

A. Proofs for density estimation

A.1. Proof of Lemma 1

Lemma 1. *Let the loss be a Bregman divergence B_F . Then, for any $\lambda \in \Lambda \subseteq \Delta_p$, if $h^* = \sum_{k=1}^p \lambda_k \mathcal{D}_k$ is in \mathcal{H} , then it is a minimizer of $h \mapsto \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$. If F is further strictly convex, then it is the unique minimizer.*

Proof. Fix $\lambda \in \Lambda$ such that $\sum_{k=1}^p \lambda_k \mathcal{D}_k$ is in \mathcal{H} . By the non-negativity of the Bregman divergence, for all h , $B_F(\sum_{k=1}^p \lambda_k \mathcal{D}_k \parallel h) \geq 0$ and equality is achieved for $h = \sum_{k=1}^p \lambda_k \mathcal{D}_k$. Thus, h^* is a minimizer of $h \mapsto B_F(\sum_{k=1}^p \lambda_k \mathcal{D}_k \parallel h)$. Since F is strictly convex, $h \mapsto B_F(\sum_{k=1}^p \lambda_k \mathcal{D}_k \parallel h)$ is strictly convex and h^* is therefore the unique minimizer.

Now, for any hypothesis h , observe that the following difference is a constant independent of h :

$$\begin{aligned} & \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h) - B_F\left(\sum_{k=1}^p \lambda_k \mathcal{D}_k \parallel h\right) \\ &= \sum_{k=1}^p \lambda_k [F(\mathcal{D}_k) - F(h) - \langle \nabla F(h), \mathcal{D}_k - h \rangle] - \left[F\left(\sum_{k=1}^p \lambda_k \mathcal{D}_k\right) - F(h) - \left\langle \nabla F(h), \sum_{k=1}^p \lambda_k \mathcal{D}_k - h \right\rangle \right] \\ &= \sum_{k=1}^p \lambda_k F(\mathcal{D}_k) - F\left(\sum_{k=1}^p \lambda_k \mathcal{D}_k\right). \end{aligned} \tag{12}$$

Thus, h^* is also the unique minimizer of $h \mapsto \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$. \square

A.2. Proof of Lemma 2

Lemma 2. *Let the loss be a Bregman divergence B_F with F strictly convex and assume that $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\}) \subseteq \mathcal{H}$. Observe that B_F is jointly convex in both arguments. Then, for any convex set $\Lambda \subseteq \Delta_p$, the solution of the optimization problem $\min_{h \in \mathcal{H}} \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$ exists and is in $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\})$.*

Proof. Let \mathcal{H}' is the closure of convex hull of \mathcal{H} . Observe that \mathcal{H}' is a convex and compact set.

$$\min_{h \in \mathcal{H}'} \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h) \leq \min_{h \in \mathcal{H}} \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h).$$

We show that minimizer over \mathcal{H}' exists and is in the $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\})$. Since $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\}) \subseteq \mathcal{H} \subseteq \mathcal{H}'$, the minimizer over \mathcal{H} also exists and is in the $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\})$.

Since B_F is convex with respect to its second argument, $h \mapsto \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$ is a convex function of h defined over the convex set \mathcal{H}' . Since any maximum of a convex function is also convex, $h \mapsto \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$ is a convex function and its minimum over the compact set \mathcal{H}' exists.

We now show that the minimizer is in $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\})$. Notice that, since $\sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$ is linear in λ , we have

$$\max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h) = \max_{\lambda \in \text{conv}(\Lambda)} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h).$$

Thus, it suffices to consider the case $\Lambda \subseteq \Delta_p$. Then, since \mathcal{H}' is a compact and convex set and since B_F is convex with respect to its second argument, by Sion's minimax theorem, we can write:

$$\min_{h \in \mathcal{H}'} \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h) = \max_{\lambda \in \Lambda} \min_{h \in \mathcal{H}'} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h).$$

Let $\lambda^{\text{opt}} = \arg\max_{\lambda \in \Lambda} \min_{h \in \mathcal{H}'} \sum_{k=1}^p \lambda_k B_F(\mathcal{D}_k \parallel h)$ and $h^* = \sum_{k=1}^p \lambda_k^{\text{opt}} \mathcal{D}_k$. By assumption, $\text{conv}(\{\mathcal{D}_1, \dots, \mathcal{D}_p\})$ is included in \mathcal{H}' , thus h^* is in \mathcal{H}' and, by Lemma 1, h^* is a minimizer of $h \mapsto \sum_{k=1}^p \lambda_k^{\text{opt}} B_F(\mathcal{D}_k \parallel h)$. In view of that, if h'

is a minimizer of $h \mapsto \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k \mathbf{B}_F(\mathcal{D}_k \parallel h)$ over \mathcal{H}' , then the following holds:

$$\begin{aligned}
 \max_{\lambda} \sum_{k=1}^p \lambda_k \mathbf{B}_F(\mathcal{D}_k \parallel h') &\geq \sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h') && \text{(def. of max)} \\
 &\geq \sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h^*) && \text{(Lemma 1)} \\
 &= \min_{h \in \mathcal{H}'} \sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h) && (h^* \text{ minimizer}) \\
 &= \max_{\lambda \in \Lambda} \min_{h \in \mathcal{H}'} \sum_{k=1}^p \lambda_k \mathbf{B}_F(\mathcal{D}_k \parallel h) && \text{(def. of } \lambda_k^{\text{opt}} \text{)} \\
 &= \min_{h \in \mathcal{H}'} \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k \mathbf{B}_F(\mathcal{D}_k \parallel h). && \text{(Sion's minimax theorem)}
 \end{aligned}$$

By the optimality of h' , the first and last expressions in this chain of inequalities are equal, which implies the equality of all intermediate terms. In particular, this implies $\sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h') = \sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h^*)$. Since F is strictly convex, by Lemma 1, the minimizer of $h \mapsto \sum_{k=1}^p \lambda_k^{\text{opt}} \mathbf{B}_F(\mathcal{D}_k \parallel h)$ is unique and $h' = h^*$. This completes the proof. \square

B. Convergence guarantee of FEDBOOST (Theorem 2)

Theorem 2. *If Properties 1 hold and $\eta = \sqrt{\frac{\sigma}{TG^2r\alpha}}$, then α^A , the output of FEDBOOST satisfies,*

$$\mathbb{E} [\mathbf{L}(\alpha^A) - \mathbf{L}(\alpha_{\text{opt}})] \leq 2\sqrt{\frac{G^2\sigma r\alpha}{T}} + \frac{\alpha_* M}{2T} \sum_{t=1}^T \sum_{k=1}^q \frac{\alpha_{k,t}^2}{\gamma_{k,t}}.$$

Proof. By Jensen's inequality,

$$\mathbf{L}(\alpha^A) \leq \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\alpha_t).$$

Hence, it suffices to bound

$$\frac{1}{T} \sum_{t=1}^T (\mathbf{L}(\alpha_t) - \mathbf{L}(\alpha)).$$

For any t ,

$$\begin{aligned}
 \mathbf{L}(\alpha_t) - \mathbf{L}(\alpha) &\leq \langle \nabla \mathbf{L}(\alpha_t), \alpha_t - \alpha \rangle && \text{(convexity of L)} \\
 &= \langle \delta_t \mathbf{L}, \alpha_t - \alpha \rangle + \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle \\
 &= \frac{1}{\eta} \langle \nabla F(\alpha_t) - \nabla F(v_{t+1}), \alpha_t - \alpha \rangle + \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle && \text{(def. of } v_{t+1} \text{)} \\
 &= \frac{1}{\eta} (\mathbf{B}_F(\alpha \parallel \alpha_t) + \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha \parallel v_{t+1})) + \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle && \text{(Bregman div. def.)} \\
 &\leq \frac{1}{\eta} (\mathbf{B}_F(\alpha \parallel \alpha_t) + \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha \parallel \alpha_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})) && \text{(13a)} \\
 &\quad + \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle, && \text{(13b)}
 \end{aligned}$$

where the last inequality follows because $\mathbf{B}_F(\alpha \parallel v_{t+1}) \geq \mathbf{B}_F(\alpha \parallel \alpha_{t+1}) + \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})$ by the generalized

Pythagorean inequality. For the first term (13a), summing over t gives the following telescoping sum,

$$\begin{aligned}
 & \sum_{t=1}^T (\mathbf{B}_F(\alpha \parallel \alpha_t) + \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha \parallel \alpha_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})) \\
 &= \mathbf{B}_F(\alpha \parallel \alpha_1) - \mathbf{B}_F(\alpha \parallel \alpha_{T+1}) + \sum_{t=1}^T \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1}) \\
 &\leq \mathbf{B}_F(\alpha \parallel \alpha_1) + \sum_{t=1}^T (\mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})).
 \end{aligned} \tag{14}$$

Now consider the summation term:

$$\begin{aligned}
 & \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1}) = F(\alpha_t) - F(\alpha_{t+1}) - \langle \nabla F(v_{t+1}), \alpha_t - \alpha_{t+1} \rangle \\
 &\leq \langle \nabla F(\alpha_t), \alpha_t - \alpha_{t+1} \rangle - \frac{\sigma}{2} \|\alpha_t - \alpha_{t+1}\|^2 - \langle \nabla F(v_{t+1}), \alpha_t - \alpha_{t+1} \rangle \quad (\text{strong convexity of } F) \\
 &= \langle \nabla F(\alpha_t) - \nabla F(v_{t+1}), \alpha_t - \alpha_{t+1} \rangle - \frac{\sigma}{2} \|\alpha_t - \alpha_{t+1}\|^2 \\
 &= \eta \langle \delta_t \mathbf{L}, \alpha_t - \alpha_{t+1} \rangle - \frac{\sigma}{2} \|\alpha_t - \alpha_{t+1}\|^2 \quad (\text{def. of } v_{t+1}) \\
 &\leq \eta \|\delta_t \mathbf{L}\|_* \|\alpha_t - \alpha_{t+1}\| - \frac{\sigma}{2} \|\alpha_t - \alpha_{t+1}\|^2 \quad (\text{Cauchy-Schwarz ineq.}) \\
 &\leq \frac{\eta^2 \|\delta_t \mathbf{L}\|_*^2}{2\sigma}.
 \end{aligned} \tag{15}$$

Combining the above inequalities,

$$\begin{aligned}
 \sum_{t=1}^T (\mathbf{L}(\alpha_t) - \mathbf{L}(\alpha)) &\leq \frac{1}{\eta} \mathbf{B}_F(\alpha \parallel \alpha_1) + \sum_{t=1}^T \left(\frac{\eta \|\delta_t \mathbf{L}\|_*^2}{2\sigma} + \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle \right) \\
 &\leq \frac{1}{\eta} \mathbf{B}_F(\alpha \parallel \alpha_1) + \frac{\eta G^2 T}{2\sigma} + \sum_{t=1}^T (\langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle).
 \end{aligned}$$

We now bound (13b) in expectation, the inner product term in the above equation. Denote by $\nabla_t \mathbf{L}(\cdot) := \sum_{j \in S_t} \frac{m_j}{m} \nabla \mathbf{L}_j(\cdot)$, where $m = \sum_{j \in S_t} m_j$. Taking the expectation over $j \in S_t$,

$$\begin{aligned}
 \mathbb{E} \left[\sum_{t=1}^T \langle \nabla \mathbf{L}(\alpha_t) - \delta_t \mathbf{L}, \alpha_t - \alpha \rangle \right] &= \sum_{t=1}^T \langle \nabla \mathbf{L}(\alpha_t) - \mathbb{E}[\delta_t \mathbf{L}], \alpha_t - \alpha \rangle \\
 &= \sum_{t=1}^T \langle \nabla \mathbf{L}(\alpha_t) - \mathbb{E}[\nabla_t \mathbf{L}(\tilde{\alpha}_t)], \alpha_t - \alpha \rangle \\
 &\leq \sum_{t=1}^T \|\nabla \mathbf{L}(\alpha_t) - \mathbb{E}[\nabla_t \mathbf{L}(\tilde{\alpha}_t)]\|_* \|\alpha_t - \alpha\| \quad (\text{Cauchy-Schwarz ineq.}) \\
 &\leq \sum_{t=1}^T \|\nabla \mathbf{L}(\alpha_t) - \mathbb{E}[\nabla_t \mathbf{L}(\tilde{\alpha}_t)]\|_* \alpha_*. \quad (\text{by Prop. 1.2.})
 \end{aligned} \tag{16}$$

To understand $\mathbb{E}[\nabla_t \mathbf{L}(\tilde{\alpha}_t)]$, we use Taylor's Theorem in several variables (Folland, 2010). Let $f = \nabla_t \mathbf{L}$. Expanding $f(\tilde{\alpha}_t)$ about α_t ,

$$f(\tilde{\alpha}_t) - f(\alpha_t) = \nabla f(\alpha_t)(\tilde{\alpha}_t - \alpha_t) + R_1(\tilde{\alpha}_t - \alpha_t),$$

where $R_1(\cdot)$ is the reminder term that can be bounded as

$$\|R_1(\tilde{\alpha}_t - \alpha_t)\|_* \leq \frac{M}{2!} \|\tilde{\alpha}_t - \alpha_t\|_2^2,$$

with $\|\nabla^2 f(\cdot)\| \leq M$. Taking expectation over $\tilde{\alpha}_t$ and using the fact that $\mathbb{E}[\tilde{\alpha}_t] = \alpha_t$, we get

$$\begin{aligned} \|\mathbb{E}[f(\tilde{\alpha}_t)] - f(\alpha_t)\| &\leq \frac{M}{2} \mathbb{E} \|\tilde{\alpha}_t - \alpha_t\|_2^2 \\ &= \frac{M}{2} \sum_{k=1}^q \mathbb{E} \left\| \frac{\alpha_{k,t} \mathbf{1}_{k,t}}{\gamma_{k,t}} - \alpha_{k,t} \right\|_2^2 \\ &= \frac{M}{2} \sum_{k=1}^q \alpha_{k,t}^2 \mathbb{E} \left\| \frac{\mathbf{1}_{k,t}}{\gamma_{k,t}} - 1 \right\|_2^2 \\ &= \frac{M}{2} \sum_{k=1}^q \alpha_{k,t}^2 \left(\frac{1 - \gamma_{k,t}}{\gamma_{k,t}} \right) \\ &\leq \frac{M}{2} \sum_{k=1}^q \frac{\alpha_{k,t}^2}{\gamma_{k,t}}. \end{aligned}$$

Combining the resulting inequalities gives

$$\mathbb{L}(\alpha^A) - \mathbb{L}(\alpha) \leq \frac{1}{\eta T} \mathbb{B}_F(\alpha \parallel \alpha_1) + \frac{\eta G^2}{2\sigma} + \frac{\alpha_* M}{2T} \sum_{t=1}^T \sum_{k=1}^q \frac{\alpha_{k,t}^2}{\gamma_{k,t}}.$$

Choosing the learning rate yields the theorem. \square

C. Convergence guarantee for AFLBOOST (Theorem 3)

Theorem 3. *Let Properties 1 and 2 hold. Let $\eta_\lambda = \sqrt{\frac{\sigma}{TG_\lambda^2 r_\lambda}}$ and $\eta_\alpha = \sqrt{\frac{\sigma}{TG_\alpha^2 r_\alpha}}$. Let α^A be the output of AFLBOOST. If $\gamma_{k,t}$ is given by 8, then $\mathbb{E}[\max_{\lambda \in \Lambda} \mathbb{L}(\alpha^A, \lambda) - \min_{\alpha \in \Delta_q} \max_{\lambda \in \Lambda} \mathbb{L}(\alpha, \lambda)]$ is at most*

$$4\sqrt{\frac{G_\alpha^2(\sigma r_\alpha + \alpha_*)}{T}} + 4\sqrt{\frac{G_\lambda^2(\sigma r_\lambda + \lambda_*)}{T}} + \frac{M(\lambda_* + \alpha_*)}{C}.$$

Proof. By Mohri et al. (2019)[Lemma 5], it suffices to bound

$$\frac{1}{T} \max_{\substack{\lambda \in \Lambda; \\ \alpha \in \Delta_q}} \left\{ \sum_{t=1}^T \mathbb{L}(\alpha_t, \lambda) - \mathbb{L}(\alpha, \lambda_t) \right\}. \quad (18)$$

Consider the following inequalities:

$$\begin{aligned} \mathbb{L}(\alpha_t, \lambda) - \mathbb{L}(\alpha, \lambda_t) &= \mathbb{L}(\alpha_t, \lambda) - \mathbb{L}(\alpha_t, \lambda_t) + \mathbb{L}(\alpha_t, \lambda_t) - \mathbb{L}(\alpha, \lambda_t) \\ &\leq \langle \nabla_\lambda \mathbb{L}(\alpha_t, \lambda_t), \lambda - \lambda_t \rangle + \langle \nabla_\alpha \mathbb{L}(\alpha_t, \lambda_t), \alpha_t - \alpha \rangle && \text{(convexity of } \mathbb{L} \text{)} \\ &= \langle \delta_{\lambda,t} \mathbb{L}, \lambda - \lambda_t \rangle + \langle \delta_{\alpha,t} \mathbb{L}, \alpha_t - \alpha \rangle \\ &+ \langle \nabla_\lambda \mathbb{L}(\alpha_t, \lambda_t) - \delta_{\lambda,t} \mathbb{L}, \lambda - \lambda_t \rangle + \langle \nabla_\alpha \mathbb{L}(\alpha_t, \lambda_t) - \delta_{\alpha,t} \mathbb{L}, \alpha_t - \alpha \rangle \end{aligned}$$

Given these inequalities, we can bound (18) using the sub-additive property of max on the previous inequality as follows:

$$\begin{aligned} \max_{\substack{\lambda \in \Lambda; \\ \alpha \in \Delta_q}} \left\{ \sum_{t=1}^T L(\alpha_t, \lambda) - L(\alpha, \lambda_t) \right\} \\ \leq \max_{\substack{\lambda \in \Lambda; \\ \alpha \in \Delta_q}} \sum_{t=1}^T \{ \langle \delta_{\lambda,t} \mathbf{L}, \lambda - \lambda_t \rangle + \langle \delta_{\alpha,t} \mathbf{L}, \alpha_t - \alpha \rangle \} \end{aligned} \quad (19a)$$

$$+ \max_{\substack{\lambda \in \Lambda; \\ \alpha \in \Delta_q}} \sum_{t=1}^T \{ \langle \lambda, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \mathbf{L} \rangle + \langle \alpha, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \mathbf{L} \rangle \} \quad (19b)$$

$$+ \sum_{t=1}^T \langle \lambda_t, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \mathbf{L} \rangle + \langle \alpha_t, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \mathbf{L} \rangle, \quad (19c)$$

which we will bound in three parts. Consider the first sub-equation (19a): similarly in arriving at (13a), it follows by definition of w_{t+1}, v_{t+1} that

$$\begin{aligned} \langle \delta_{\lambda,t} \mathbf{L}, \lambda - \lambda_t \rangle + \langle \delta_{\alpha,t} \mathbf{L}, \alpha_t - \alpha \rangle &\leq \frac{1}{\eta_{\lambda}} (\mathbf{B}_F(\lambda \parallel \lambda_t) + \mathbf{B}_F(\lambda_t \parallel w_{t+1}) - \mathbf{B}_F(\lambda \parallel \lambda_{t+1}) - \mathbf{B}_F(\lambda_{t+1} \parallel w_{t+1})) \\ &\quad + \frac{1}{\eta_{\alpha}} (\mathbf{B}_F(\alpha \parallel \alpha_t) + \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha \parallel \alpha_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})). \end{aligned}$$

Summing over t , this gives the following by similar argument as in (14) for all λ, α :

$$\begin{aligned} \sum_{t=1}^T \langle \delta_{\lambda,t} \mathbf{L}, \lambda - \lambda_t \rangle + \langle \delta_{\alpha,t} \mathbf{L}, \alpha_t - \alpha \rangle &\leq \frac{1}{\eta_{\lambda}} (\mathbf{B}_F(\lambda \parallel \lambda_1) + \sum_{t=1}^T \mathbf{B}_F(\lambda_t \parallel w_{t+1}) - \mathbf{B}_F(\lambda_{t+1} \parallel w_{t+1})) \\ &\quad + \frac{1}{\eta_{\alpha}} (\mathbf{B}_F(\alpha \parallel \alpha_1) + \sum_{t=1}^T \mathbf{B}_F(\alpha_t \parallel v_{t+1}) - \mathbf{B}_F(\alpha_{t+1} \parallel v_{t+1})) \end{aligned}$$

In view of the inequality resulting from (15), for all λ, α , this is bounded by

$$\begin{aligned} \sum_{t=1}^T \langle \delta_{\lambda,t} \mathbf{L}, \lambda - \lambda_t \rangle + \langle \delta_{\alpha,t} \mathbf{L}, \alpha_t - \alpha \rangle &\leq \frac{1}{\eta_{\lambda}} \mathbf{B}_F(\lambda \parallel \lambda_1) + \frac{1}{\eta_{\alpha}} \mathbf{B}_F(\alpha \parallel \alpha_1) + \sum_{t=1}^T \frac{\eta_{\lambda}^2 \|\delta_{\lambda,t} \mathbf{L}\|_*^2 + \eta_{\alpha}^2 \|\delta_{\alpha,t} \mathbf{L}\|_*^2}{2\sigma} \\ &= \frac{1}{\eta_{\lambda}} \mathbf{B}_F(\lambda \parallel \lambda_1) + \frac{1}{\eta_{\alpha}} \mathbf{B}_F(\alpha \parallel \alpha_1) + \frac{T(\eta_{\lambda} G_{\lambda}^2 + \eta_{\alpha} G_{\alpha}^2)}{2\sigma}. \end{aligned}$$

Next, we proceed with the bound for third sub-equation (19c) in expectation via similar argument followed to arrive at (15):

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \langle \lambda_t, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \mathbf{L} \rangle + \langle \alpha_t, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \mathbf{L} \rangle \right] \\ = \sum_{t=1}^T \langle \lambda_t, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \mathbb{E}[\nabla_{t,\lambda} L(\tilde{\alpha}_t, \lambda_t)] \rangle + \langle \alpha_t, \nabla_{\alpha} L(\tilde{\alpha}_t, \lambda_t) - \mathbb{E}[\nabla_{t,\alpha} L(\tilde{\alpha}_t, \lambda_t)] \rangle, \end{aligned}$$

where $\nabla_{t,\lambda} L(\cdot) := \sum_{j \in S_t} \frac{m_j}{m} \nabla_{\lambda} L_j(\cdot)$, and similarly for $\nabla_{t,\alpha} L(\cdot)$. Similar to the proof of (15) and (9), it can be shown that

$$\sum_{t=1}^T \langle \lambda_t, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \mathbb{E}[\nabla_{t,\lambda}] \rangle \leq \frac{MT\lambda_*}{2C}.$$

Similarly,

$$\mathbb{E} \left[\sum_{t=1}^T \langle \alpha_t, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \mathbf{L} \rangle \right] \leq \frac{MT\alpha_*}{2C}.$$

Combining the two bounds, we have

$$\mathbb{E}\left[\sum_{t=1}^T \langle \lambda_t, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \rangle + \langle \alpha_t, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \rangle\right] \leq \frac{MT(\alpha_* + \lambda_*)}{2C}.$$

We now consider the second sub-equation term (19b), focusing on the first summand with the max over λ and bound this by the Cauchy-Schwarz inequality, then Jensen's inequality:

$$\begin{aligned} & \mathbb{E}\left[\max_{\lambda \in \Lambda} \left\{ \sum_{t=1}^T \langle \lambda, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \rangle \right\}\right] \\ & \leq \mathbb{E}\left[\max_{\lambda \in \Lambda} \left\{ \sum_{t=1}^T \langle \lambda, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \mathbb{E}[\delta_{\lambda,t}] \rangle \right\}\right] + \mathbb{E}\left[\max_{\lambda \in \Lambda} \left\{ \sum_{t=1}^T \langle \lambda, \delta_{\lambda,t} - \mathbb{E}[\delta_{\lambda,t}] \rangle \right\}\right] \\ & \leq \frac{MT\lambda_*}{2C} + \lambda_* G_{\lambda} \sqrt{T}, \end{aligned}$$

where λ_* denotes the max over the compact set Λ . Similarly, we can obtain the following inequality:

$$\mathbb{E}\left[\max_{\alpha \in \Delta_q} \sum_{t=1}^T \langle \alpha, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \rangle\right] \leq \frac{MT\alpha_*}{2C} + \alpha_* G_{\alpha} \sqrt{T}.$$

Thus, combining the inequalities gives

$$\max_{\substack{\lambda \in \Lambda; \\ \alpha \in \Delta_q}} \sum_{t=1}^T \{ \langle \lambda, \nabla_{\lambda} L(\alpha_t, \lambda_t) - \delta_{\lambda,t} \rangle + \langle \alpha, \nabla_{\alpha} L(\alpha_t, \lambda_t) - \delta_{\alpha,t} \rangle \} = \frac{MT(\lambda_* + \alpha_*)}{2C} + \alpha_* G_{\alpha} \sqrt{T} + \lambda_* G_{\lambda} \sqrt{T}.$$

Combining the bounds for (19a), (19b), and (19c), the following bound holds:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \eta_{\alpha} L(\alpha_t, \lambda) - \eta_{\lambda} L(\alpha, \lambda_t) \\ & \leq \frac{1}{T} \left(\frac{\mathbf{B}_F(\lambda \parallel \lambda_1)}{\eta_{\lambda}} + \frac{\mathbf{B}_F(\alpha \parallel \alpha_1)}{\eta_{\alpha}} \right) + \frac{T(\eta_{\lambda} G_{\lambda}^2 + \eta_{\alpha} G_{\alpha}^2)}{2\sigma} + \frac{M(\lambda_* + \alpha_*)}{C} + \frac{\alpha_* G_{\alpha} + \lambda_* G_{\lambda}}{\sqrt{T}}. \end{aligned}$$

□

D. Additional density estimation experiments on synthetic data

Continuing the experimental validation of FEDBOOST as described in 5.1, we examine the effect of modulating the communication budget C on a density estimation task using the same setup as before, but with a power-law distributed synthetic dataset with parameter $p = 1000$.

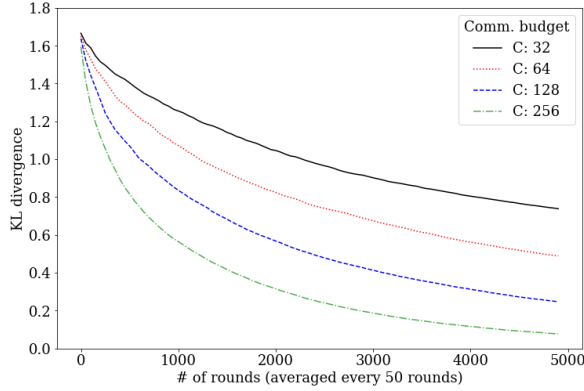


Figure 6. Comparison of loss curves as a function of C using *uniform sampling* in density estimation on synthetic data.

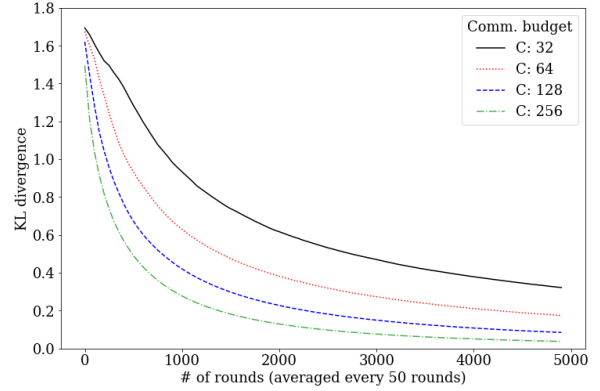


Figure 7. Comparison of convergence as C varies using *weighted random sampling* for density estimation.

We use a hand-tuned step size $\eta = 0.001$ for all values of C , and include ℓ_1 regularization in the experiment using *weighted random sampling* (Fig. 7). The experimental setup is otherwise the same for both Fig. 6 and 7. Across all values of C , *weighted random sampling* of the h_k achieves lower loss than using *uniform sampling*, which validates that using weighted random sampling reduces the communication-dependent term of FEDBOOST.