

## Supplementary Material: Bayesian Coreset Optimization for Personalized Federated Learning

In this supplementary material we discuss extensively on the proofs. involved for the theoretical analysis for CORESET-PFEDBAYES along with more fine-grained experimental details and corresponding baselines.

### 9 PROOFS

Here we discuss the proofs involved with particular propositions and theorems specified in the Theoretical Contributions of this paper. Utilising the assumptions taken in [Zhang et al. \(2022b\)](#), [Polson & Ročková \(2018\)](#) we consider the analysis for equal-width Bayesian Neural network.

**Assumption 1:** The widths of the neural network are equal width i.e.  $s_i = M$ .

**Assumption 2:** Each individual client  $i \in [N]$  has equal coreset size of samples  $n_k < n$ .

**Assumption 3:** Parameters  $s_0, n$  (total client dataset size)  $n_k$  (coreset client dataset size),  $M, L$  (number of DNN layers as per Section [3](#)) are large enough such that the sequence  $\sigma_n^2$  is bounded as follows

$$\sigma_n^2 = \frac{T}{8n} A \leq \Omega^2,$$

where  $\tau = \Omega M$  and

$$A = \log^{-1}(3s_0M) \cdot (2\tau)^{-2(L+1)} \left[ \left( s_0 + 1 + \frac{1}{\tau - 1} \right)^2 + \frac{1}{(2\tau)^2 - 1} + \frac{2}{(2\tau - 1)^2} \right]^{-1}.$$

Here  $T$  indicates the total number of parameters as defined in Section [3](#)

Similarly, utilising the coreset regime, we have the following:

$$\sigma_{n_k}^2 = \frac{T}{8n_k} A \leq \Omega^2,$$

Since  $n_k \ll n$ , hence  $\sigma_{n_k}^2 \gg \sigma_n^2$

**Assumption 4:** We consider 1-Lipschitz continuous activation function  $\sigma(\bullet)$

We also define here a few terms as defined in [\(Zhang et al., 2022b\)](#) which would be useful for our following proof proposals as well.

**Definition 2.** Preliminaries and Definitions required for theoretical proofs under PFEDBAYES

$$\begin{aligned} d^2(\mathcal{P}_{\theta}^i, \mathcal{P}^i) &= \mathbb{E}_{X^i} \left( 1 - e^{-\frac{[f_{\theta}^i(X^i) - f^i(X^i)]^2}{8\sigma_n^2}} \right) \\ r_n &= ((L+1)T/n) \log M + (T/n) \log \left( s_0 \sqrt{n/T} \right) \\ \xi_n^i &= \inf_{\theta \in \Theta(L, \mathcal{S}), \|\theta\|_{\infty} \leq \Omega} \|f_{\theta}^i - f^i\|_{\infty}^2, \\ \varepsilon_n &= n^{-\frac{1}{2}} \sqrt{(L+1)T \log M + T \log \left( s_0 \sqrt{n/T} \right)} \log^{\delta}(n) = \sqrt{r_n} \log^{\delta}(n), \end{aligned}$$

where  $\delta > 1$

Here  $r_n$  indicates the variational error incurred due to the Bayesian approximation to the true posterior distribution in Equation [1](#) and  $\xi_n^i$  indicates the approximation error incurred during regression w.r.t the actual function to be learnt.

Similarly for the coreset size  $n_k$  we define the following:

**Definition 3.** *Preliminaries and Definitions required for theoretical proofs under CORESET-PFEDBAYES*

$$\begin{aligned} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) &= \mathbb{E}_{X^i} \left( 1 - e^{-\frac{[f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)]^2}{8\sigma_\epsilon^2}} \right) \\ \xi_{n_k}^i &= \inf_{\theta \in \Theta(L, \mathcal{S}), \|\theta\|_\infty \leq \Omega} \|f_{\theta, \mathbf{w}}^i - f^i\|_\infty^2 \\ r_{n_k} &= ((L+1)T/n_k) \log M + (T/n_k) \log \left( s_0 \sqrt{n_k/T} \right) \\ \varepsilon_{n_k} &= n_k^{-\frac{1}{2}} \sqrt{(L+1)T \log M + T \log \left( s_0 \sqrt{n_k/T} \right)} \log^\delta(n_k) = \sqrt{r_{n_k}} \log^\delta(n_k) \end{aligned}$$

**Lemma 1.** *The Hellinger Distance from Definition [1](#) is symmetrical in its arguments  $\mathcal{P}_\theta^i$  and  $\mathcal{P}^i$ .*

*Proof.* It is easy to show that,

$$d^2(\mathcal{P}_\theta^i, \mathcal{P}^i) = \mathbb{E}_{X^i} \left( 1 - e^{-\frac{[f_\theta^i(X^i) - f^i(X^i)]^2}{8\sigma_\epsilon^2}} \right) \quad (10)$$

$$= \mathbb{E}_{X^i} \left( 1 - e^{-\frac{[f^i(X^i) - f_\theta^i(X^i)]^2}{8\sigma_\epsilon^2}} \right) \quad (11)$$

$$= d^2(\mathcal{P}^i, \mathcal{P}_\theta^i) \quad (12)$$

□

### PROOF. OF THEOREM [1](#)

**Theorem 1.** *The difference in the upper bound incurred in the overall generalization error of CORESET-PFEDBAYES as compared w.r.t that of PFEDBAYES is always upper bounded by a closed form positive function that depends on the coreset weights and coreset size-  $\mathfrak{S}(\mathbf{w}, n_k)$ . generalization error in the original full data setup*

$$\left[ \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_\theta^i, \mathcal{P}^i) \hat{q}^i(\theta) d\theta \right]_{u.b.} - \left[ \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta \right]_{u.b.} \leq \mathfrak{S}(\mathbf{w}, n_k)$$

*Proof.* Let us define  $\log \eta(\mathcal{P}_\theta^i, \mathcal{P}^i) = l_n(\mathcal{P}_\theta^i, \mathcal{P}^i) / \zeta + n d^2(\mathcal{P}_\theta^i, \mathcal{P}^i)$ .

Using Theorem 3.1 of [Pati et al. \(2018\)](#) with probability at most  $e^{-C n_k \varepsilon_{n_k}^2}$ , where  $C$  is a constant, with high probability for CORESET-PFEDBAYES we have

$$\int_{\Theta} \eta(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) z^*(\theta) d\theta \leq e^{C n_k \varepsilon_{n_k}^2} \quad (13)$$

Similarly with high probability at most  $e^{-C n \varepsilon_n^2}$  for the vanilla PFEDBAYES

$$\int_{\Theta} \eta(\mathcal{P}_\theta^i, \mathcal{P}^i) z^*(\theta) d\theta \leq e^{C n \varepsilon_n^2} \quad (14)$$

Using Lemma A.1 from [Zhang et al. \(2022b\)](#) we know that for any probability measure  $\mu$  and any measurable function  $h$  with  $e^h \in L_1(\mu)$ ,

$$\log \int e^{h(\eta)} \mu(d\eta) = \sup_{\rho} \left[ \int h(\eta) \rho(d\eta) - \mathbb{D}_{KL}(\rho \| \mu) \right]$$

Further, we let  $l_n(P^i, P_\theta^i)$  is the log-likelihood ratio of  $P^i$  and  $P_\theta^i$

$$l_n(P^i, P_\theta^i) = \log \frac{\mathcal{P}^i(\mathbf{D}^i)}{\mathcal{P}_\theta^i(\mathbf{D}^i)}.$$

Hence,

$$\begin{aligned} nd^2(P_\theta^i, P^i) &= l_n(P_\theta^i, P^i)/\zeta - \log \eta(P_\theta^i, P^i) \\ &= l_n(P^i, P_\theta^i)/\zeta - \log \eta(P^i, P_\theta^i) \quad \text{since } d^2(P_\theta^i, P^i) = d^2(P^i, P_\theta^i) \text{ from Lemma 1} \end{aligned}$$

This follows from 9

Similarly, for the weighted likelihood based Hellinger Distance,

$$n_k d^2(P_{\theta, \mathbf{w}}^i, P^i) = l_n(P_{\theta, \mathbf{w}}^i, P^i)/\zeta - \log \eta(P_{\theta, \mathbf{w}}^i, P^i) \quad (15)$$

By using Lemma A.1 with  $h(\eta) = \log \eta(P_\theta^i, P^i)$ ,  $\mu = z^*(\theta)$  and  $\rho = \hat{q}^i(\theta)$ , we obtain

$$\begin{aligned} \int_{\Theta} d^2(P_\theta^i, P^i) \hat{q}^i(\theta) d\theta &\leq \frac{1}{n} \left[ \frac{1}{\zeta} \int_{\Theta} l_n(P^i, P_\theta^i) \hat{q}^i(\theta) d\theta + \mathbb{D}_{KL}(\hat{q}^i(\theta) \| z^*(\theta)) + \log \int_{\Theta} \eta(P_\theta^i, P^i) z^*(\theta) d\theta \right] \\ &\leq \frac{1}{n} \left[ \frac{1}{\zeta} \int_{\Theta} l_n(P^i, P_\theta^i) \hat{q}^i(\theta) d\theta + \mathbb{D}_{KL}(\hat{q}^i(\theta) \| z^*(\theta)) \right] + C\varepsilon_n^2 \end{aligned}$$

$$\int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) q^i(\hat{\theta}; \mathbf{w}) d\theta \leq \frac{1}{n_k} \left[ \frac{1}{\zeta} \int_{\Theta} l_n(P^i, P_{\theta, \mathbf{w}}^i) q^i(\hat{\theta}; \mathbf{w}) d\theta + \mathbb{D}_{KL}(q^i(\hat{\theta}; \mathbf{w}) \| z^*(\theta)) \right] + C\varepsilon_{n_k}^2$$

Utilising analysis under Supplementary in (Bai et al., 2020), there exists an upper bound for the term

$$\int_{\Theta} l_n(P^i, P_\theta^i) \hat{q}^i(\theta) d\theta \leq C''(nr_n + n\xi_n^i) \quad (16)$$

Lemma 2 from (Zhang et al., 2022b) provides the upper bound for the KL divergence term

$$\mathbb{D}_{KL}(\hat{q}^i(\theta) \| z^*(\theta)) \leq C'(nr_n) \quad (17)$$

Therefore we can write the following expression that captures the weighted Hellinger distance displacement given in our coreset framework CORESET-PFEDBAYES as compared to PFEDBAYES

$$\begin{aligned}
& \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta}^i, \mathcal{P}^i) q^i(\hat{\theta}) d\theta - \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) q^i(\hat{\theta}, \mathbf{w}) d\theta \\
& \leq \frac{1}{N} \sum_{i=1}^N \frac{1}{n} \left[ \frac{1}{\zeta} \int_{\Theta} l_n(P^i, P_{\theta}^i) \hat{q}^i(\theta) d\theta + \mathbb{D}_{KL}(\hat{q}^i(\theta) \| z^*(\theta)) \right] + C\varepsilon_n^2 - \\
& \frac{1}{N} \sum_{i=1}^N \frac{1}{n_k} \left[ \frac{1}{\zeta} \int_{\Theta} l_n(P^i, P_{\theta, \mathbf{w}}^i) q^i(\hat{\theta}; \mathbf{w}) d\theta + \mathbb{D}_{KL}(q^i(\hat{\theta}; \mathbf{w}) \| z^*(\theta)) \right] - C\varepsilon_{n_k}^2 \\
& \hspace{15em} \text{Using Eq:(17) and Eq:(16)} \\
& \leq C\varepsilon_n^2 - C\varepsilon_{n_k}^2 + n \left( C' \zeta r_n + \frac{C''}{N} \sum_{i=1}^N \xi_n^i \right) - n_k \left( C' \zeta r_{n_k} + \frac{C''}{N} \sum_{i=1}^N \xi_{n_k}^i \right) \\
& \leq C(\varepsilon_n^2 - \varepsilon_{n_k}^2) + \zeta C' (nr_n - n_k r_{n_k}) + \frac{C''}{N} \sum_{i=1}^N (n\xi_n^i - n_k \xi_{n_k}^i) \\
& = \underbrace{C(\varepsilon_n^2 - \varepsilon_{n_k}^2)}_{\text{Estimation error Type I Drift}} + \underbrace{\zeta C' (nr_n - n_k r_{n_k})}_{\text{Estimation error Type II Drift}} + \underbrace{\frac{C''}{N} \sum_{i=1}^N (n\xi_n^i - n_k \xi_{n_k}^i)}_{\text{Approximation Error Drift}} \\
& \hspace{15em} = \underbrace{\mathfrak{S}(\mathbf{w}, n_k)}_{\geq 0}
\end{aligned}$$

Where  $\mathfrak{S}(\mathbf{w}, n_k) = C(\varepsilon_n^2 - \varepsilon_{n_k}^2) + \zeta C' (nr_n - n_k r_{n_k}) + \frac{C''}{N} \sum_{i=1}^N (n\xi_n^i - n_k \xi_{n_k}^i)$  where each of the coefficients of the closed form function are constants related to  $s_0, \beta, \mathbf{A}, L, M, \zeta$  and  $n_k$

Using Lemma 2, 3 and with suitable assumptions on the Approximation drift error such that we see that each of the individual error terms are positive, there by indicating  $\mathfrak{S}(\mathbf{w}, n_k) \geq 0$   $\square$

**Lemma 2.** *The Estimation error Type II Drift is a positive quantity i.e.  $nr_n > n_k r_{n_k}$ .*

*Proof.* By Definition,

$$r_n = ((L+1)T/n) \log M + (T/n) \log \left( s_0 \sqrt{n/T} \right)$$

and

$$r_{n_k} = ((L+1)T/n_k) \log M + (T/n_k) \log \left( s_0 \sqrt{n_k/T} \right)$$

Hence

$$\frac{r_n}{r_{n_k}} = \frac{n_k}{n} \times \frac{((L+1)T) \log M + (T) \log \left( s_0 \sqrt{n/T} \right)}{((L+1)T) \log M + (T) \log \left( s_0 \sqrt{n_k/T} \right)}$$

$$\frac{r_n}{r_{n_k}} = \frac{n_k}{n} \times \frac{(L+1) \log M + \log \left( s_0 / \sqrt{T} \right) + \log(\sqrt{n})}{(L+1) \log M + \log \left( s_0 / \sqrt{T} \right) + \log(\sqrt{n_k})}$$

Considering  $(L+1) \log M + \log \left( s_0 / \sqrt{T} \right)$  as a constant  $\mathfrak{G}$  we have

$$\frac{r_n}{r_{n_k}} = \frac{n_k}{n} \times \frac{\mathfrak{G} + \log(\sqrt{n})}{\mathfrak{G} + \log(\sqrt{n_k})}$$

Thus,

$$\frac{nr_n}{n_k r_{n_k}} = \frac{\mathfrak{G} + \log(\sqrt{n})}{\mathfrak{G} + \log(\sqrt{n_k})}$$

It is clear since  $\log(\bullet)$  is an increasing function for  $n > n_k$  we have  $nr_n > n_k r_{n_k}$ .  $\square$

**Lemma 3.** *The Estimation error Type I Drift is a positive quantity i.e.  $\varepsilon_n^2 > \varepsilon_{n_k}^2$ .*

*Proof.* From the definition under Assumption 3

$$\varepsilon_{n_k} = n_k^{-\frac{1}{2}} \sqrt{(L+1)T \log M + T \log(s_0 \sqrt{n_k/T})} \log^\delta(n_k) = \sqrt{r_{n_k}} \log^\delta(n_k)$$

Hence  $\varepsilon_{n_k}^2 = r_{n_k} \log^{2\delta}(n_k)$ . Similarly,  $\varepsilon_n^2 = r_n \log^{2\delta}(n)$

$$\frac{\varepsilon_n^2}{\varepsilon_{n_k}^2} = \frac{r_n \log^{2\delta}(n)}{r_{n_k} \log^{2\delta}(n_k)} = \frac{nr_n \frac{\log^{2\delta}(n)}{n}}{n_k r_{n_k} \frac{\log^{2\delta}(n_k)}{n_k}}$$

From Lemma 2 we know that  $nr_n > n_k r_{n_k}$ , hence

$$\frac{\varepsilon_n^2}{\varepsilon_{n_k}^2} > \frac{\frac{\log^{2\delta}(n)}{n}}{\frac{\log^{2\delta}(n_k)}{n_k}} > 1$$

This follows due to the increasing nature of the function.  $\square$

## PROOF OF THEOREM 2

**Theorem 2.** *The convergence rate of the generalization error under  $L^2$  norm of CORESET-PFEDBAYES is minimax optimal up to a logarithmic term (in order  $n_k$ ) for bounded functions ( $\beta$ -Hölder-smooth functions)  $\{f^i\}_{i=1}^N$ ,  $\{f_\theta^i\}_{i=1}^N$  and  $\{f_{\theta,w}^i\}_{i=1}^N$  where  $C_2$ ,  $C_3$  and  $\delta'$  are constants and  $\Lambda$  being the intrinsic dimension of each client's data:*

$$\frac{C_F}{N} \sum_{i=1}^N \int_{\theta} \|f_{\theta,w}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta; \mathbf{w}) d\theta \leq C_2 n_k^{-\frac{2\beta}{2\beta+\Lambda}} \log^{2\delta'}(n_k).$$

and

$$\left\{ \|f_{\theta,w}^i\|_{\infty} \leq F \right\}_{i=1}^N \inf_{\{ \|f^i\|_{\infty} \leq F \}_{i=1}^N} \frac{C_F}{N} \sum_{i=1}^N \int_{\theta} \|f_{\theta,w}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta; \mathbf{w}) d\theta \geq C_3 n_k^{-\frac{2\beta}{2\beta+\Lambda}}$$

where  $n_k$  denotes the coreset size per client dataset and  $n$  denotes the original per client dataset

size and  $\frac{d^2(P_{\theta,w}^i, P^i)}{\|f_{\theta,w}^i(X^i) - f^i(X^i)\|_{L^2}^2} \geq \frac{1 - \exp\left(-\frac{4F^2}{8\sigma_\varepsilon^2}\right)}{4F^2} \triangleq C_F$ .

We present the choice of  $T$  for a typical class of functions. We already assumed that  $\{f^i\}$  are  $\beta$ -Hölder-smooth functions (Definition 4. (Nakada & Imaizumi, 2020)) and the intrinsic dimension of data is  $\Lambda$ .

From our above theorem result from Theorem: 1 we say the following:

$$\frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta \leq C \varepsilon_{n_k}^2 + C' r_{n_k} + \frac{C''}{N\zeta} \sum_{i=1}^N \xi_{n_k}^i \quad (18)$$

Utilising Corollary 6 in (Nakada & Imaizumi, 2020), the approximation error is upper-bounded as follows

$$\|f_{\theta, \mathbf{w}}^i - f^i\|_{\infty} \leq C_0 T^{-\frac{\beta}{\Lambda}}$$

where  $C_0 > 0$  is a constant related to  $s_0, \beta$  and  $\Lambda$

Thus from the above definitions 2 and 3, we have the following

$$\xi_n^i, \xi_{n_k}^i \leq C_0 T^{-\frac{2\beta}{\Lambda}}, i = 1, \dots, N$$

Utilising the above upper bound in 18 and substituting  $T = C_1 n^{\frac{\Lambda}{2\beta+\Lambda}}$ , we get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta &\leq C \varepsilon_{n_k}^2 + C' r_{n_k} + \frac{C''}{N\zeta} \sum_{i=1}^N C_0 T^{-\frac{2\beta}{\Lambda}} \\ &\leq C r_{n_k} \log^{2\delta}(n_k) + C' r_{n_k} + \frac{C''}{N\zeta} \sum_{i=1}^N C_0 T^{-\frac{2\beta}{\Lambda}} \because \varepsilon_{n_k}^2 = r_{n_k} \log^{2\delta}(n_k) \\ &\leq C_2 n_k^{-\frac{2\beta}{2\beta+\Lambda}} \log^{2\delta'}(n_k) \text{ [ substituting } T \text{ in } r_{n_k} \text{]} \end{aligned}$$

where  $\delta' > \delta > 1$ , and  $C_1, C_2 > 0$  are constants related to  $s_0, \beta, \Lambda, L, M, \zeta$  and  $n_k$ .

Similar to Theorem 1.1 from (Bai et al., 2020) and Theorem 1 from (Zhang et al., 2022b) norm, we can write the following

$$\begin{aligned} \frac{C_F}{N} \sum_{i=1}^N \int_{\Theta} \|f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta, \mathbf{w}) d\theta \\ \leq \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta \\ \leq C_2 n_k^{-\frac{2\beta}{2\beta+\Lambda}} \log^{2\delta'}(n_k). \end{aligned}$$

Now, using the minimax lower bound under  $L^2$  norm in Theorem 8 of (Nakada & Imaizumi, 2020), we see that for coresets regime the same formulation holds similar to our original setting as shown in (Zhang et al., 2022b)

$$\inf_{\{ \|f_{\theta, \mathbf{w}}^i\|_{\infty} \leq F \}_{i=1}^N} \inf_{\{ \|f^i\|_{\infty} \leq F \}_{i=1}^N} \frac{C_F}{N} \sum_{i=1}^N \int_{\Theta} \|f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta; \mathbf{w}) d\theta \geq C_3 n_k^{-\frac{2\beta}{2\beta+\Lambda}}$$

where  $C_3 > 0$  is a constant.

Combining the above two equations, the convergence rate of the generalization error of the coresets weighted objective is minimax optimal upto a logarithmic term for bounded functions  $\{f_{\theta, \mathbf{w}}^i\}_{i=1}^N$  and  $\{f^i\}_{i=1}^N$ .

**PROOF. OF THEOREM 3**

**Theorem 3.** *The lower bound (l.b.) incurred for the deviation for the weighted coreset CORESET-PFEDBAYES (5) generalization error is always higher than the lower bound of that for the original PFEDBAYES objective (1) with a delta difference (Error I - Error II) as  $\mathcal{O}(n_k^{-\frac{2\beta}{2\beta+\Lambda}})$*

$$\underbrace{\left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta, \mathbf{w}) d\theta \right]}_{\text{Coreset weighted objective Generalization Error (Error I)}} \Big|_{l.b.} > \underbrace{\left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta) d\theta \right]}_{\text{Vanilla objective Generalization Error (Error II)}} \Big|_{l.b.}$$

*Proof.* As we know  $n_k < n$  hence  $C_3 n_k^{-\frac{2\beta}{2\beta+\Lambda}} > C_3 n^{-\frac{2\beta}{2\beta+\Lambda}}$  ( $\because C_3$  is a constant independent of  $n$  or  $n_k$ ), which therefore means that inequality holds in the lower bound (l.b.) of the two expressions (shown by the previous proposition 2).

$$\left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta, \mathbf{w}) d\theta \right] \Big|_{l.b.} > \left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta) d\theta \right] \Big|_{l.b.}$$

Let us denote  $\Delta_{deviation}^{l.b.}$  as follows

$$\Delta_{deviation}^{l.b.} = \left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta, \mathbf{w}}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta, \mathbf{w}) d\theta \right] \Big|_{l.b.} - \left[ \sum_{i=1}^N \int_{\Theta} \|f_{\theta}^i(X^i) - f^i(X^i)\|_{L^2}^2 \hat{q}^i(\theta) d\theta \right] \Big|_{l.b.}$$

And the  $\Delta_{deviation}^{l.b.}$  term is given by  $\left( C_3 n_k^{-\frac{2\beta}{2\beta+\Lambda}} - C_3 n^{-\frac{2\beta}{2\beta+\Lambda}} \right) \approx \mathcal{O}(n_k^{-\frac{2\beta}{2\beta+\Lambda}})$ .  $\square$

**PROOF. OF THEOREM 4**

**Theorem 4.** *The lower bound incurred in the overall generalization error across all  $N$  clients of CORESET-PFEDBAYES is always higher compared to that of the generalization error in the original full data setup*

$$\left[ \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta \right] \Big|_{l.b.} \geq \left[ \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta}^i, \mathcal{P}^i) \hat{q}^i(\theta) d\theta \right] \Big|_{l.b.}$$

*Proof.* It is easy to show since from Theorem 3, we know the lower bounds for the individual terms and also since  $n > n_k$  holds, hence we can rewrite as follows:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta, \mathbf{w}}^i, \mathcal{P}^i) \hat{q}^i(\theta; \mathbf{w}) d\theta &- \frac{1}{N} \sum_{i=1}^N \int_{\Theta} d^2(\mathcal{P}_{\theta}^i, \mathcal{P}^i) \hat{q}^i(\theta) d\theta \\ &\geq C_3 n_k^{-\frac{2\beta}{2\beta+\Lambda}} - C_3 n^{-\frac{2\beta}{2\beta+\Lambda}} \\ &\geq 0 \end{aligned}$$

The implication of this proof states that the overall error incurred due to coreset weighted deviation is always more than that of the original deviation which can be measured approximately in order of  $n_k$ , the coreset sample size.  $\square$

**Proposition 1.** *The gradient of the first term in Equation 7 i.e.*

$$\nabla_{\mathbf{w}} \mathbb{D}_{KL}(\hat{q}^i(\theta; \mathbf{w}) \| \hat{q}^i(\theta))$$

*is given by the following expression*

$$\int_{\Theta} \nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \left[ \log \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) + 1 - \log \hat{q}^i(\boldsymbol{\theta}) \right] d\boldsymbol{\theta}$$

where

$$\nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) = \frac{\hat{q}^i(\boldsymbol{\theta}; \mathbf{w})}{\varrho^i(\boldsymbol{\theta}_{i,m}; \mathbf{w})} g'_m(\mathbf{w}) + g_m(\mathbf{w}) \nabla_w \prod_{k \neq m}^T \varrho^i(\boldsymbol{\theta}_{i,k}; \mathbf{w}) \quad \text{and} \quad q^i(\boldsymbol{\theta}; \mathbf{w}) = \prod_{m=1}^T \varrho^i(\boldsymbol{\theta}_{i,m}; \mathbf{w})$$

*Proof.*

$$\begin{aligned} & \nabla_w \mathbb{D}_{KL}(\hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \| \hat{q}^i(\boldsymbol{\theta})) \\ &= \nabla_w \mathbb{E}_{\hat{q}^i(\boldsymbol{\theta}; \mathbf{w})} \left[ \log \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) - \log \hat{q}^i(\boldsymbol{\theta}) \right] \\ &= \nabla_w \left[ \int_{\Theta} \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \log \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) d\boldsymbol{\theta} - \int_{\Theta} \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \log \hat{q}^i(\boldsymbol{\theta}) d\boldsymbol{\theta} \right] \\ &= \left[ \int_{\Theta} \nabla_w \left( \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \log \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \right) d\boldsymbol{\theta} - \int_{\Theta} \nabla_w \left( \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \log \hat{q}^i(\boldsymbol{\theta}) \right) d\boldsymbol{\theta} \right] \\ &= \left[ \int_{\Theta} \left( \log q^i(\boldsymbol{\theta}; \mathbf{w}) \nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) + \nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \right) d\boldsymbol{\theta} - \int_{\Theta} \log \hat{q}^i(\boldsymbol{\theta}) \nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) d\boldsymbol{\theta} \right] \\ &= \int_{\Theta} \nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) \left[ \log \hat{q}^i(\boldsymbol{\theta}; \mathbf{w}) + 1 - \log \hat{q}^i(\boldsymbol{\theta}) \right] d\boldsymbol{\theta} \end{aligned} \quad (19)$$

In order to compute the gradient  $\nabla_w \hat{q}^i(\boldsymbol{\theta}; \mathbf{w})$ , the following objective can be utilized.

Let  $z^*(\boldsymbol{\theta})$  be the optimal variable solution to Equation (5).

$$\begin{aligned} & \nabla_{q^i(\boldsymbol{\theta})} F_i^w(z^*) \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} = 0 \\ \implies & \underbrace{\nabla_{q^i(\boldsymbol{\theta})} \int_{\Theta} -\log \mathcal{P}_{\theta, w}(\mathcal{D}^i) q^i(\boldsymbol{\theta}) d\boldsymbol{\theta}}_{\text{First Part}} \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} + \underbrace{\zeta \nabla_{q^i(\boldsymbol{\theta})} \mathbb{D}_{KL}(q^i(\boldsymbol{\theta}) \| z^*(\boldsymbol{\theta}))}_{\text{Second Part}} \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} = 0 \end{aligned}$$

For the **first part**,

$$\begin{aligned} & \nabla_{q^i(\boldsymbol{\theta})} \int_{\Theta} -\log \mathcal{P}_{\theta, w}(\mathcal{D}^i) q^i(\boldsymbol{\theta}) d\boldsymbol{\theta} \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} \\ &= \int_{\Theta} \nabla_{q^i(\boldsymbol{\theta})} \left[ -\log \mathcal{P}_{\theta, w}(\mathcal{D}^i) q^i(\boldsymbol{\theta}) \right] d\boldsymbol{\theta} \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} \\ &= \int_{\Theta} \underbrace{\left[ -q^i(\boldsymbol{\theta}) \nabla_{q^i(\boldsymbol{\theta})} \log \mathcal{P}_{\theta, w}(\mathcal{D}^i) + \log \mathcal{P}_{\theta, w}(\mathcal{D}^i) \right]}_{\text{Modified First part}} d\boldsymbol{\theta} \Big|_{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})} \end{aligned} \quad (20)$$

By the assumption that the distribution  $q^i(\boldsymbol{\theta})$  satisfies mean-field decomposition i.e.

$$\begin{aligned} q^i(\boldsymbol{\theta}) &= \prod_{m=1}^T \mathcal{N}(\theta_{i,m}, \sigma_n^2) \\ &= \prod_{m=1}^T \varrho^i(\boldsymbol{\theta}_{i,m}) \end{aligned} \quad (21)$$

Let us denote  $\mathcal{M}_w = \mathcal{P}_{\theta,w}(\mathcal{D}^i)$ .

Therefore, we extract out the following portion from (20):  $\nabla_{q^i(\theta)} \log \mathcal{P}_{\theta,w}(\mathcal{D}^i)$

$$\nabla_{q^i(\theta)} \log \mathcal{P}_{\theta,w}(\mathcal{D}^i) = \nabla_{q^i(\theta)} \log \mathcal{M}_w \quad (22)$$

We now consider the individual partial differentials here

$$\frac{\partial}{\partial \rho^i(\theta_{i,m})} \log \mathcal{M}_w = \frac{1}{\mathcal{M}_w} \frac{\partial \mathcal{M}_w}{\partial w} \frac{\partial w}{\partial \rho^i(\theta_{i,m})} \quad (23)$$

Thus, we can rewrite (20) from the perspective of individual components of  $q^i(\theta)$  as follows:

$$\begin{aligned} & \int_{\Theta} \left[ -q^i(\theta) \frac{1}{\mathcal{M}_w} \frac{\partial \mathcal{M}_w}{\partial w} \frac{\partial w}{\partial \rho^i(\theta_{i,m})} + \log \mathcal{P}_{\theta,w}(\mathcal{D}^i) \right] d\theta \Big|_{q^i(\hat{\theta}; \mathbf{w})} \\ = & \underbrace{\int_{\Theta} \left[ -q^i(\theta) \frac{1}{\mathcal{P}_{\theta,w}(\mathcal{D}^i)} \frac{\partial \mathcal{P}_{\theta,w}(\mathcal{D}^i)}{\partial w} \frac{\partial w}{\partial \rho^i(\theta_{i,m})} + \log \mathcal{P}_{\theta,w}(\mathcal{D}^i) \right] d\theta}_{\text{Modified First part}} \Big|_{q^i(\hat{\theta}; \mathbf{w})} \end{aligned} \quad (24)$$

Now, we can rewrite the **second part** as follows:

$$\begin{aligned} & \zeta \nabla_{q^i(\theta)} \mathbb{D}_{KL}(q^i(\theta) || z^*(\theta)) \Big|_{q^i(\hat{\theta}; \mathbf{w})} \\ = & \zeta \nabla_{q^i(\theta)} \left[ \int_{\Theta} q^i(\theta) \log q^i(\theta) - q^i(\theta) \log(z^*(\theta)) d\theta \right] \Big|_{q^i(\hat{\theta}; \mathbf{w})} \\ = & \zeta \int_{\Theta} \nabla_{q^i(\theta)} \left[ q^i(\theta) \log q^i(\theta) - q^i(\theta) \log(z^*(\theta)) \right] d\theta \Big|_{q^i(\hat{\theta}; \mathbf{w})} \\ = & \zeta \underbrace{\int_{\Theta} \left( \log q^i(\theta) + 1 - \log(z^*(\theta)) \right) d\theta}_{\text{Modified Second Part}} \Big|_{q^i(\hat{\theta}; \mathbf{w})} \end{aligned} \quad (25)$$

□

Combining both the first and second part we get

$$\begin{aligned} & \int_{\Theta} \left[ -q^i(\hat{\theta}; \mathbf{w}) \frac{1}{\mathcal{P}_{\theta,w}(\mathcal{D}^i)} \frac{\partial \mathcal{P}_{\theta,w}(\mathcal{D}^i)}{\partial w} \frac{\partial w}{\partial \rho^i(\theta_{i,m})} + \log \mathcal{P}_{\theta,w}(\mathcal{D}^i) \right] d\theta \\ & \quad + \zeta \int_{\Theta} \left( \log q^i(\hat{\theta}; \mathbf{w}) + 1 - \log(z^*(\theta)) \right) d\theta = 0 \\ \implies & \zeta \int_{\Theta} \left( \log q^i(\hat{\theta}; \mathbf{w}) + 1 - \log(z^*(\theta)) + \log \mathcal{P}_{\theta,w}(\mathcal{D}^i) \right) d\theta \\ & = \int_{\Theta} \left( q^i(\hat{\theta}; \mathbf{w}) \frac{1}{\mathcal{P}_{\theta,w}(\mathcal{D}^i)} \frac{\partial \mathcal{P}_{\theta,w}(\mathcal{D}^i)}{\partial w} \frac{\partial w}{\partial \rho^i(\theta_{i,m})} \right) d\theta \end{aligned} \quad (26)$$

Let us assume without loss of generality that each of the individual components of the optimal coresot weighted client distribution  $q^i(\hat{\theta}; \mathbf{w})$  can be denoted as some function  $g(\mathbf{w})$ . More, specifically,

$$\begin{aligned}\varrho^i(\boldsymbol{\theta}_{i,j}; \mathbf{w}) &= g_j(\mathbf{w}) \\ \nabla_{\mathbf{w}} \varrho^i(\boldsymbol{\theta}_{i,j}; \mathbf{w}) &= g'_j(\mathbf{w})\end{aligned}\quad (27)$$

Thus we can reuse the above expression to simplify (26)

$$g'_m(\mathbf{w}) = \frac{\int_{\Theta} \left( q^i(\hat{\boldsymbol{\theta}}; \mathbf{w}) \frac{1}{\mathcal{P}_{\theta, \mathbf{w}}(\mathcal{D}^i)} \frac{\partial \mathcal{P}_{\theta, \mathbf{w}}(\mathcal{D}^i)}{\partial w} \right) d\boldsymbol{\theta}}{\zeta \int_{\Theta} \left( \log q^i(\hat{\boldsymbol{\theta}}; \mathbf{w}) + 1 - \log(z^*(\boldsymbol{\theta})) + \log \mathcal{P}_{\theta, \mathbf{w}}(\mathcal{D}^i) \right) d\boldsymbol{\theta}} \quad (28)$$

We now go back to utilizing the above derived expression in our main Eq. (19) to replace  $\nabla_{\mathbf{w}} q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})$

$$\begin{aligned}& \nabla_{\mathbf{w}} q^i(\hat{\boldsymbol{\theta}}; \mathbf{w}) \\ &= \nabla_{\mathbf{w}} \prod_{k=1}^T \varrho^i(\boldsymbol{\theta}_{i,k}; \mathbf{w}) \\ &= \nabla_{\mathbf{w}} \prod_{k=1}^T g_k(\mathbf{w}) \\ &= \prod_{k \neq m}^T g_k(\mathbf{w}) \nabla_{\mathbf{w}} g_m(\mathbf{w}) + g_m(\mathbf{w}) \nabla_{\mathbf{w}} \prod_{k \neq m}^T g_k(\mathbf{w}) \\ &= \prod_{k \neq m}^T g_k(\mathbf{w}) g'_m(\mathbf{w}) + g_m(\mathbf{w}) \nabla_{\mathbf{w}} \prod_{k \neq m}^T g_k(\mathbf{w}) \\ &= \frac{q^i(\hat{\boldsymbol{\theta}}; \mathbf{w})}{\varrho^i(\boldsymbol{\theta}_{i,m}; \mathbf{w})} g'_m(\mathbf{w}) + g_m(\mathbf{w}) \nabla_{\mathbf{w}} \prod_{k \neq m}^T \varrho^i(\boldsymbol{\theta}_{i,k}; \mathbf{w})\end{aligned}\quad (29)$$

Thus, we now have a closed form solution to computing the gradient of the KL divergence  $\mathbb{D}(q^i(\hat{\boldsymbol{\theta}}; \mathbf{w}) \| q^i(\hat{\boldsymbol{\theta}}))$  w.r.t the coreset weight parameters  $\mathbf{w}$ .

**Proposition 2.** *The gradient of the second term in Equation 8 w.r.t  $\mathbf{w}$  i.e.*

$$\nabla_{\mathbf{w}} \|P_{\theta}(\mathcal{D}^i) - P_{\theta, \mathbf{w}}(\mathcal{D}^i)\|_{\hat{\pi}, 2}^2$$

is given by the following expression

$$-2\mathcal{P}_{\Phi}^T \left( \mathcal{P} - \mathcal{P}_{\Phi} \mathbf{w} \right)$$

where  $\mathcal{P} = \sum_{j=1}^n \hat{g}_j$  and  $\mathcal{P}_{\Phi} = [\hat{g}_1, \hat{g}_2, \dots, \hat{g}_n]$

*Proof.* First, we reformulate the given expression in terms

$$\begin{aligned}& \|P_{\theta}(\mathcal{D}^i) - P_{\theta, \mathbf{w}}(\mathcal{D}^i)\|_{\hat{\pi}, 2}^2 \\ &= \mathbb{E}_{\theta \sim \hat{\pi}} [(P_{\theta}(\mathcal{D}^i) - P_{\theta, \mathbf{w}}(\mathcal{D}^i))^2]\end{aligned}$$

We define  $g_j = \mathcal{P}_\theta(\mathcal{D}_j^i) - \mathbb{E}_{\theta \sim \hat{\pi}} \mathcal{P}_\theta(\mathcal{D}_j^i)$

As a result the equivalent optimization problem becomes minimizing  $\left\| \sum_{j=1}^n g_j - \sum_{j=1}^n w_j g_j \right\|_{\hat{\pi}, 2}^2$

Further, using Monte Carlo approximation, given  $S$  samples  $\{\theta_j\}_{j=1}^S$ ,  $\theta_j \sim \hat{\pi}$ , the  $L^2(\hat{\pi})$ -norm can be approximated as follows

$$\left\| \sum_{j=1}^n \hat{g}_j - \sum_{j=1}^n w_j \hat{g}_j \right\|_2^2$$

where

$$\hat{g}_j = \frac{1}{\sqrt{S}} [\mathcal{P}_{\theta_1}(\mathcal{D}_j^i) - \mathcal{P}(\bar{\mathcal{D}}_j^i), \mathcal{P}_{\theta_2}(\mathcal{D}_j^i) - \mathcal{P}(\bar{\mathcal{D}}_j^i), \dots, \mathcal{P}_{\theta_S}(\mathcal{D}_j^i) - \mathcal{P}(\bar{\mathcal{D}}_j^i)] \text{ and } \mathcal{P}(\bar{\mathcal{D}}_j^i) = \frac{1}{S} \sum_{k=1}^S \mathcal{P}_{\theta_k}(\mathcal{D}_j^i)$$

We can write the above problem in matrix notation as follows

$$f(\mathbf{w}) := \|\mathcal{P} - \mathcal{P}_\Phi \mathbf{w}\|_2^2$$

where  $\mathcal{P} = \sum_{j=1}^n \hat{g}_j$  and  $\mathcal{P}_\Phi = [\hat{g}_1, \hat{g}_2, \dots, \hat{g}_n]$

Thus we have the gradient w.r.t  $\mathbf{w}$  as follows:

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = -2\mathcal{P}_\Phi^T (\mathcal{P} - \mathcal{P}_\Phi \mathbf{w}) \quad (30)$$

□

## 10 EXPERIMENTS

All the experiments have been done using the following configuration: Nvidia RTX A4000(16GB) and Apple M2 Pro 10 cores and 16GB memory.

### 10.1 PROPOSAL FOR A MODIFIED OBJECTIVE IN EQUATION 8

$$\{\mathbf{w}_i^*\} \triangleq \arg \min_{\mathbf{w}} \mathbb{D}_{KL}(\hat{q}^i(\boldsymbol{\theta}, \mathbf{w}) \| \hat{q}^i(\boldsymbol{\theta})) + \|P_{\theta}(\mathcal{D}^i) - P_{\theta, \mathbf{w}_i}(\mathcal{D}^i)\|_{\hat{\pi}, 2}^2 \quad \|\mathbf{w}_i\|_0 \leq k \quad (31)$$

We discuss here the utility of our proposed modified client side objective function via an ablation study where we want to gauge the inclusion of the first term in our objective function as just including the coresets loss.

Through experimental analysis, we find that just including the coresets loss optimization results in early saturation, possibly hinting towards getting stuck in local minima, but however inclusion the KL Divergence loss and forcing the coresets weighted local distribution of the client and the normal local distribution of the client to be similar leads to better stability in the training loss and better convergence.

### 10.2 COMMUNICATION COMPLEXITY ANALYSIS FOR DIFFERENT CORESET SIZES

Here we showcase an analysis for different coresets sample size for different datasets and how it affects on the final accuracy and the total number of communication rounds in the Federated Learning setting. This showcases cost-effectiveness of our approach where by using only a small number of communication rounds our proposed approach is able to attain near-optimal performance as per the table below. In addition Fig: 5 substantiates the cost-effectiveness of our approach.

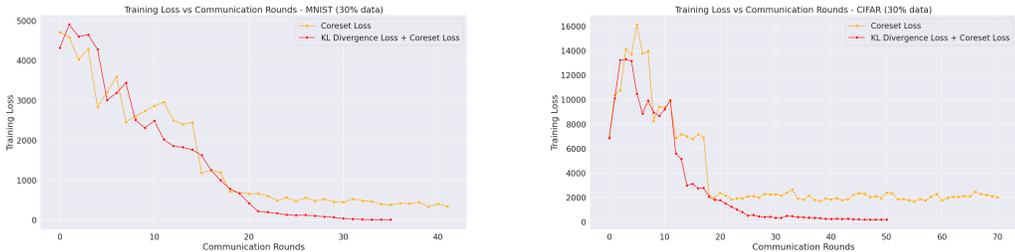


Figure 4: Ablation Study on using KL divergence between two local distribution w.r.t just using coresets weights

Table 3: Comparative results of test accuracies across different coreset sample complexity

Method (Percentage = sampling fraction)	MNIST		FashionMNIST		CIFAR	
	Test Accuracy	Communication Rounds	Test Accuracy	Communication Rounds	Test Accuracy	Communication Rounds
PFEDBAYES (Full)	98.79	194	93.01	215	83.46	266
RANDOMSUBSET (50%)	80.2	135	87.12	172	48.31	183
CORESET-PFEDBAYES (k = 50%)	92.48	98	89.55	93	69.66	112
CORESET-PFEDBAYES (k = 30%)	90.17	84	88.16	72	59.12	70
CORESET-PFEDBAYES (k = 15%)	88.75	62	85.15	38	55.66	32
CORESET-PFEDBAYES (k = 10%)	85.43	32	82.64	24	48.25	16

(a) We report test accuracies across different sample complexity for datasets like MNIST, CIFAR, Fashion-MNIST. Full indicates training on full dataset and 50% is on using half the data size after randomly sampling 50% of the training set.

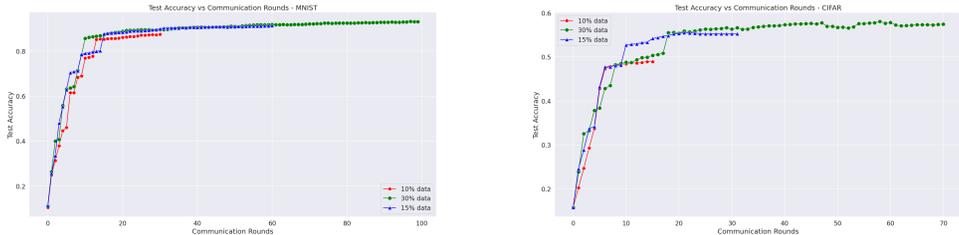


Figure 5: Communication Rounds across Different Sample Size - Convergence analysis

### 10.3 COMPUTING LIKELIHOOD OBJECTIVE USING AIHT

Here, we showcase how we utilised the Accelerated Iterative Hard Thresholding algorithm (A-IHT) for computing the likelihood.

### 10.4 MEDICAL DATASET EXPERIMENT DETAILS

Owing to the rise of Federated Learning based approaches in the medical setting due to privacy-preserving features, we chose to perform our experiments on 3 medical datasets in addition to our main experiments.

For our Federated Learning setup, we considered the setting where we have only 2 clients and one global server.

For each of the 3 datasets in the medical dataset setting, we consider each client has X-ray images of symptomatic type **A**/ type **B** and Normal images . We perform a classification task at each client.

### 10.5 BASLINE COMPARISONS: DIVERSITY BASED SUBMODULAR OPTIMIZATION FUNCTIONS

For our second set of experiments, we chose different diversity based submodular optimization functions, specifically the following functions whose definition have been provided here

**Definition 4. Log-determinant Function** is a diversity-based submodular function. It is non-monotone in nature. Let  $\mathbf{L}$  denote a positive semidefinite kernel matrix and  $\mathbf{L}_{\mathbf{S}}$  denote the subset of rows and columns indexed by set  $\mathbf{S}$ . Log-determinant function  $f$  is specified as:

$$f(\mathbf{S}) = \log\det(\mathbf{L}_{\mathbf{S}}) \quad (32)$$

The log-det function models diversity and is closely related to a determinantal point process.

**Definition 5. Disparity Sum Function** characterizes diversity by considering the sum of distances between every pair of points in a subset  $\mathbf{S}$ . For any two points  $i, j \in \mathbf{S}$ , let  $d_{ij}$  denote the distance between them.

$$f(\mathbf{S}) = \sum_{i,j \in \mathbf{S}} d_{ij} \quad (33)$$

The aim is to select a subset  $\mathbf{S}$  such that  $f(\mathbf{S})$  is maximized. Disparity sum is not a submodular function.

**Definition 6. Disparity Min Function** characterizes diversity by considering the minimum distance between any two non-similar points in a subset  $\mathbf{S}$ .

$$f(\mathbf{S}) = \min_{i,j \in \mathbf{S}, i \neq j} d_{ij} \quad (34)$$

The aim is to select a subset  $\mathbf{S}$  such that  $f(\mathbf{S})$  is maximized. Disparity min is not a submodular function.

For the above experiments we utilise the *Submodlib library*<sup>3</sup> for our implementation [Kaushal et al. \(2022\)](#).

## 10.6 EXPERIMENT CONFIGURATION

### 10.6.1 MNIST EXPERIMENT CONFIGURATION

For both CORESET-PFEDBAYES and corresponding baseline PFEDBAYES, we use a fully connected DNN model with 3 layers [784,100,10] on MNIST dataset.

**Learning rate hyperparameters:** As per [Zhang et al. \(2022b\)](#)'s proposal i.e. PFEDBAYES the learning rates for personalized (client model) and global model ( $\eta_1, \eta_2$ ) are set to 0.001 since these choices result in the best setting for PFEDBAYES. To compare against the stable best hyperparameters of PFEDBAYES, we also fix the same for our proposal CORESET-PFEDBAYES.

**Personalization Hyperparameter:** The  $\zeta$  parameter adjusts the degree of personalization in the case of clients. Again for a fair comparison against our baseline PFEDBAYES, we fix the  $\zeta$  parameter for our proposal CORESET-PFEDBAYES to the best setting given by the baseline. In [Zhang et al. \(2022b\)](#) the authors tune  $\zeta \in \{0.5, 1, 5, 10, 20\}$  and find that  $\zeta = 10$  results in the best setting. We, therefore, fix the personalization parameter  $\zeta = 10$ .

### 10.6.2 MEDICAL DATASETS EXPERIMENT CONFIGURATION

We discuss here the detailed configuration and models used for our further experiments.

Here we specifically consider the setting where we only have 2 clients and a single global server. Each of the 2 clients are assigned with data from only 2 classes along with a shared class for classification purpose.

For example, client 1 has class *A* and *Normal* (shared class) images while client 2 has class *B* and the remaining *Normal* images.

**COVID-19 Radiography Database:** Client 1 has COVID-19 x-ray images while client 2 has lung opacity x-ray images. Normal X-ray images are shared across both clients. Fig. 6 depicts the dataset distribution. For random subset selection, we randomly choose  $\lambda = 0.1$  fraction of samples on the client side. For diversity-based subset selection, we first convert each of the  $299 \times 299$  images into a  $[512 \times 1]$  vector embeddings using a ResNet architecture. Diversity functions are then applied to

<sup>3</sup>Submodlib decile library

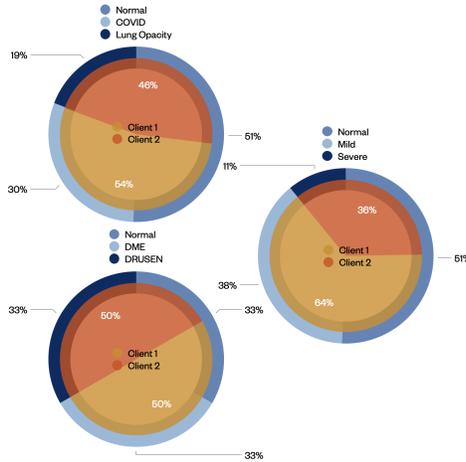


Figure 6: Data distribution for Medical Datasets

these embeddings to retrieve a final subset of diverse and representative embeddings. Eventually, we decode back to the original space using the chosen representative indices.

**APTOS 2019 Blindness Detection:** Unlike the COVID-19 radiography dataset, the APTOS dataset has 3 RGB channels and a higher resolution. We rescale the dimension of images to 299x299 for maintaining uniformity across all datasets. The same model configuration is followed as in the COVID-19 radiography dataset.

**OCTMnist:** The OCTMnist dataset is a large dataset with single-channel images of a higher resolution. We have resized the images to 299x299 resolution for our experiments. The Normal class has above 50,000 train images itself, with the other two classes having close to 10,000 train images. Due to this class imbalance, we have randomly selected 8,000 images from each class for our experiments. Post which we again use a ResNet architecture to reduce the feature dimensions, which we then feed into the CORESET-PFEDBAYES pipeline.

**Baseline : Independent Learning** In this scenario, each of the 2 clients solve the classification problem independently without any involvement of a server as opposed to federated learning. Thus there is no sharing of model weights to a common server as compared to the federated setting.

**Baseline : Independent Learning on other client’s test data** In this scenario, similar to the independent learning setup, we report the metrics for a particular client not only on its own test data but also on the other client’s test data by training on the individual client’s own training data.

For all the experiments for the medical dataset analysis across all the baselines, we report the class-wise accuracy in Table 2.

**Definition 7. Submodular Functions** are set functions which exhibit diminishing returns. Let  $\mathbf{V}$  denote the ground-set of  $n$  data points  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbb{R}^d$ . More formally,  $\mathbf{V} = \{x_i\}_{i=1}^n$ . Let  $\mathbf{A} \subseteq \mathbf{B}$  where  $\mathbf{A}, \mathbf{B} \subset \mathbf{V}$  and  $v \in \mathbf{V}$ . A submodular function  $f : 2^{\mathbf{V}} \mapsto \mathbb{R}$  satisfies the diminishing returns property as follows:

$$f(\mathbf{A} \cup v) - f(\mathbf{A}) \geq f(\mathbf{B} \cup v) - f(\mathbf{B}) \tag{35}$$

10.7 ALGORITHM FOR ACCELERATED IHT

We first present the accelerated IHT algorithm as proposed in Zhang et al. (2021) in Algorithm 10.7.

**Algorithm 2** Accelerated IHT (A-IHT) for Bayesian Coreset Optimization

---

**Input Objective**  $f(w) = \|y - \Phi w\|_2^2$ ; sparsity  $k$

- 1:  $t = 0, z_0 = 0, w_0 = 0$
- 2: repeat
- 3:  $\mathcal{Z} = \text{supp}(z_t)$
- 4:  $\mathcal{S} = \text{supp}(\Pi_{\mathcal{N}_k \setminus \mathcal{Z}}(\nabla f(z_t))) \cup \mathcal{Z}$  where  $|\mathcal{S}| \leq 3k$
- 5:  $\tilde{\nabla}_t = \nabla f(z_t)|_{\mathcal{S}}$
- 6:  $\mu_t = \arg \min_{\mu} f(z_t - \mu \tilde{\nabla}_t) = \frac{\|\tilde{\nabla}_t\|_2^2}{2\|\Phi \tilde{\nabla}_t\|_2^2}$
- 7:  $w_{t+1} = \Pi_{\mathcal{C}_k \cap \mathbb{R}_+^n}(z_t - \mu_t \tilde{\nabla}_t)$
- 8:  $\tau_{t+1} = \arg \min_{\tau} f(w_{t+1} + \tau(w_{t+1} - w_t))$   
 $= \frac{\langle y - \Phi w_{t+1}, \Phi(w_{t+1} - w_t) \rangle}{2\|\Phi(w_{t+1} - w_t)\|_2^2}$
- 9:  $z_{t+1} = w_{t+1} + \tau_{t+1}(w_{t+1} - w_t)$
- 10:  $t = t + 1$
- 11: until Stop criteria met
- 12: return  $w_t$

---

The algorithm **Accelerated IHT** above is proposed by [Zhang et al. \(2021\)](#). We share a high level view of the algorithm include some of the important features.

**Step Size Selection** The authors propose that given the quadratic objective of the coreset optimization, they perform exact line search to obtain the best step size per iteration.  $\frac{\|\tilde{\nabla}_t\|_2^2}{2\|\Phi \tilde{\nabla}_t\|_2^2}$

**Momentum** The authors propose adaptive momentum acceleration as is evident from line 8 of the pseudocode. At the end during the next update, Nesterov Accelerated Gradient is applied as shown in line 9.

## 11 CODE

We share our code on GitHub at [Link](#)