

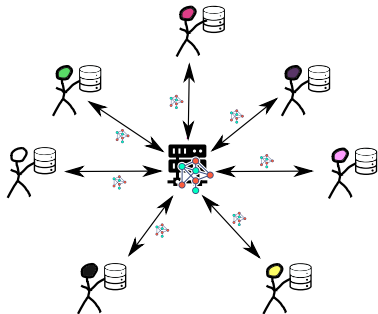
Analyse Raffinée de Federated Averaging et Extrapolation de Richardson-Romberg Fédérée

Paul Mangold (CMAP, École Polytechnique)

Journées de Statistique de la SFdS

Mercredi 4 juin 2025

Federated Learning



Collaborative optimization problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N f_c(x) \quad , \quad f_c(x) = \mathbb{E}_Z[F_c(x; Z)]$$

Central Challenges: data and computational heterogeneity
+ slow and difficult-to-establish communication

Federated Averaging

Federated Averaging¹

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N \mathbb{E}_Z[F_c(x; Z)]$$

At each global iteration

- For $c = 1$ à N in parallel

- Receive $x^{(t)}$, set $x_c^{(t,0)} = x^{(t)}$

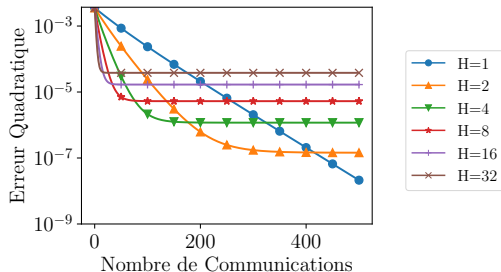
- For $h = 0$ to $H - 1$

$$x_c^{(t,h+1)} = x_c^{(t,h)} - \gamma \nabla F_c(x_c^{(t,h)}; Z_c^{(t,h+1)})$$

- Aggregate local models

$$x^{(t+1)} = \frac{1}{N} \sum_{c=1}^N x_c^{(t,H)}$$

With deterministic gradients:



¹B. McMahan et al. "Communication-efficient learning of deep networks from decentralized data". In: **AISTATS**. 2017.

Classical analyses of this algorithm

(For L -smooth, μ -strongly convex functions)

Choose your favorite heterogeneity measure

¹X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS (2017)*.

²A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv preprint arXiv:2209.02257 (2022)*.

³J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR 2024 (2024)*.

Classical analyses of this algorithm

(For L -smooth, μ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order¹: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$

¹X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS (2017)*.

²A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv preprint arXiv:2209.02257 (2022)*.

³J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR 2024 (2024)*.

Classical analyses of this algorithm

(For L -smooth, μ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order¹: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$
- second-order²: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla_c^2 f(x^*) - \nabla^2 f(x^*)\|^2$

¹X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS (2017)*.

²A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv preprint arXiv:2209.02257 (2022)*.

³J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR 2024 (2024)*.

Classical analyses of this algorithm

(For L -smooth, μ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order¹: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$
- second-order²: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla_c^2 f(x^*) - \nabla^2 f(x^*)\|^2$
- average drift³: $\zeta = \left\| \frac{1}{NH} \sum_{c=1}^N \sum_{h=0}^{H-1} \nabla f(x_c^{(h)}) - \nabla f(x^*) \right\|^2$

¹X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS* (2017).

²A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv preprint arXiv:2209.02257* (2022).

³J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR* 2024 (2024).

Classical analyses of this algorithm

(For L -smooth, μ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order¹: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$
- second-order²: $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla_c^2 f(x^*) - \nabla^2 f(x^*)\|^2$
- average drift³: $\zeta = \left\| \frac{1}{NH} \sum_{c=1}^N \sum_{h=0}^{H-1} \nabla f(x_c^{(h)}) - \nabla f(x^*) \right\|^2$

Show **convergence to a neighborhood** of x^*

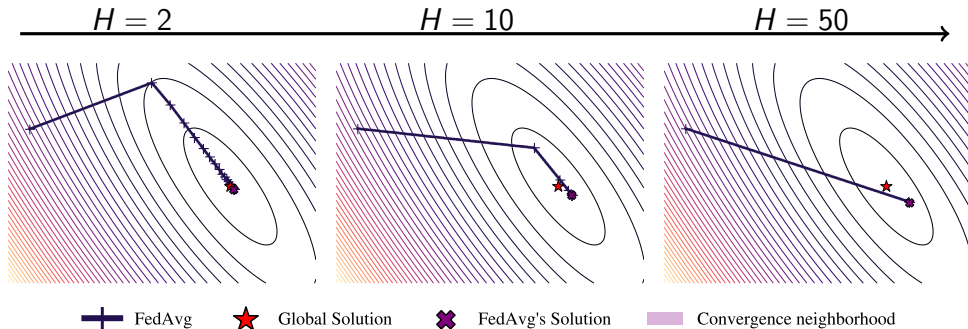
$$\|x^{(T)} - x^*\|^2 \lesssim (1 - \gamma\mu)^{HT} \|x^{(0)} - x^*\|^2 + \chi(\gamma, H, \zeta) \quad (\text{for some function } \chi)$$

¹X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS* (2017).

²A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv preprint arXiv:2209.02257* (2022).

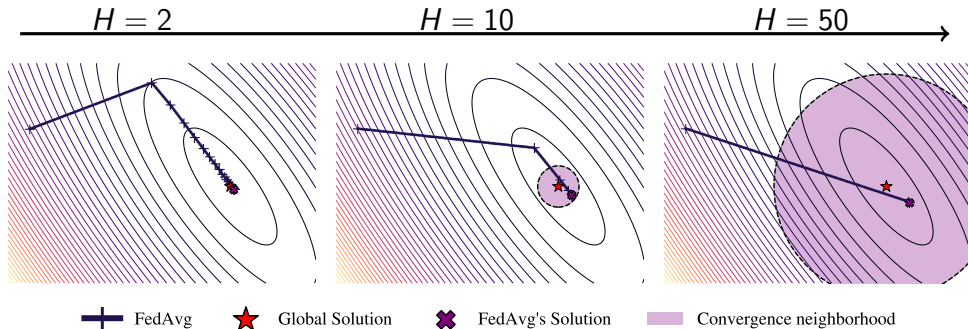
³J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR* 2024 (2024).

$$\|x^{(T)} - x^*\|^2 \lesssim (1 - \gamma\mu)^{HT} \|x^{(0)} - x^*\|^2 + \chi(\gamma, H, \zeta) \quad (\text{for some function } \chi)$$



When the number of local iterations increases, bias increases

$$\|x^{(T)} - x^*\|^2 \lesssim (1 - \gamma\mu)^{HT} \|x^{(0)} - x^*\|^2 + \chi(\gamma, H, \zeta) \quad (\text{for some function } \chi)$$



When the number of local iterations increases, bias increases

Remark: It seems that iterates converge in some way?

Federated Averaging as Fixed Point Iteration

Remark that, starting with $x_c^{(t)}, y_c^{(t)} \in \mathbb{R}^d$,

$$x_c^{(t,h+1)} - y_c^{(t,h+1)} = x_c^{(t,h)} - y_c^{(t,h)} - \gamma(\nabla f_c(x_c^{(t,h)}) - \nabla f_c(y_c^{(t,h)}))$$

Thus

$$\|x_c^{(t+1)} - y_c^{(t+1)}\| \leq (1 - \gamma\mu)^H \|x_c^{(t)} - y_c^{(t)}\|$$

¹G. Malinovskiy et al. "From local SGD to local fixed-point methods for federated learning". In: ICML. 2020.

Federated Averaging as Fixed Point Iteration

Remark that, starting with $x_c^{(t)}, y_c^{(t)} \in \mathbb{R}^d$,

$$x_c^{(t,h+1)} - y_c^{(t,h+1)} = x_c^{(t,h)} - y_c^{(t,h)} - \gamma(\nabla f_c(x_c^{(t,h)}) - \nabla f_c(y_c^{(t,h)}))$$

Thus

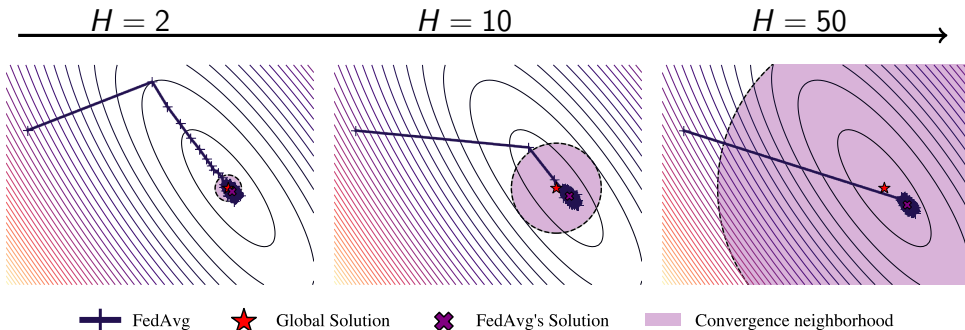
$$\|x_c^{(t+1)} - y_c^{(t+1)}\| \leq (1 - \gamma\mu)^H \|x_c^{(t)} - y_c^{(t)}\|$$

\Rightarrow deterministic FedAvg converges to a unique point¹

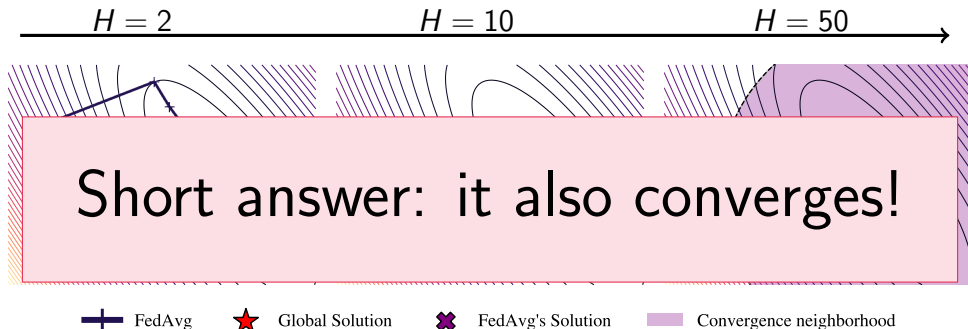
¹G. Malinovskiy et al. "From local SGD to local fixed-point methods for federated learning". In: ICML. 2020.

Open Question: What about the Stochastic Case?

Open Question: What about the Stochastic Case?



Open Question: What about the Stochastic Case?



FedAvg (with stochastic gradients) converges!¹

(For thrice derivable, L -smooth, μ -strongly convex functions)

- FedAvg converges to a stationary distribution $\pi^{(\gamma, H)}$
 - denoting $x^{(t)} \sim \psi_{x^{(t)}}$, we have

$$\mathcal{W}_2(\psi_{x^{(t)}}; \pi^{(\gamma, H)}) \leq (1 - \gamma\mu)^{Ht} \mathcal{W}_2(\psi_{x^{(0)}}; \pi^{(\gamma, H)})$$

- where \mathcal{W}_2 is the second order Wasserstein distance

¹P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

FedAvg (with stochastic gradients) converges!¹

(For thrice derivable, L -smooth, μ -strongly convex functions)

- FedAvg converges to a stationary distribution $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \boxed{\frac{\gamma}{N} C(x^*)} + O(\gamma^{3/2} H)$$

¹P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

FedAvg (with μ and σ) converges! ¹

(F

Linear speed-up !

variance decreases in $1/N$

$C(x^*)$ is ∇F^Z 's covariance at x^*

ngly convex functions)

- FedAvg converges to a stationary distribution $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \frac{\gamma}{N} C(x^*) + O(\gamma^{3/2} H)$$

¹P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

FedAvg (with stochastic gradients) converges!¹

(For thrice derivable, L -smooth, μ -strongly convex functions)

- FedAvg converges to a stationary distribution $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is
- We can now give an **exact expansion of the bias**

$$\int x \pi^{(\gamma, H)}(dx) = x^* + \frac{\gamma(H-1)}{2N} \sum_{c=1}^N \nabla^2 f(x^*)^{-1} (\nabla^2 f_c(x^*) - \nabla^2 f(x^*)) \nabla f_c(x^*) \\ - \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2} H)$$

¹P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

Heterogeneity bias

vanishes when $\nabla^2 f_c(x^*) = \nabla^2 f(x^*)$
or when $\nabla f_c(x^*) = \nabla f(x^*)$
or when $H = 1$ (one local update)

Stochasticity bias

$A = I \otimes \nabla^2 f(x^*) + \nabla^2 f(x^*) \otimes I$
 $C(x^*)$ is ∇F^Z 's covariance at x^*

- FedAvg's iterates covariance is
- We can now give an **exact expansion of the bias**

$$\int x \pi^{(\gamma, H)}(dx) = x^* + \frac{\gamma(H-1)}{2N} \sum_{c=1}^N \nabla^2 f(x^*)^{-1} (\nabla^2 f_c(x^*) - \nabla^2 f(x^*)) \nabla f_c(x^*) \\ - \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2} H)$$

¹P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

Correcting the Bias

Novel Algorithm: Federated Richardson-Romberg Extrapolation

Run FedAvg twice:

- with step size γ : global iterates $x_\gamma^{(t)}$
- with step size 2γ : global iterates $x_{2\gamma}^{(t)}$

We can combine the iterates

$$\chi_{\text{RR}}^{(t)} = 2x_\gamma^{(t)} - x_{2\gamma}^{(t)}$$

Correcting the Bias

Novel Algorithm: Federated Richardson-Romberg Extrapolation

Run FedAvg twice:

- with s
- with s

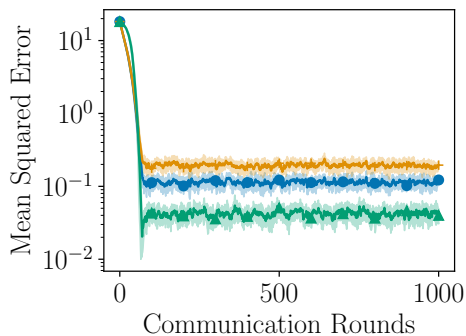
$$\text{Theorem: } \mathbb{E}[\chi_{\text{RR}}^{(t)}] = x_{\star} + O(\gamma^2 H^2 + \gamma^{3/2} H)$$

→ bias is effectively reduced!!

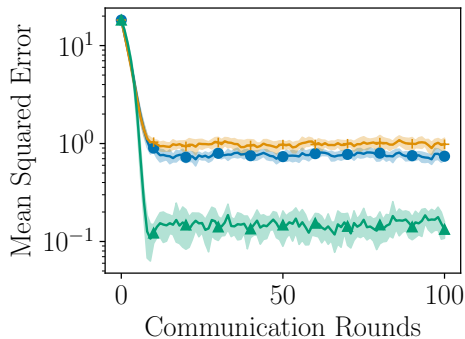
We can co

$$\chi_{\text{RR}}^{(t)} = 2x_{\gamma}^{(t)} - x_{2\gamma}^{(t)}$$

Numerical Illustration: FedAvg



(a) $H = 10$



(b) $H = 100$

Blue: FedAvg, Orange: Scaffold, Green: Federated Richardson-Romberg

Conclusion

- FedAvg converges (even with stochastic gradients)
- This allows to derive new analyses for these problems, with exact first-order expression for bias
- FedAvg (and friends) are still biased!
 - there is still a lot to do in federated optimization
 - ... especially when gradients are stochastic,
... which is the most interesting setting!
- Similar results hold for Scaffold (see next slide)

PS: Extension to Scaffold

Little teaser :)

Similar results hold for Scaffold, and we prove that:

- Scaffold's iterates converge
- Scaffold eliminates heterogeneity bias *but not stochasticity bias*
- new convergence rate for Scaffold:

$$\mathbb{E} [\|x^{(T)} - x^*\|^2] \lesssim \left(1 - \frac{\gamma\mu}{4}\right)^{HT} \|x^{(0)} - x^*\|^2 + O\left(\frac{\gamma}{N} + \gamma^{3/2}\right)$$

Thank you!

Papers related to this presentation:

- P. Mangold et al. “Refined Analysis of Federated Averaging’s Bias and Federated Richardson-Romberg Extrapolation”. In: **AISTATS**. 2025
- P. Mangold et al. “Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up”. In: **ICML**. 2025