

The Problem and Assumptions

We want to solve the finite-sum optimization problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}. \quad (1)$$

- Problem (1) has many applications in machine learning, data science and engineering;
- We focus on the regime when n is very large. This is typically the case in big data settings (e.g., massively distributed and federated learning).

Assumptions:

- f is μ -strongly convex;
- all f_i have Lipschitz continuous Hessians with respect to spectral (L_*), Frobenius (L_F) and infinity (L_∞) norms;
- x^* is the solution for Problem (1).

Main goal

Our goal is to develop a communication efficient Newton-type method for federated learning.

Newton's method

Newton's step: $x^{k+1} = x^k - (\nabla^2 f(x^k))^{-1} \nabla f(x^k)$.

- Pros:**
- Fast **local quadratic** convergence rate
 - Rate is **independent on the condition number**
- Cons:**
- Requires $\mathcal{O}(d^2)$ floats to be communicated by each worker to the server, where d is typically very large

Newton Star

Newton Star step: $x^{k+1} = x^k - (\nabla^2 f(x^*))^{-1} \nabla f(x^k)$.

- Pros:**
- Fast **local quadratic** convergence rate
 - Rate is **independent on the condition number**
 - Requires $\mathcal{O}(d)$ floats to be communicated by each worker;
- Cons:**
- Cannot be implemented in practice.

Newton Zero

Newton Zero step: $x^{k+1} = x^k - (\nabla^2 f(x^0))^{-1} \nabla f(x^k)$.

- Pros:**
- Fast **local linear** convergence rate
 - Rate is **independent on the condition number**
 - Requires $\mathcal{O}(d)$ floats to be communicated by each worker;

Newton Triangle

FedNL and its four extensions interpolates between these three special Newton-type method — Newton (N), Newton Star (NS) and Newton Zero (NO).

FedNL

How to address the communication bottleneck?

- Compressed communication

In FedNL we maintain a sequence of matrices $\mathbf{H}_i^k \in \mathbb{R}^{d \times d}$, for all $i = 1, \dots, n$ throughout the iterations $k \geq 0$, with the goal of learning $\mathbf{H}_i(x^*)$ for all i :

$$\mathbf{H}_i^k \rightarrow \nabla^2 f_i(x^*) \quad \text{as } k \rightarrow +\infty.$$

Using $\mathbf{H}_i^k \approx \nabla^2 f_i(x^*)$, we can estimate the Hessian $\nabla^2 f(x^*)$ via

$$\nabla^2 f(x^*) \approx \mathbf{H}^k := \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^k.$$

FedNL: Extensions

- **FedNL-PP:** FedNL with partial participation;
- **FedNL-LS:** FedNL with globalization via Line Search;
- **FedNL-CR:** FedNL with globalization via Cubic Regularization [2];
- **FedNL-BC:** FedNL with Bidirectional Compression.

Compression operators

Unbiased Compressors. By $\mathbb{B}(\omega)$ we denote the class of (possibly randomized) unbiased compression operators $\mathcal{C} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$ with variance parameter $\omega \geq 0$ satisfying

$$\mathbb{E} \mathcal{C}(M) = M, \quad \mathbb{E} \|\mathcal{C}(M) - M\|_F^2 \leq \omega \|M\|_F^2, \quad M \in \mathbb{R}^{d \times d}.$$

Contractive Compressors. By $\mathbb{C}(\delta)$ we denote the class of deterministic contractive compression operators $\mathcal{C} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$ with contraction parameter $\delta \in [0, 1]$ satisfying

$$\|\mathcal{C}(M)\|_F \leq \|M\|_F, \quad \|\mathcal{C}(M) - M\|_F^2 \leq (1 - \delta) \|M\|_F^2, \quad \forall M \in \mathbb{R}^{d \times d}.$$

Learning mechanism

Learning the matrices: the idea

We design a learning rule for matrices \mathbf{H}_i^k via the **DIANA trick** [1]:

$$\mathbf{H}_i^{k+1} = \mathbf{H}_i^k + \alpha \mathcal{C}_i^k (\nabla^2 f_i(x^k) - \mathbf{H}_i^k),$$

where $\alpha > 0$ is a learning rate, and \mathcal{C}_i^k is a freshly sampled compressor by node i at iteration k .

Main features of a family of FedNL methods important for Federated Learning

- supports **heterogeneous data** setting
- uses **adaptive stepsizes**
- supports **unbiased Hessian compression** (e.g., Rand- K)
- fast local rate: **independent of the condition number**
- has global convergence guarantees via **line search**
- supports smart **uplink gradient compression** at the devices
- applies to general **finite-sum problems**
- **privacy is enhanced** (training data is not sent to the server)
- supports **contractive Hessian compression** (e.g., Top- K)
- supports **partial participation**
- has global convergence guarantees via **cubic regularization**
- supports smart **downlink model compression** by the server

Table: Convergence results for a family of FedNL methods.

Method	Convergence			Rate independent on the condition number
	result [†]	type	rate	
NO	$r_k \leq \frac{1}{2^k} r_0$	local	linear	✓
NS	$r_{k+1} \leq c r_k^2$	local	quadratic	✓
FedNL	$r_k \leq \frac{1}{2^k} r_0$	local	linear	✓
	$\Phi_1^k \leq \theta^k \Phi_1^0$	local	linear	✓
FedNL-PP	$r_{k+1} \leq c \theta^k r_k$	local	superlinear	✓
	$\mathcal{W}^k \leq \theta^k \mathcal{W}^0$	local	linear	✓
FedNL-LS	$\Phi_2^k \leq \theta^k \Phi_2^0$	local	linear	✓
	$r_{k+1} \leq c \theta^k r_k$	local	linear	✓
FedNL-LS	$\Delta_k \leq \theta^k \Delta_0$	global	linear	✗
	$\Delta_k \leq c/k$	global	sublinear	✗
FedNL-CR	$\Delta_k \leq \theta^k \Delta_0$	global	linear	✗
	$\Phi_1^k \leq \theta^k \Phi_1^0$	local	linear	✓
FedNL-BC	$r_{k+1} \leq c \theta^k r_k$	local	superlinear	✓
	$\Phi_3^k \leq \theta^k \Phi_3^0$	local	linear	✓

[†] Refer to the precise statements of the theorems in [3] for the exact values.

Experiments

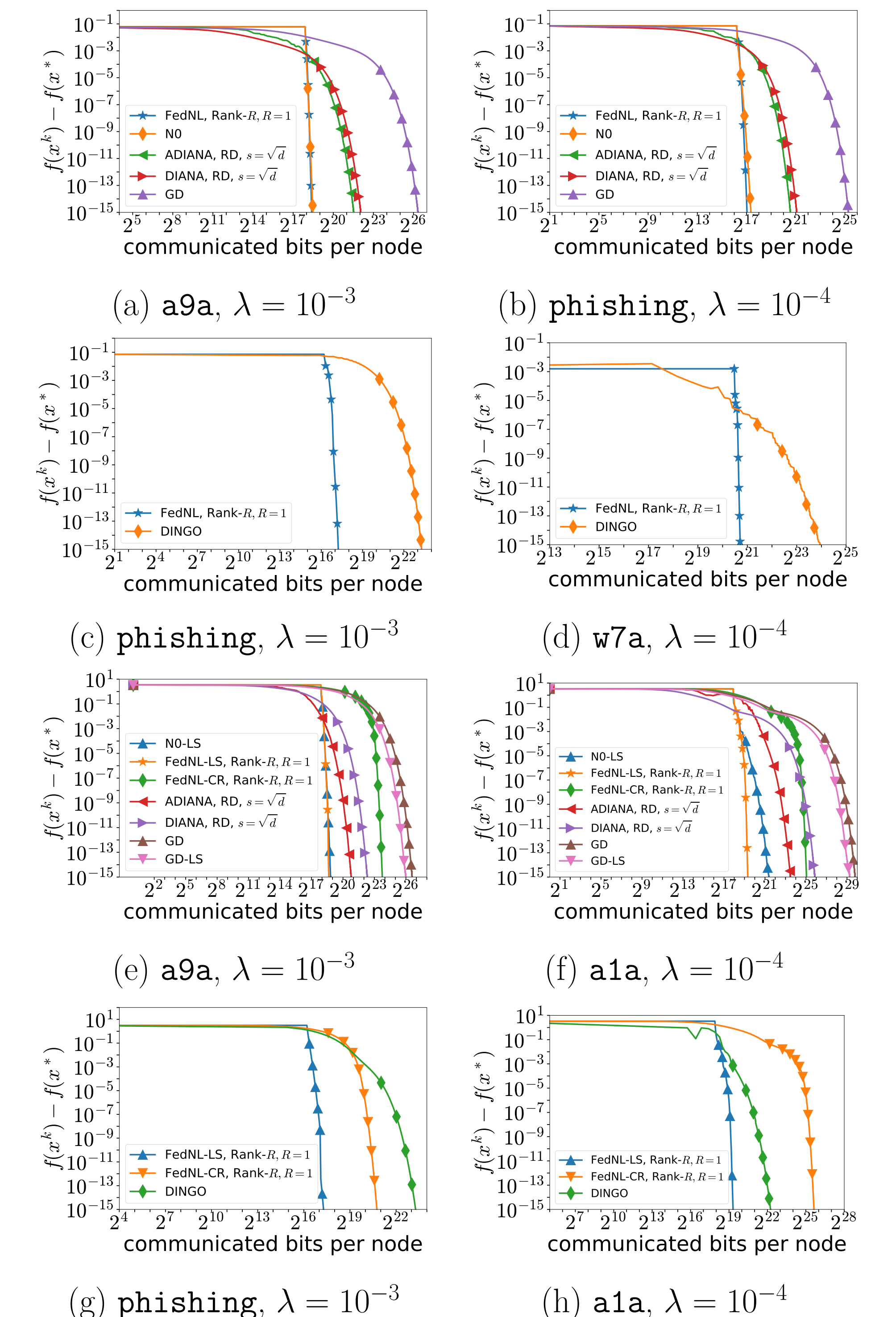


Figure 1: **First row:** Local comparison of FedNL and N0 with ADIANA, DIANA, and GD in terms of communication complexity. **Second row:** Local comparison of FedNL with DINGO (second row) in terms of communication complexity. **Third row:** Global comparison of FedNL-LS, N0-LS, and FedNL-CR with ADIANA, DIANA, GD, and GD-LS in terms of communication complexity. **Fourth row:** Global comparison of FedNL-LS and FedNL-CR with DINGO in terms of communication complexity.

References

- [1] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- [2] Yurii Nesterov and Boris T. Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1) : 177 – 205, 2006.
- [3] Mher Safaryan, Rustem Islamov, Xun Qian, and Peter Richtárik. FedNL: Making Newton-Type Methods Applicable to Federated Learning. *arXiv preprint arXiv:2106.02969*, 2021.