

## 4 Centralized Learning

---

**4.1** In this problem we will demonstrate the linear performance gain promised by relation (4.46) in code for Example 4.2. Formulate a model  $w^o$  of your choice, and generate data according to the observation model (4.39), ensuring agents are homogenous by sampling  $\mathbf{h}_k$  and  $\mathbf{v}_k$  from identical distributions. Reasonable choices are  $\mathbf{h}_k \sim \mathcal{N}(0, \sigma_h^2 I_M)$  and  $\mathbf{v}_k \sim \mathcal{N}(0, \sigma_v^2)$ . Implement recursions (4.42)–(4.43) and plot the evolution of the error over time. Compare the performance for  $K = 1$ ,  $K = 10$  and  $K = 100$  agents and verify whether you observe linear gains in performance.

**Solution.** The solution is provided as a Jupyter notebook in the separate file `Problem_4.1.ipynb`.

**4.2** Consider local empirical risk minimization problems with unbalanced datasets of the form:

$$J_k(w) = \frac{1}{N_k} \sum_{n=1}^{N_k} Q(w; x_{k,n})$$

and the global empirical risk minimization problem:

$$J(w) = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^{N_k} Q(w; x_{k,n})$$

where we defined  $N = \sum_{k=1}^K N_k$ . Show that  $J(w)$  can be written as a *weighted* consensus problem of the form:

$$J(w) = \sum_{k=1}^K p_k J_k(w)$$

where  $p_k$  are suitably chosen weights satisfying  $\sum_{k=1}^K p_k = 1$ .

**Solution.** We begin reformulating  $J(w)$  as:

$$\begin{aligned}
 J(w) &= \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^{N_k} Q(w; x_{k,n}) \\
 &= \frac{1}{N} \sum_{k=1}^K N_k \cdot \frac{1}{N_k} \sum_{n=1}^{N_k} Q(w; x_{k,n}) \\
 &= \sum_{k=1}^K \frac{N_k}{N} J_k(w) \\
 &= \sum_{k=1}^K p_k J_k(w)
 \end{aligned}$$

where we defined:

$$p_k = \frac{N_k}{N}$$

It follows that the weights  $p_k$  are given by the fraction of data available to agent  $k$ . Local objectives are weighted higher for agents who have more data, which is reasonable since they better approximate the true data distribution.

## 5 Federated Learning

---

**5.1** Suppose in the federated averaging recursion we do not employ normalized step-sizes as suggested by (5.26), but instead let  $\mu_k = \mu$ . Relation (5.25) suggests that in that case we can no longer interpret federated averaging as employing an unbiased gradient approximation to the consensus problem (5.1). Find an expression for which the approximation is unbiased, and use this insight to determine the limiting point of the federated averaging algorithm without step-size normalization.

**Solution.** From (5.25), we have:

$$\begin{aligned}\mathbb{E} \left\{ \widehat{\nabla J}^{\text{fed}}(\mathbf{w}_{i-1}) | \mathbf{w}_{i-1} \right\} &= \frac{1}{K} \sum_{k=1}^K \frac{\mu_k}{\mu} E_k \nabla J_k(\mathbf{w}_{i-1}) \\ &= \frac{1}{K} \sum_{k=1}^K \frac{\mu_k}{\mu} E_k \nabla J_k(\mathbf{w}_{i-1}) \\ &\stackrel{\mu_k = \mu}{=} \frac{1}{K} \sum_{k=1}^K E_k \nabla J_k(\mathbf{w}_{i-1})\end{aligned}$$

We identify  $\mathbb{E} \left\{ \widehat{\nabla J}^{\text{fed}}(\mathbf{w}_{i-1}) | \mathbf{w}_{i-1} \right\}$  as the gradient of:

$$J'(w) \triangleq \frac{1}{K} \sum_{k=1}^K E_k J_k(w)$$

We conclude that the gradient federated averaging algorithm without step-size normalization approximates the gradient of a weighted objective  $J(w)$ , where the weighting is determined by the number  $E_k$  of local steps taken.

**5.2** Prove Lemma 5.1.

**Solution.** Using a similar technique as in (5.23)–(5.24), we introduce the selection indicator:

$$\mathbb{1}_k = \begin{cases} 1, & \text{if the } k\text{-th sample is picked,} \\ 0, & \text{otherwise.} \end{cases} \quad (5.1)$$

We can then write:

$$\hat{\mathbf{a}} = \frac{1}{L} \sum_{\ell=1}^L \mathbf{a}_\ell = \frac{1}{L} \sum_{k=1}^K \mathbb{1}_k \mathbf{a}_k \quad (5.2)$$

It holds that:

$$\mathbb{E} \mathbb{1}_k = \frac{L}{K} \quad (5.3)$$

$$\begin{aligned} \mathbb{E} \{ \mathbb{1}_k \mathbb{1}_\ell \} &= \mathbb{E} \{ \mathbb{1}_k \mathbb{1}_\ell | \mathbb{1}_\ell = 0 \} \cdot \Pr \{ \mathbb{1}_\ell = 0 \} + \mathbb{E} \{ \mathbb{1}_k \mathbb{1}_\ell | \mathbb{1}_\ell = 1 \} \cdot \Pr \{ \mathbb{1}_\ell = 1 \} \\ &= \mathbb{E} \{ \mathbb{1}_k \mathbb{1}_\ell | \mathbb{1}_\ell = 1 \} \cdot \Pr \{ \mathbb{1}_\ell = 1 \} \\ &= \mathbb{E} \{ \mathbb{1}_k | \mathbb{1}_\ell = 1 \} \cdot \frac{L}{K} \\ &= \frac{L-1}{K-1} \cdot \frac{L}{K} \end{aligned} \quad (5.4)$$

We can then conclude for the mean:

$$\mathbb{E} \hat{\mathbf{a}} = \mathbb{E} \frac{1}{L} \sum_{k=1}^K \mathbb{1}_k \mathbf{a}_k = \frac{1}{L} \sum_{k=1}^K \frac{L}{K} \mathbb{E} \mathbf{a}_k = \bar{\mathbf{a}} \quad (5.5)$$

For the variance:

$$\begin{aligned} &\mathbb{E} \|\hat{\mathbf{a}} - \bar{\mathbf{a}}\|^2 \\ &= \mathbb{E} \|\hat{\mathbf{a}}\|^2 - \|\bar{\mathbf{a}}\|^2 \\ &= \frac{1}{L^2} \sum_{k=1}^K \mathbb{E} \{ \mathbb{1}_k^2 \} \mathbb{E} \{ \|\mathbf{a}_k\|^2 \} + \frac{1}{L^2} \sum_{k=1}^K \sum_{\ell \neq k}^K \mathbb{E} \{ \mathbb{1}_k \mathbb{1}_\ell \} \mathbb{E} \{ \mathbf{a}_k^\top \mathbf{a}_\ell \} - \|\bar{\mathbf{a}}\|^2 \\ &= \frac{1}{L^2} \sum_{k=1}^K \frac{L}{K} \left( \sigma_{a_k}^2 + \|\bar{\mathbf{a}}_k\|^2 \right) + \frac{1}{L^2} \sum_{k=1}^K \sum_{\ell \neq k}^K \frac{L-1}{K-1} \frac{L}{K} \bar{\mathbf{a}}_k^\top \bar{\mathbf{a}}_\ell - \|\bar{\mathbf{a}}\|^2 \\ &= \frac{1}{KL} \sum_{k=1}^K \left( \sigma_{a_k}^2 + \|\bar{\mathbf{a}}_k\|^2 \right) + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \sum_{\ell \neq k}^K \bar{\mathbf{a}}_k^\top \bar{\mathbf{a}}_\ell - \|\bar{\mathbf{a}}\|^2 \end{aligned} \quad (5.6)$$

We are almost done, and now aim to reformulate the expression to be consistent

with Lemma 5.1 by completing the square:

$$\begin{aligned}
& \mathbb{E} \|\hat{\mathbf{a}} - \bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \left(1 - \frac{L-1}{K-1}\right) \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 \\
&\quad + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \sum_{\ell \neq k} \bar{\mathbf{a}}_k^\top \bar{\mathbf{a}}_\ell - \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \left(\frac{K-1}{K-1} - \frac{L-1}{K-1}\right) \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 \\
&\quad + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \sum_{\ell \neq k} \bar{\mathbf{a}}_k^\top \bar{\mathbf{a}}_\ell - \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 \\
&\quad + \frac{1}{KL} \frac{L-1}{K-1} \sum_{k=1}^K \sum_{\ell \neq k} \bar{\mathbf{a}}_k^\top \bar{\mathbf{a}}_\ell - \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \frac{1}{KL} \frac{L-1}{K-1} \left\| \sum_{k=1}^K \bar{\mathbf{a}}_k \right\|^2 - \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \frac{K}{L} \frac{L-1}{K-1} \|\bar{\mathbf{a}}\|^2 - \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 + \left( \frac{K}{L} \frac{L-1}{K-1} - \frac{L}{L} \frac{K-1}{K-1} \right) \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}}_k\|^2 - \frac{1}{L} \frac{K-L}{K-1} \|\bar{\mathbf{a}}\|^2 \\
&= \frac{1}{KL} \sum_{k=1}^K \sigma_{a_k}^2 + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K \|\bar{\mathbf{a}} - \bar{\mathbf{a}}_k\|^2
\end{aligned} \tag{5.7}$$

Where the last line can be verified by expanding  $\sum_{k=1}^K \|\bar{\mathbf{a}} - \bar{\mathbf{a}}_k\|^2$  and using the fact that  $\sum_{k=1}^K \bar{\mathbf{a}}_k = \bar{\mathbf{a}}$ .

**5.3** Consider the federated averaging algorithm for the least-mean square prob-

lem, which takes the form:

$$\begin{aligned}\phi_{k,e} &= \phi_{k,e-1} + \frac{\mu}{E_k} \frac{1}{B_k} \sum_{b=1}^{B_k} \mathbf{h}_{k,i,e,b} \left( \gamma_{k,i,e,b} - \mathbf{h}_{k,i,e,b}^\top \phi_{k,e-1} \right) \\ \psi_{k,i} &= \phi_{k,E_k} \\ \mathbf{w}_i &= \frac{1}{L} \sum_{k \in \mathcal{L}_i} \psi_{k,i}\end{aligned}$$

Determine an expression for the limiting performance  $\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2$  of the algorithm as a function of data statistics and tuning parameters.

**Solution.** The local gradient approximation takes the form of a mini-batch of size  $B_k$  on the least-mean square problem. We have encountered both mini-batching and least-mean square in Chapter 3, and can hence conclude from (3.35), (3.36) and (3.46) that:

$$\begin{aligned}\sigma_k^2 &= \frac{\sigma_{\mathbf{v}_k}^2 \text{Tr}(R_{\mathbf{h}_k})}{B_k} \\ \sigma_k^2 &= \frac{\mathbb{E} \|\mathbf{h}_k \mathbf{h}_k^\top - R_{\mathbf{h}_k}\|^2}{B_k}\end{aligned}$$

For the strong-convexity and Lipschitz constants, we can differentiate the least-mean square objective (see, e.g., (5.44)) twice to find:

$$\nabla^2 J_k(w) = R_{\mathbf{h}_k}$$

and hence:

$$\begin{aligned}\nu_k &= \lambda_{\min}(R_{\mathbf{h}_k}) \\ \delta_k &= \lambda_{\max}(R_{\mathbf{h}_k})\end{aligned}$$

We can fuse the local variances using (5.40) to find:

$$\begin{aligned}\sigma_{\text{fed}}^2 &= \frac{1}{KL} \sum_{k=1}^K \left( 2 \frac{\mathbb{E} \|\mathbf{h}_k \mathbf{h}_k^\top - R_{\mathbf{h}_k}\|^2}{B_k E_k} \|w_k^o - w^o\|^2 + \frac{\sigma_{\mathbf{v}_k}^2 \text{Tr}(R_{\mathbf{h}_k})}{B_k E_k} \right) \\ &\quad + \frac{1}{KL} \frac{K-L}{K-1} \sum_{k=1}^K 4\lambda_{\max}(R_{\mathbf{h}_k})^2 \|w_k^o - w^o\|^2\end{aligned}$$

Putting everything together, we conclude from (5.57):

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \leq \frac{4\mu\sigma_{\text{fed}}^2}{\nu} = \frac{4\mu\sigma_{\text{fed}}^2}{\lambda_{\max}(R_{\mathbf{h}_k})}$$

where the expression for  $\sigma_{\text{fed}}^2$  is given in the previous line.