

On the Convergence of Differentially Private Federated Learning on Non-Lipschitz Objectives via Clipping and Normalized Client Updates

Abolfazl Hashemi, Purdue ECE
FLOW Seminar, April 20th, 2022

<https://arxiv.org/abs/2106.07094>



Rudrajit Das (UT CS)



Sujay Sanghavi (UT ECE)



Inderjit S. Dhillon (UT CS)

Collaborative ML via Federated Learning (FL)

- Decentralized data (as opposed to traditional distributed learning)
- Different (resp., identical) D_i 's: **heterogeneous** (resp., **homogeneous**) setting.

- Despite the locality of data storage in FL, information-sharing opens the door to the possibility of sabotaging the security of personal data through communication.
- Can the server **optimize** while preserving a **strong notion of privacy** of clients' data?

Differential Privacy (DP)

- DP is a popular privacy-quantifying framework for training of ML models.
- Goal: Learning nothing about an individual while learning useful information about a whole population
- **Reducing** learning algorithm's **sensitivity** to an individual's data
- To protect the individuals' privacy, one **adds a controlled amount of random noise** to the results of our analysis.

Differential Privacy (DP) Background

Neighboring datasets

Two datasets $D \in \mathcal{D}_c$ and $D^0 \in \mathcal{D}_c$ are said to be neighboring if they differ in exactly one sample, and we denote this by $d_H(D, D^0) = 1$.

$(\epsilon; \delta)$ -DP [DMNS06]

Given a collection of datasets \mathcal{D}_c and a query function $h : \mathcal{D}_c \rightarrow \mathcal{X}$, a randomized mechanism $M : \mathcal{X} \rightarrow \mathcal{Y}$ is said to be $(\epsilon; \delta)$ -DP, if for any two neighboring datasets

$$P(M(h(D)) \in R) \leq e^\epsilon P(M(h(D^0)) \in R) + \delta$$

- When $\delta = 0$, it is commonly known as pure DP. Otherwise, it is known as approximate DP.
- Setting M to **additive random Gaussian noise** – known as the Gaussian mechanism – is a customary approach to provide DP.

How Much Noise is Adequate?

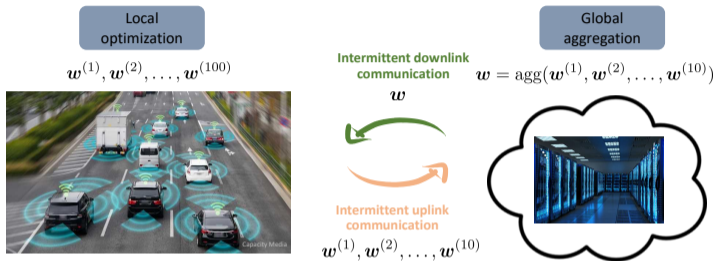
Gaussian mechanism [DR+14]

Let $\epsilon := \sup_{D, D'} \|h(D) - h(D')\|_p$. If we set $M(h(D)) = h(D) + Z$; where $Z \sim \mathcal{N}(0, \frac{2 \log(1/\delta)}{\epsilon^2})$, then the mechanism M is (ϵ, δ) -DP.

- Noise power increases with sensitivity and desired privacy guarantees.
- Similar results exist for Laplace mechanism and discretized/truncated distributions

Going back to FL: How to apply DP in FL?

DP in FL: FedAvg Revisited



- Local update: For E steps, do GD, i.e., $w_k^{(i)} \leftarrow w_k^{(i)} - \eta \nabla f_i(w_k^{(i)})$
- Each client communicates its update $u_k^{(i)} = \frac{w_k^{(i)} - w_k}{E}$ to server w.p. $\frac{r}{n}$ (total of K rounds)
- Global update: $w_{k+1} = w_k + \frac{\eta}{r} \sum_{i \in S_k} u_k^{(i)}$
- Output: w_k with $k \sim \text{Unif}[0; K - 1]$.

$u_k^{(i)}$ contains local gradient information ! needs to be made private!

DP-FedAvg with Clipping

- Maximum sensitivity, ϵ , grows with norm of $u_k^{(i)} = \frac{w_k w_{k;E}^{(i)}}{k}$
- Assuming G -Lipschitzness, i.e., $\sup_{\theta_2} \|r g(\cdot)\|_2 \leq G$, this norm is at most GE !
controlled additive noise
- In general, need to limit how large $u_k^{(i)}$ can get to **remove unbounded impact of one client.**

Clipped Updates: $u_k^{(i)} \min \left(1; \frac{C}{k u_k^{(i)}} \right) + \frac{u_k^{(i)}}{k}; \quad \frac{u_k^{(i)}}{k} \sim N(0_d; r^2 I_d)$

Theorem (Based on [ACG+16])

There exists an absolute constant $q > 0$ s.t. for $\epsilon = O(1)$, DP-FedAvg will be $(\epsilon; \delta)$ -DP as long as

$$\epsilon^2 = qKC^2 \frac{\log(1/\delta)}{n^2 \epsilon^2}.$$

Convergence of DP-FedAvg with Clipping

Relevant Notations and Definitions

Convexity

A function $g : \mathbb{R} \rightarrow \mathbb{R}$ is convex if $g(\alpha x + (1 - \alpha)y) \leq \alpha g(x) + (1 - \alpha)g(y)$ for any $x, y \in \mathbb{R}$ and $0 \leq \alpha \leq 1$.

Smoothness

A function $g : \mathbb{R} \rightarrow \mathbb{R}$ is said to be L -smooth if for all $x, y \in \mathbb{R}$, $\|g(x) - g(y) - L(x - y)\| \leq \frac{L}{2}\|x - y\|^2$. If g is twice differentiable, then for all $x, y \in \mathbb{R}$:

$$g(y) = g(x) + h \nabla g(x) \cdot (y - x) + \frac{L}{2} \|y - x\|^2$$

Heterogeneity

Let $w^* \in \arg \min_{w \in \mathbb{R}^d} f(w)$ and $\mu_i := f_i(w^*) - \min_{w \in \mathbb{R}^d} f_i(w) \geq 0$. Then the heterogeneity of the system is quantified by some increasing function of the μ_i 's.

Convergence Without Assuming Lipschitzness

Theorem 1

Suppose the f_i 's are convex and L -smooth over \mathbb{R}^d . Define $0 < \epsilon := \frac{\rho \sqrt{qd \log(1/\epsilon)}}{n^\alpha} < 1$. There exists constant local and global learning rates η and η_g , and a lower bound on the clipping threshold C_{low} , such that for any $C \geq C_{\text{low}}$ and in $K = O\left(\frac{Lkw_0}{C^2} \frac{w}{k}\right)$ rounds

$$E \frac{1}{n} \sum_{i=1}^n \min_{w_k} f_i(w_k) - f_i(w^*) \leq O\left(\frac{C}{LE} kr f_i(w_k)k\right) \leq O\left(\frac{C}{E} kw_0\right) + E \frac{1}{n} \sum_{i=1}^n \dots$$

- ϵ : the price of privacy – increases as the level of privacy increases (i.e., α and ρ decrease).
- Non-vanishing convergence error in Private FL
- The second term, **effect of heterogeneity**, can be reduced arbitrarily by increasing K , but the first term, **effect of initialization** remains.

Friendlier Results Under Lipschitzness

- Under G -Lipschitzness we can set $C = GE$ since $\|u_k^{(i)}\| \leq GE$
- This means **no clipping occurs** and we can get rid of $O\left(\frac{C}{LE} \sum_{i=1}^n \|f_i(w_k)\|\right)$ from the convergence criterion to obtain

$$\mathbb{E}[f(w_k)] - f(w^*) \leq \frac{8}{5} G \|w_0 - w^*\| + \frac{3}{4} E \sum_{i=1}^n \frac{1}{n} \|f_i(w_k)\|$$

- With $E = O(1)$ our bound matches the lower bound for the centralized convex and Lipschitz case with respect to the dependence on

**Are multiple local steps ($E > 1$)
beneficial or detrimental?**

Effect of Multiple Local Steps

- In a nutshell, increasing E mitigates the effect of initialization at the cost of increasing the effect of heterogeneity; the "best" value of E depends on which one is more dominant, and also the privacy level.
- Let us quantify this with an additional assumption.

Assumption 1

(i) For any $w \in \mathbb{R}^d$ and each $i \in [n]$, we have $\|kr f_i(w) - r f_i(w)\| \leq L \|kr f_i(w)\|$, for some $0 < L < 1$ and $\frac{1}{2L}$. (ii) Additionally, each f_i is G -Lipschitz over \mathbb{R}^d .

- For small enough ϵ , $\|kr f_i(w) - r f_i(w)\| \leq \epsilon (kr^2 f_i(w) - r f_i(w))$; so, we are basically assuming $\|kr^2 f_i(w) - r f_i(w)\| \leq \epsilon \|kr f_i(w)\|$ which is weaker than strong convexity.

Effect of Multiple Local Steps

- Recall $\epsilon = O\left(\frac{\rho \sqrt{d \log(1/\delta)}}{n}\right)$ is the privacy cost.

Proposition 1

Under Assumption 1, there exists a choice of C (depending on E), s.t. we get the following convergence guarantee:

$$\mathbb{E}[f(w_k)] - f(w^*) \leq 2Gkw_0 - wk + \frac{6}{5}E \sum_{i=1}^n \frac{1}{n} \sum_i \frac{11Gkw_0 - wk}{48} \frac{1}{L^2} \epsilon^2$$

- So if $\frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sum_i \epsilon^2 < O\left(\frac{G}{L^2} kw_0 - wk\right)$, then having a large value of E is beneficial; in particular, setting the maximum permissible value of E , which is $\frac{1}{2}$, is the best (in terms of smallest suboptimality gap). Otherwise, having a small value of E is better; specifically, $E = 1$ is the best.

So far we discussed

- Convergence of DP-FedAvg with clipping with and without Lipschitzness
- Role of ϵ in DP-FedAvg with clipping

But clipping has a potential issue!

$$\text{clip}(z; c) := z \min \left(1, \frac{c}{\|z\|} \right)$$

- As our FL algorithm converges, norm of model update $u_k^{(i)} = \frac{w_k w_{k;E}^{(i)}}{\|w_k\|}$ becomes small !
no clipping occurs
- But $\not\propto c$ regardless (can we adaptively reduce C ?)
- Hence, we enter a **low SNR regime** where added noise is dominant and hurts the convergence

Normalized Updates in DP-FL

Client-Update Normalization (Instead of Clipping)

- We propose to use

$$\text{norm}(z; c) := \frac{cz}{\|z\|} \quad \text{vs.} \quad \text{clip}(z; c) := z \min\left(1, \frac{c}{\|z\|}\right)$$

in DP-FedAvg, i.e.,

$$\text{Normalized Updates: } \frac{c u_k^{(i)}}{\|u_k^{(i)}\|} + \frac{\epsilon}{\|u_k^{(i)}\|}; \quad \frac{\epsilon}{\|u_k^{(i)}\|} \sim N(0_d; r^{-2} I_d)$$

- This ensures the **updates are uniformly bounded** and at the same time **noise will not overpower the update direction**, leading to better convergence and accuracy.
- For smaller C , normalization and clipping become equivalent.

Theoretical Comparison of Clipping and Normalization

See Section 5.1 and Remark 1 in the paper for a precise comparison. In summary:

- Our theory shows **normalization enjoys a smaller effect of initialization on convergence.**
- Not easy to characterize whether the effect of heterogeneity is smaller for normalization or clipping.
- But recall, the **effect of heterogeneity can be controlled by increasing K** , the number of rounds.
- Hence, **normalization has a better asymptotic convergence.**

Theory-guided Recommendation

Do normalization if we can afford training for large K .

Experiments: Synthetic Convex Quadratic Problems

- We set $(\epsilon; \delta) = (5; 10^{-6})$, $K = 500$, $E = 20$, and $n = 100$.
- $f_i(w) = \frac{1}{2}(w - w_i)^T Q_i (w - w_i)$
- w_i and $Q_i = A_i A_i^T$, where A_i is 200×20 , are formed randomly.
- Two different initialization
 - **I1:** $w_0 = w + z$, and
 - **I2:** $w_0 = w + \frac{z}{5}$,
- Finally, we consider full-device participation and vary ϵ and C .

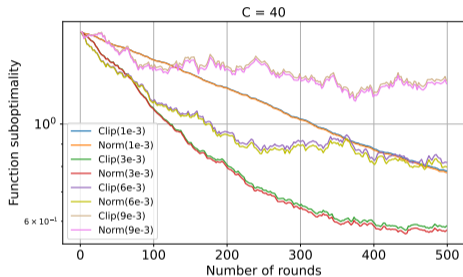
Synthetic Problems: 2D Trajectories

(a) I1: $C = 50$ and $\eta = 0.003$

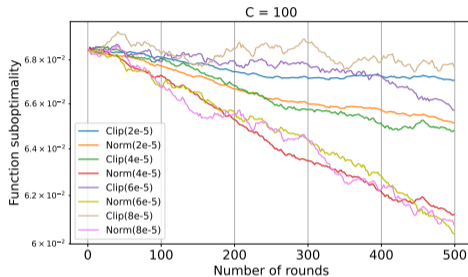
(b) I1: $C = 100$ and $\eta = 0.001$

- **DP-NormFedAvg** reaches closer to the optimum than **DP-FedAvg** with clipping.

Synthetic Problems: Convergence Curves



(a) I1: $C = 40$

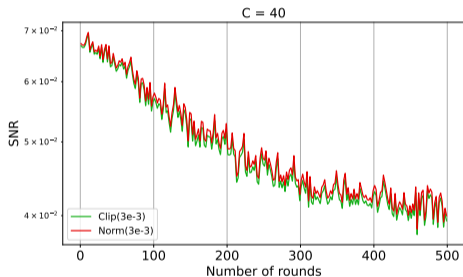


(b) I2: $C = 100$

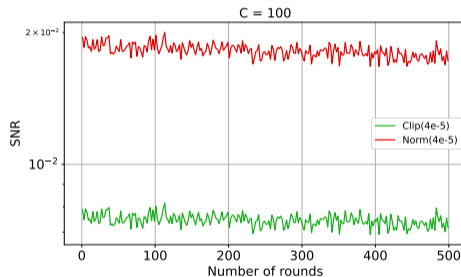
- Normalization does significantly better than clipping for large C .
- For smaller C normalization and clipping are nearly equivalent (as expected).

Synthetic Problems: SNR Comparison

- SNR: The ratio of average clipped/normalized per-client update and average per-client noise



(a) I1: $C = 40$



(b) I2: $C = 100$

- SNR of normalization is never lower than that of clipping, explaining the superiority of the former.

Logistic Regression on FMNIST, CIFAR-10 and CIFAR-100

- Average test accuracy over the last 5 rounds

FMNIST	$(5, 10^{-5})$ -DP	$(1.5, 10^{-5})$ -DP
Clipping	75.59%	56.90%
Normalization	77.72%	57.80%
FedAvg (w/o privacy)	83.43%	

CIFAR-10	$(5, 10^{-5})$ -DP	$(1.5, 10^{-5})$ -DP
Clipping	82.63%	81.53%
Normalization	84.21%	82.42%
FedAvg (w/o privacy)	85.64%	

CIFAR-100	$(5, 10^{-5})$ -DP	$(1.5, 10^{-5})$ -DP
Clipping	56.53%	41.33%
Normalization	59.36%	42.76%
FedAvg (w/o privacy)	64.61%	

Goal

Convergence analysis of private federated learning

- Established convergence of DP-FedAvg with clipping without Lipschitzness on smooth convex functions
- Effect of heterogeneity can be controlled while effect of initialization remains (cannot hope to do better)
- Role of local steps E : If $\frac{1}{n} \sum_{i=1}^n \lambda_i < O \left(\frac{G^2}{L^2} \right) \frac{1}{k} \frac{1}{\epsilon}$ having a large value of E is beneficial (under a suitable hessian assumption).

Theory-guided Recommendation

Normalized client updates instead of clipping for DP-FedAvg

- Ensures updates enjoy higher SNR
- Theoretical advantage over clipping in mitigating the effect of initialization

Thank you!

On the Convergence of Differentially Private Federated Learning on Non-Lipschitz Objectives via Clipping and Normalized Client Updates

<https://arxiv.org/abs/2106.07094>

Hiring Postdocs and PhD Students at Purdue ECE!

Abolfazl Hashemi (email: abolfazl@purdue.edu)