

Refined Analysis of Constant Step Size Federated Averaging and Federated Richardson-Romberg Extrapolation

Paul Mangold¹, Alain Durmus¹, Aymeric Dieuleveut¹, Sergey Samsonov², Eric Moulines^{1,3}
¹CMAP, Ecole polytechnique, France ²HSE University, Russia ³MBZUAI, UAE

Federated Learning

Take N functions $f_c(\theta) = \mathbb{E}[F_c^{Z_c}(\theta)]$, with $Z_c \sim \xi_c$

Goal: solve collaboratively

$$\theta^* \in \arg \min_{\theta \in \mathbb{R}^d} f(\theta) = \frac{1}{N} \sum_{c=1}^N f_c(\theta)$$

- (a) f_c is μ -strongly convex with $\mu > 0$, F_c^z if convex
- (b) F_c^z is L -smooth and $\nabla^2 f_c(\theta) \preceq L \text{Id}$
- (c) the third derivative of f_c is uniformly bounded

Like SGD [1], FedAvg [2] is a Markov Chain

FEDAVG's iterates are a Markov chain with kernel κ :

$$\kappa(\theta_t, B) \triangleq \mathbb{E}[\mathbf{1}_B(\theta_{t+1}) | \theta_t] \text{ for } B \in \mathcal{B}(\mathbb{R}^d)$$

Theorem: FEDAVG's global iterates, with step γ and H local updates, converge to $\pi^{(\gamma, H)}$ with rate

$$\mathbf{W}_2^2(\rho \kappa^t, \pi^{(\gamma, H)}) \leq (1 - \gamma \mu)^{Ht} \mathbf{W}_2^2(\rho, \pi^{(\gamma, H)})$$

Its expectation and covariance in the limit are

$$\text{Limit point} \quad \bar{\theta}_{\text{sto}}^{(\gamma, H)} \triangleq \int \vartheta \pi^{(\gamma, H)}(d\vartheta)$$

$$\text{Covariance} \quad \bar{\Sigma}_{\text{sto}}^{(\gamma, H)} \triangleq \int \{\vartheta - \theta^*\}^{\otimes 2} \pi^{(\gamma, H)}(d\vartheta)$$

FedAvg as a Markov chain: it converges to a stationary distribution

We give an exact expression of its bias and a new method to reduce bias

$$\begin{aligned} \text{FEDAVG's bias : } \bar{\theta}_{\text{sto}}^{(\gamma, H)} - \theta^* &= \underbrace{\frac{\gamma(H-1)}{2N} \sum_{c=1}^N \nabla^2 f(\theta^*)^{-1} (\nabla^2 f_c(\theta^*) - \nabla^2 f(\theta^*)) \nabla f_c(\theta^*)}_{\text{Heterogeneity Bias}} \\ &\quad - \underbrace{\frac{\gamma}{2N} \nabla^2 f(\theta^*)^{-1} \nabla^3 f(\theta^*) \mathbf{AC}(\theta^*)}_{\text{Stochastic Bias}} + O(\gamma^2 H^2 + \gamma^{3/2} H) \end{aligned}$$

FedAvg's Linear Speed-Up

The variance of FEDAVG at stationarity is:

$$\bar{\Sigma}_{\text{sto}}^{(\gamma, H)} = \frac{\gamma}{N} \mathbf{AC}(\theta^*) + O(\gamma^2 H^2 + \gamma^{3/2} H)$$

where $\mathbf{A} \triangleq (\text{Id} \otimes \nabla^2 f(\theta^*) + \nabla^2 f(\theta^*) \otimes \text{Id})^{-1}$

$$\mathbf{C}(\theta^*) \triangleq \mathbb{E} \left[\frac{1}{N} \sum_{c=1}^N \left\{ \nabla F_c^{Z_c}(\theta^*) - \nabla f_c(\theta^*) \right\}^{\otimes 2} \right]$$

New method: Federated Richardson-Romberg

New algorithm: run FEDAVG twice with step sizes γ and 2γ and combine iterates

$$\vartheta_t^{(\gamma, H)} \triangleq 2\theta_t^{(\gamma, H)} - \theta_t^{(2\gamma, H)}$$

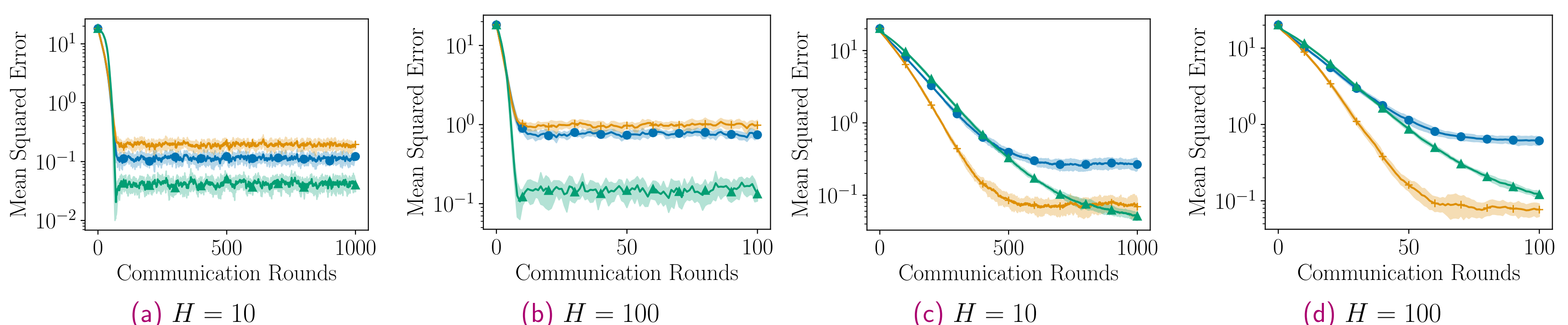
It converges to $\bar{\vartheta}_{\text{sto}}^{(\gamma, H)} = 2\bar{\theta}_{\text{sto}}^{(\gamma, H)} - \bar{\theta}_{\text{sto}}^{(2\gamma, H)}$ with bias

$$\bar{\vartheta}_{\text{sto}}^{(\gamma, H)} - \theta^* = O(\gamma^2 H^2 + \gamma^{3/2} H)$$

→ reduces **both bias** without control variates

Numerical Illustration

Logistic regression: (a-b) homogeneous data with large gradient variance, (c-d) heterogeneous data with small variance
 Blue = FedAvg [2], Green = FedAvg with Richardson-Romberg [ours], Orange = Scaffold [3]



References

- [1] A. Dieuleveut, A. Durmus, and F. Bach. "Bridging the gap between constant step size stochastic gradient descent and Markov chains". In: *The Annals of Statistics* 48.3 (2020), pp. 1348–1382.
- [2] B. McMahan et al. "Communication-efficient learning of deep networks from decentralized data". In: *Artificial intelligence and statistics*. PMLR. 2017, pp. 1273–1282.
- [3] S. P. Karimireddy et al. "Scaffold: Stochastic controlled averaging for federated learning". In: *International conference on machine learning*. PMLR. 2020, pp. 5132–5143.