), \quad (32)$$

and let the server aggregate these perturbed covariances by weighted averaging:

$$\tilde{\Sigma} = \sum_{j=1}^K \frac{n_j}{N} \tilde{\Sigma}_j = \sum_{j=1}^K \frac{n_j}{N} \Sigma_j + \eta_{\text{agg}}, \quad (33)$$

where  $\eta_{\text{agg}} = \sum_{j=1}^K (n_j/N) \xi_j$  denotes the aggregated noise. Because the  $\xi_j$  are zero-mean, independent Gaussian perturbations, the expected squared Frobenius norm of the aggregated noise satisfies

$$\mathbb{E}\|\eta_{\text{agg}}\|_F^2 = \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 d_{\Sigma} \sigma_{\text{DP},j}^2, \quad (34)$$

where  $d_{\Sigma}$  is the dimension of the vectorized covariance matrix (i.e. the number of entries after flattening).

Replacing  $\Sigma_k$  by its noisy counterpart in the previous Lipschitz-based bound and taking expectation over the DP noise gives

$$\begin{aligned} \text{tr}(\text{Var}(\tilde{g}_k^{\text{DP}})) &\leq L^2 \left( \text{tr}(\Sigma_k) + \|\mu_k\|^2 + \mathbb{E}[\text{tr}(\eta_{\text{agg}})] \right) \\ &= L^2 \left( \text{tr}(\Sigma_k) + \|\mu_k\|^2 \right. \\ &\quad \left. + \sum_{j=1}^K \left(\frac{n_j}{N}\right)^2 d_{\Sigma} \sigma_{\text{DP},j}^2 \right), \end{aligned}$$

which is exactly the bound stated in Eq. (10). This completes the proof.  $\square$

## B.4. Proof of Proposition 3

*Proof.* Let  $g_k = \mathbb{E}_{x \sim P_k}[\nabla \ell(w; x)]$  and  $g = \sum_{j=1}^K (n_j/N) g_j$  denote the local and global gradients. Using the Edgeworth expansion, the expected gradient under distribution  $P_k$  admits the approximation

$$g_k \approx g^{\text{Gauss}} + \frac{\gamma_k}{6} M(w),$$

where  $M(w)$  is the third-order derivative term of the loss. Similarly,

$$g \approx g^{\text{Gauss}} + \frac{\bar{\gamma}}{6} M(w).$$

Subtracting gives the skewness-induced deviation

$$g_k - g \approx \frac{\gamma_k - \bar{\gamma}}{6} M(w),$$

and therefore

$$\|g_k - g\|^2 \propto |\gamma_k - \bar{\gamma}|^2.$$

**Effect of DP noise.** Let the DP-perturbed gradients be

$$g_k^{\text{DP}} = g_k + \xi_k, \quad g^{\text{DP}} = g + \xi,$$

where  $\xi_k, \xi$  are zero-mean Gaussian and independent of skewness effects. Then

$$g_k^{\text{DP}} - g^{\text{DP}} = (g_k - g) + (\xi_k - \xi).$$

The cross term vanishes in expectation, so the deterministic skewness contribution remains dominant:

$$\mathbb{E}\|g_k^{\text{DP}} - g^{\text{DP}}\|^2 = \|g_k - g\|^2 + \mathcal{O}(\sigma_{\text{DP}}^2) \propto |\gamma_k - \bar{\gamma}|^2.$$

This completes the proof.  $\square$

## B.5. Proof of Proposition 4

*Proof.* Let  $x \sim P_k$  with mean  $\mu_k$ , covariance  $\Sigma_k$ , and kurtosis  $\kappa_k$ . For smooth loss  $\ell(w; x)$ , linearizing the stochastic gradient around  $\mu_k$  gives

$$\nabla \ell(w; x) \approx A(w)(x - \mu_k) + b(w),$$

so the Gaussian reference variance satisfies

$$\text{Var}^{\text{Gauss}}(\nabla \ell(w; x)) \propto \text{tr}(\Sigma_k).$$

**Kurtosis correction.** Applying GramCharlier expansion and retaining the leading non-Gaussian term gives the variance adjustment

$$\text{Var}(\nabla \ell(w; x)) \approx \text{Var}^{\text{Gauss}}(\nabla \ell(w; x)) \left( 1 + \frac{\kappa_k - 3}{4} \right).$$

Taking traces yields

$$\text{tr}(\text{Var}(\nabla \ell(w; x))) \propto \left( 1 + \frac{\kappa_k - 3}{4} \right) \text{tr}(\Sigma_k).$$

**Effect of DP noise.** With DP-perturbed gradient

$$\tilde{g}_k^{\text{DP}} = g_k + \epsilon_{\text{DP}},$$

the additive isotropic noise contributes independent variance

$$\text{tr}(\text{Var}(\epsilon_{\text{DP}})) = d\sigma_{\text{DP}}^2.$$

Combining the kurtosis-induced multiplicative factor with the DP noise term gives the desired relationship.  $\square$

## B.6. Proof of Theorem 2

*Proof.* The proof follows the standard descent analysis under  $L$ -smoothness, combined with a decomposition of the stochastic update direction into its bias and variance components, and an explicit upper bound on the statistical disparity term.

**Descent inequality under  $L$ -smoothness.** For any  $L$ -smooth non-convex objective  $F$ , the update  $w_{t+1} = w_t - \eta d_t$  satisfies

$$F(w_{t+1}) \leq F(w_t) + \langle \nabla F(w_t), w_{t+1} - w_t \rangle + \frac{L}{2} \|w_{t+1} - w_t\|^2. \quad (35)$$

Substituting the update rule gives

$$F(w_{t+1}) \leq F(w_t) - \eta \langle \nabla F(w_t), d_t \rangle + \frac{L\eta^2}{2} \|d_t\|^2. \quad (36)$$

**Expectation and bias–variance decomposition.** Taking expectation of (36) and expanding  $d_t = \mathbb{E}[d_t] + (d_t - \mathbb{E}[d_t])$  yields

$$\begin{aligned} \mathbb{E}[F(w_{t+1})] &\leq \mathbb{E}[F(w_t)] - \eta \langle \nabla F(w_t), \mathbb{E}[d_t] \rangle \\ &\quad + \frac{L\eta^2}{2} \left( \|\mathbb{E}[d_t]\|^2 + \mathbb{E}\|d_t - \mathbb{E}[d_t]\|^2 \right). \end{aligned} \quad (37)$$

Introduce and subtract  $\eta \|\nabla F(w_t)\|^2$  to reveal a descent term:

$$\begin{aligned} \mathbb{E}[F(w_{t+1})] &\leq \mathbb{E}[F(w_t)] - \eta \|\nabla F(w_t)\|^2 \\ &\quad + \eta \langle \nabla F(w_t), \nabla F(w_t) - \mathbb{E}[d_t] \rangle \\ &\quad + \frac{L\eta^2}{2} \left( \|\mathbb{E}[d_t]\|^2 + \mathbb{E}\|d_t - \mathbb{E}[d_t]\|^2 \right). \end{aligned} \quad (38)$$

Applying the Cauchy–Schwarz and Young inequalities gives

$$\begin{aligned} \langle \nabla F(w_t), \nabla F(w_t) - \mathbb{E}[d_t] \rangle &\leq \frac{1}{2} \|\nabla F(w_t)\|^2 \\ &\quad + \frac{1}{2} \|\nabla F(w_t) - \mathbb{E}[d_t]\|^2. \end{aligned} \quad (39)$$

Combining terms yields the per-iteration inequality

$$\begin{aligned} \mathbb{E}[F(w_{t+1})] &\leq \mathbb{E}[F(w_t)] - \frac{\eta}{2} \|\nabla F(w_t)\|^2 \\ &\quad + \frac{\eta}{2} \|\nabla F(w_t) - \mathbb{E}[d_t]\|^2 \\ &\quad + \frac{L\eta^2}{2} \left( \|\mathbb{E}[d_t]\|^2 + \mathbb{E}\|d_t - \mathbb{E}[d_t]\|^2 \right). \end{aligned} \quad (40)$$

**Bounding the bias term via statistical disparity.** Using Jensen’s inequality,

$$\|\nabla F(w_t) - \mathbb{E}[d_t]\|^2 \leq \frac{1}{K} \sum_{k=1}^K \|\nabla F(w_t) - g_k(w_t)\|^2. \quad (41)$$

The deviation  $\|\nabla F - g_k\|$  can be bounded by discrepancy of the empirical statistics of client  $k$ , giving

$$\begin{aligned} \|\nabla F(w_t) - \mathbb{E}[d_t]\|^2 &\leq \frac{1}{K} \sum_{k=1}^K \left( L_\mu^2 \Delta_{\mu,k} + L_\sigma^2 \|\Sigma_k - \Sigma\|_F^2 \right. \\ &\quad \left. + L_\gamma |\gamma_k - \gamma| + L_\kappa |\kappa_k - \kappa| \right) \\ &= \Gamma_{\text{stat}}. \end{aligned} \quad (42)$$

**Bounding the variance term and DP noise.** The variance of  $d_t$  decomposes as

$$\mathbb{E}\|d_t - \mathbb{E}[d_t]\|^2 = \mathbb{E}\|\tilde{g}_t - \mathbb{E}[\tilde{g}_t]\|^2 + \mathbb{E}\|\varepsilon_t\|^2, \quad (43)$$

where  $\tilde{g}_t$  is the stochastic mini-batch gradient and  $\varepsilon_t$  is DP Gaussian noise.

Define the minimal client-side variance by  $\sigma_{\min}^2$ , and note that DP noise satisfies

$$\mathbb{E}\|\varepsilon_t\|^2 = \sigma_{\text{DP}}^2 = \frac{8C^2 \log(1.25/\delta)}{\epsilon^2}. \quad (44)$$

Thus

$$\mathbb{E}\|d_t - \mathbb{E}[d_t]\|^2 \leq \sigma_{\min}^2 + \sigma_{\text{DP}}^2. \quad (45)$$

**Final convergence rate.** Summing (40) over  $t = 0, \dots, T-1$  and telescoping the left-hand side gives

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla F(w_t)\|^2 \leq \frac{2\Delta_F}{\eta T} + \eta L (\sigma_{\min}^2 + \Gamma_{\text{stat}} + \sigma_{\text{DP}}^2), \quad (46)$$

which matches the claim in Theorem 2.  $\square$