

FedNS: A Fast Sketching Newton-Type Algorithm for Federated Learning

Jian Li¹, Yong Liu^{2,*} and Weiping Wang¹

¹ Institute of Information Engineering, Chinese Academy of Sciences

² Gaoling School of Artificial Intelligence, Renmin University of China



中国科学院信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING, CAS

Backgrounds

Federated algorithms	Communication Rounds	Communication complexity
First-order methods (e.g. FedAvg & FedProx)	$\mathcal{O}(1/t)$	$\mathcal{O}(M)$
Distributed Newton's Methods	$\mathcal{O}(\log 1/t)$	$\mathcal{O}(M^2)$

Note: The computational complexities are computed in terms of regularized least squared loss to obtain a δ -accurate solution, i.e., $L(\mathbf{w}_t) - L(\mathbf{w}_{\mathcal{D},\lambda}) \leq \delta$. M is the number of model parameters.

Existing federated learning algorithms:

- First-order methods: Low communication burdens but slow convergence.
- Second-order methods: Fast convergence but high communication burdens.

Motivation: How to take advantage of both first-order and second-order algorithms, such that the federated algorithms can achieve low communication burdens and fast convergence.

Contributions: We propose a federated Newton sketch method, named FedNS, which shares both first-order and second-order information across devices, i.e., local gradients and sketched square-root Hessian.

1) **On the algorithmic front.** FedNS: a simple federated method with Newton sketch. FedNDES: a communication-efficient Newton-type federated algorithm with line-search.

2) **On the statistical front.** Convergence analysis for two federated Newton sketching algorithms: FedNS and its line-search variant FedNDES. These methodologies delineate not only **super-linear convergence** rates but also entail a **small communication complexity**.

Preliminaries

$$\mathbf{w}_{\mathcal{D},\lambda} = \arg \min_{\mathbf{w} \in \mathcal{H}_K} \underbrace{\frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i)}_{L(\mathcal{D}, \mathbf{w})} + \lambda \alpha(\mathbf{w}),$$

Centralized Newton's Method

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mu \mathcal{H}_{\mathcal{D},t}^{-1} \mathbf{g}_{\mathcal{D},t}, \quad \text{with}$$

$$\mathbf{g}_{\mathcal{D},t} := \nabla L(\mathcal{D}, \mathbf{w}_t) + \lambda \nabla \alpha(\mathbf{w}_t),$$

$$\mathcal{H}_{\mathcal{D},t} := \nabla^2 L(\mathcal{D}, \mathbf{w}_t) + \lambda \nabla^2 \alpha(\mathbf{w}_t).$$

Newton's Method with Partial Sketching

$$\nabla^2 L(\mathcal{D}, \mathbf{w}_t) =$$

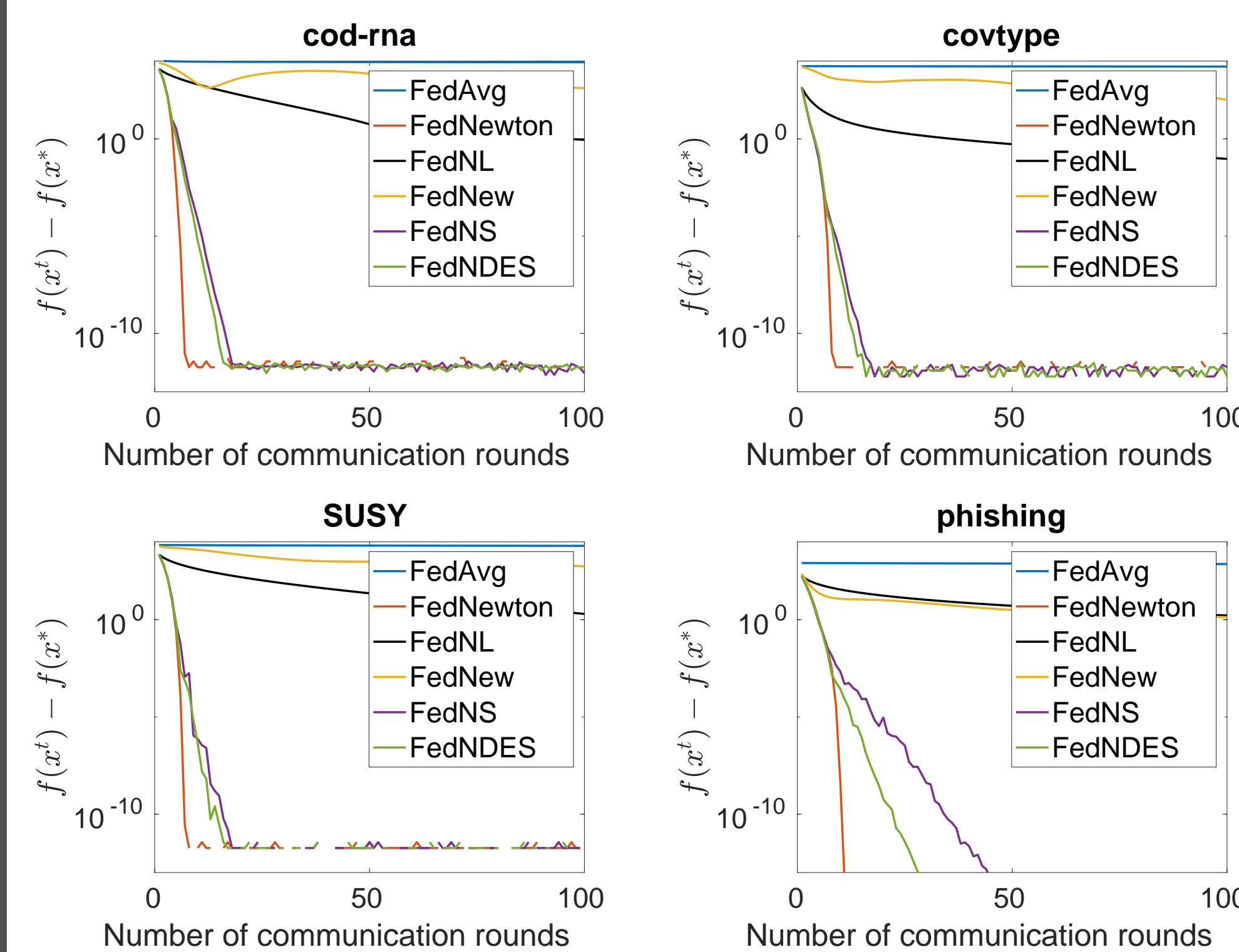
$$(\mathbf{S} \nabla^2 L(\mathcal{D}, \mathbf{w}_t)^{1/2})^\top (\mathbf{S} \nabla^2 L(\mathcal{D}, \mathbf{w}_t)^{1/2}).$$

Federated Newton's Method (FedNewton)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mu \mathcal{H}_{\mathcal{D},t}^{-1} \mathbf{g}_{\mathcal{D},t} \quad \text{with}$$

$$\mathcal{H}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \mathcal{H}_{\mathcal{D}_j,t}, \quad \mathbf{g}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \mathbf{g}_{\mathcal{D}_j,t}$$

Experiments



There is significant gaps between the convergence speeds of our proposed methods FedNS, FedNDES and the existing Newton-type FL methods.

The proposed FedNS and FedNDES converge nearly to FedNewton, while FedNew and FedNL converge slowly closed to FedAvg.

FedNDES converges faster than FedNS and the final predictive accuracies of FedNDES are higher.

FedNS

Algorithm 1: Federated Learning with Newton Sketch (FedNS)

Input: Feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^M$, start point \mathbf{w}_0 , the termination iterations T and the step-size μ .

Output: The global estimator \mathbf{w}_T .

- Local machines:** Compute the local feature mapping data matrix $\phi(\mathbf{X}_j)$.
- for** $t = 1$ to T **do**
- Local machines:** Sample the sketch matrix $\mathbf{S}_j^t \in \mathbb{R}^{k \times n_j}$ from the SRHT. Compute local sketch Hessian matrices $\Upsilon_{\mathcal{D}_j,\lambda} = \mathbf{S}_j^t \nabla^2 L(\mathcal{D}_j, \mathbf{w}_t)^{1/2}$ and local gradients $\mathbf{g}_{\mathcal{D}_j,t}$. Upload them to the global server (\uparrow).
- Global server:** Compute the global Hessian matrix $\tilde{\mathbf{H}}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \Upsilon_{\mathcal{D}_j,\lambda}^\top \Upsilon_{\mathcal{D}_j,\lambda} + \lambda \nabla^2 \alpha(\mathbf{w}_t)$ and the global gradient $\mathbf{g}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \mathbf{g}_{\mathcal{D}_j,t}$ and update the global estimator

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \mu \tilde{\mathbf{H}}_{\mathcal{D},t}^{-1} \mathbf{g}_{\mathcal{D},t}$$

and communicate it to local machines (\downarrow).

end for

Local sketch Hessian

$$\Upsilon_{\mathcal{D}_j,\lambda} = \mathbf{S}_j \nabla^2 L(\mathcal{D}_j, \mathbf{w}_t)^{1/2}.$$

Global sketch Hessian

$$\tilde{\mathbf{H}}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \Upsilon_{\mathcal{D}_j,\lambda}^\top \Upsilon_{\mathcal{D}_j,\lambda} + \lambda \nabla^2 \alpha(\mathbf{w}_t).$$

FedNDES

Algorithm 2: Dimension-efficient federated Newton (FedNDES)

Input: Feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^M$, start point \mathbf{w}_0 , accuracy tolerance $\delta > 0$, line-search parameters (a, b) , threshold sketch sizes \bar{m}_1 and \bar{m}_2 , and the decrement parameter η .

Output: The global estimator \mathbf{w}_T .

- Local machines:** Compute the local feature mapping data matrix $\phi(\mathbf{X}_j)$. Initialize and $k_t = \bar{m}_1$.
- for** $t = 1, \dots, T$ **do**
- Local machines:** Sample the sketch matrix $\mathbf{S}_j^t \in \mathbb{R}^{k \times n_j}$ from the SRHT. Compute local sketch square root Hessian $\Upsilon_{\mathcal{D}_j,\lambda} = \mathbf{S}_j^t \nabla^2 L(\mathcal{D}_j, \mathbf{w}_t)^{1/2}$ and local gradients $\mathbf{g}_{\mathcal{D}_j,t}$. Upload them to the global server (\uparrow).
- Global server:** Compute the global Hessian matrix $\tilde{\mathbf{H}}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \Upsilon_{\mathcal{D}_j,\lambda}^\top \Upsilon_{\mathcal{D}_j,\lambda} + \lambda \nabla^2 \alpha(\mathbf{w}_t)$, the global gradient $\mathbf{g}_{\mathcal{D},t} = \sum_{j=1}^m \frac{n_j}{N} \mathbf{g}_{\mathcal{D}_j,t}$ and the approximate Newton step $\Delta \mathbf{w}_t = -\tilde{\mathbf{H}}_{\mathcal{D},t}^{-1} \mathbf{g}_{\mathcal{D},t}$. Compute the approximate Newton decrement

$$\tilde{\lambda}(\mathbf{w}_t) = \mathbf{g}_{\mathcal{D},t}^\top \Delta \mathbf{w}_t.$$

If $\tilde{\lambda}(\mathbf{w}_t)^2 \leq \frac{3}{4} \delta$ return the model \mathbf{w}_t . Otherwise send $\Delta \mathbf{w}_t$ and $\tilde{\lambda}(\mathbf{w}_t)$ to local workers.

- Local machines:** Line search from $\mu_j = 1$: **while** $L(\mathcal{D}_j, \mathbf{w}_t + \mu_j \Delta \mathbf{w}_t) > L(\mathcal{D}_j, \mathbf{w}_t) + a \mu_j \tilde{\lambda}(\mathbf{w}_t)$, **then** $\mu_j \leftarrow b \mu_j$. Send μ_j to the global server.
- Global server:** Let $\mu = \min_{j \in \mathcal{M}} \mu_j$. Update the global estimator

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \mu \Delta \mathbf{w}_t.$$

If $\tilde{\lambda}(\mathbf{w}_t) > \eta$, set $k = \bar{m}_1$. Otherwise, set $k = \bar{m}_2$. Communicate the global model \mathbf{w}_t and the sketch size k to local machines (\downarrow).

end for

Comparison

Related Work	Heterogeneous setting	Sketch size k	Iterations t	Communication per round	Total communication complexity
FedAvg (Li et al. 2020c; Su, Xu, and Yang 2021)	✓	–	$\mathcal{O}(\frac{1}{\delta})$	$\mathcal{O}(M)$	$\mathcal{O}(\frac{M}{\delta})$
FedProx (Li et al. 2020a; Su, Xu, and Yang 2021)	✓	–	$\mathcal{O}(\frac{1}{\delta})$	$\mathcal{O}(M)$	$\mathcal{O}(\frac{M}{\delta})$
DistributedNewton (Ghosh, Maity, and Mazumdar 2020)	×	–	$\mathcal{O}(\log \frac{1}{\delta})$	$\mathcal{O}(M)$	$\mathcal{O}(M \log \frac{1}{\delta})$
LocalNewton (Gupta et al. 2021)	×	–	$\mathcal{O}(\log \frac{1}{\delta})$	$\mathcal{O}(M)$	$\mathcal{O}(M \log \frac{1}{\delta})$
FedNL (Safaryan et al. 2022)	✓	–	$\mathcal{O}(\log \frac{1}{\delta})$	$\mathcal{O}(M)$	$\mathcal{O}(M \log \frac{1}{\delta})$
SHED (Fabbro et al. 2022)	✓	–	$\mathcal{O}(\log \frac{1}{\delta})$	–	$\mathcal{O}(M^2)$
FedNewton	✓	–	$\mathcal{O}(\log \log \frac{1}{\delta})$	$\mathcal{O}(M^2)$	$\mathcal{O}(M^2 \log \log \frac{1}{\delta})$
FedNS (Algorithm 1)	✓	M	$\mathcal{O}(\log \log \frac{1}{\delta})$	$\mathcal{O}(kM)$	$\mathcal{O}(kM \log \log \frac{1}{\delta})$
FedNDES (Algorithm 2)	✓	d_λ	$\mathcal{O}(\log \log \frac{1}{\delta})$	$\mathcal{O}(kM)$	$\mathcal{O}(kM \log \log \frac{1}{\delta})$

Note: The computational complexities are computed in terms of regularized least squared loss to obtain a δ -accurate solution, i.e., $L(\mathbf{w}_t) - L(\mathbf{w}_{\mathcal{D},\lambda}) \leq \delta$. The convergence analysis for FedAvg is provided in (Li et al. 2020c; Su, Xu, and Yang 2021) and that for FedProx is provided in (Li et al. 2020a; Su, Xu, and Yang 2021).

The proposed algorithms achieve faster convergence with relatively small communication burdens.