

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>



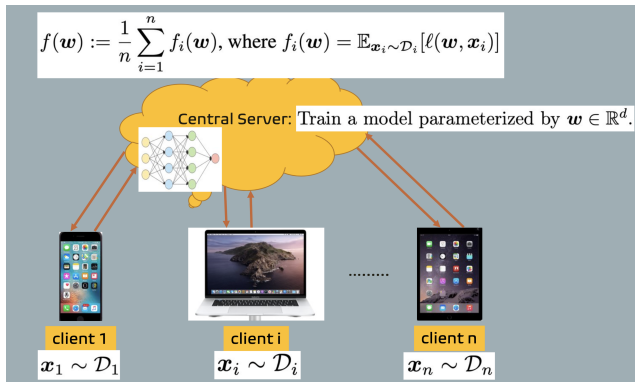
Rudrajit Das (UT CS)



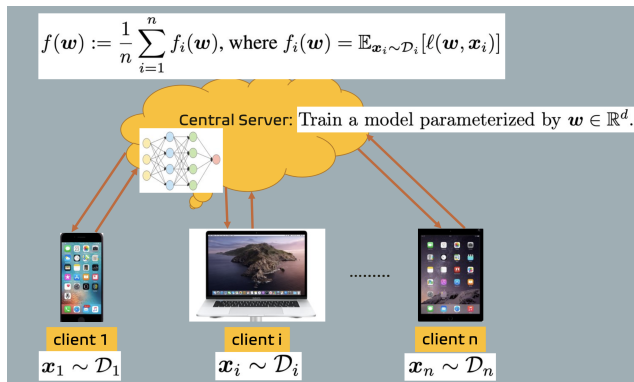
Sujay Sanghavi (UT ECE)



Inderjit S. Dhillon (UT CS)

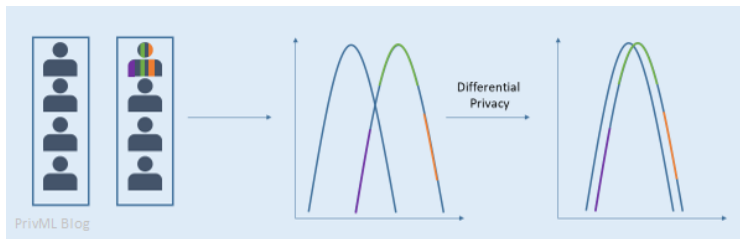


- Decentralized data (as opposed to traditional distributed learning)
- Different (resp., identical) \mathcal{D}_i 's: **heterogeneous** (resp., **homogeneous**) setting.



- 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?

- 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.



Neighboring datasets

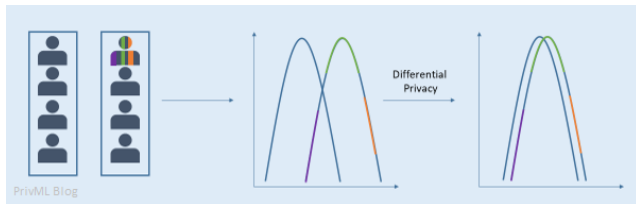
Two datasets $\mathcal{D} \in D_c$ and $\mathcal{D}' \in D_c$ are said to be neighboring if they differ in exactly one sample, and we denote this by $|\mathcal{D} - \mathcal{D}'| = 1$.

(ϵ, δ) -DP [DMNS06]

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

$$\mathbb{P}(\mathcal{M}(h(\mathcal{D})) \in \mathcal{R}) \leq e^\epsilon \mathbb{P}(\mathcal{M}(h(\mathcal{D}')) \in \mathcal{R}) + \delta.$$

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

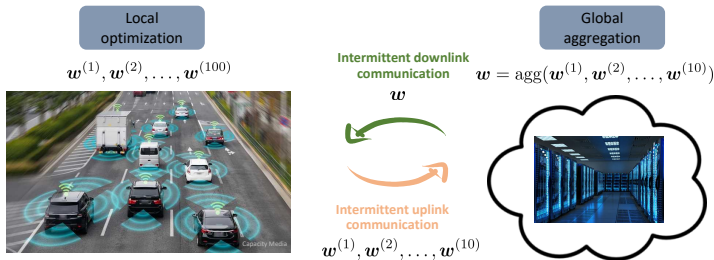


Gaussian mechanism [DR+14]

Let $\Delta := \sup_{\mathcal{D}, \mathcal{D}' \in \mathcal{D}_c: |\mathcal{D} - \mathcal{D}'|=1} \|h(\mathcal{D}) - h(\mathcal{D}')\|$. If we set $\mathcal{M}(h(\mathcal{D})) = h(\mathcal{D}) + \mathbf{Z}$, where $\mathbf{Z} \sim \mathcal{N}\left(\vec{0}_p, \frac{2 \log(1.25/\delta) \Delta^2}{\epsilon^2} \mathbf{I}_p\right)$, then the mechanism \mathcal{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?



- Local update: For E steps, do GD, i.e., $\mathbf{w}_{k,\tau+1}^{(i)} \leftarrow \mathbf{w}_{k,\tau}^{(i)} - \eta_k \nabla f_i(\mathbf{w}_{k,\tau}^{(i)})$
- Each client communicates its **update** $\mathbf{u}_k^{(i)} = \frac{\mathbf{w}_k - \mathbf{w}_{k,E}^{(i)}}{\eta_k}$ to server w.p. $\frac{r}{n}$ (total of K rounds)
- Global update: $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \frac{\beta_k}{r} \sum_{i \in \mathcal{S}_k} \mathbf{u}_k^{(i)}$
- Output: $\mathbf{w}_{\tilde{k}}$ with $\tilde{k} \sim \text{Unif}[0, K - 1]$.

$\mathbf{u}_k^{(i)}$ contains **local gradient** information \rightarrow needs to be made **private**!



DP-FedAvg with Clipping

- Maximum sensitivity, Δ , grows with norm of $\mathbf{u}_k^{(i)} = \frac{\mathbf{w}_k - \mathbf{w}_{k,E}^{(i)}}{\eta_k}$
- Assuming G -Lipschitzness, i.e., $\sup_{\theta \in \Theta} \|\nabla g(\theta)\|_2 \leq G$, this norm is at most $GE \rightarrow$ controlled additive noise
- In general, need to limit how large $\mathbf{u}_k^{(i)}$ can get to **remove unbounded impact of one client.**

Clipped Updates: $\mathbf{u}_k^{(i)} \min\left(1, \frac{C}{\|\mathbf{u}_k^{(i)}\|}\right) + \zeta_k^{(i)}, \quad \zeta_k^{(i)} \sim \mathcal{N}(0_d, r\sigma^2 \mathbf{I}_d)$

Theorem (Based on [ACG+16])

There exists an absolute constant $q > 0$ s.t. for $\varepsilon = \mathcal{O}(1)$, DP-FedAvg will be (ε, δ) -DP as long as

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

Convergence of DP-FedAvg with Clipping

Convexity

A function $g : \Theta \rightarrow \mathbb{R}$ is convex if $g(\lambda\theta + (1 - \lambda)\theta') \leq \lambda g(\theta) + (1 - \lambda)g(\theta')$ for any $\theta, \theta' \in \Theta$ and $0 \leq \lambda \leq 1$.

Smoothness

A function $g : \Theta \rightarrow \mathbb{R}$ is said to be L -smooth if for all $\theta, \theta' \in \Theta$, $\|\nabla g(\theta) - \nabla g(\theta')\|_2 \leq L\|\theta - \theta'\|_2$. If g is twice differentiable, then for all $\theta, \theta' \in \Theta$:

$$g(\theta') \leq g(\theta) + \langle \nabla g(\theta), \theta' - \theta \rangle + \frac{L}{2} \|\theta' - \theta\|_2^2.$$

Heterogeneity

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

Theorem 1

Suppose the f_i 's are convex and L -smooth over \mathbb{R}^d . Define $0 < \rho := \frac{\sqrt{qd \log(1/\delta)}}{n\varepsilon} < 1$. There exists constant local and global learning rates η and β , and a lower bound on the clipping threshold C_{low} , such that for any $C \geq C_{\text{low}}$ and in $K = \mathcal{O}\left(\frac{L\|\mathbf{w}_0 - \mathbf{w}^*\|}{C\rho^2}\right)$ rounds

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \min \left(f_i(\mathbf{w}_{\tilde{k}}) - f_i(\mathbf{w}^*), \mathcal{O}\left(\frac{C}{LE} \|\nabla f_i(\mathbf{w}_{\tilde{k}})\| \right) \right) \right] \leq \mathcal{O} \left(\frac{C}{E} \|\mathbf{w}_0 - \mathbf{w}^*\| + E \left(\frac{1}{n} \sum_{i=1}^n \Delta_i^* \right) \right) \rho.$$

- ρ : the price of privacy – increases as the level of privacy increases (i.e., ε and δ 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.

- Under G -Lipschitzness we can set $C = GE$ since $\|\mathbf{u}_k^{(i)}\| \leq GE$
- This means **no clipping occurs** and we can get rid of $\mathcal{O}\left(\frac{C}{LE}\|\nabla f_i(\mathbf{w}_{\tilde{k}})\|\right)$ from the convergence criterion to obtain

$$\mathbb{E}[f(\mathbf{w}_{\tilde{k}})] - f(\mathbf{w}^*) \leq \frac{8}{5} \left(G\|\mathbf{w}_0 - \mathbf{w}^*\| + \frac{3}{4} E \left(\frac{1}{n} \sum_{i=1}^n \Delta_i^* \right) \right) \rho.$$

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

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

- 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 $\mathbf{w} \in \mathbb{R}^d$ and each $i \in [n]$, we have $\|\nabla f_i(\mathbf{w} - \eta \nabla f_i(\mathbf{w})) - \nabla f_i(\mathbf{w})\| \geq \eta \lambda \|\nabla f_i(\mathbf{w})\|$, for some $0 < \lambda \leq L$ and $\eta \leq \frac{\rho}{2L}$. (ii) Additionally, each f_i is G -Lipschitz over \mathbb{R}^d .

- For small enough η , $\|\nabla f_i(\mathbf{w} - \eta \nabla f_i(\mathbf{w})) - \nabla f_i(\mathbf{w})\| = \Theta(\eta \|\nabla^2 f_i(\mathbf{w}) \nabla f_i(\mathbf{w})\|)$; so, we are basically assuming $\|\nabla^2 f_i(\mathbf{w}) \nabla f_i(\mathbf{w})\| \geq \Omega(\lambda \|\nabla f_i(\mathbf{w})\|)$ which is weaker than strong convexity.

- Recall $\rho = \mathcal{O}\left(\frac{\sqrt{d \log(1/\delta)}}{n\epsilon}\right)$ 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:

$$\mathbb{E}[f(\mathbf{w}_{\tilde{k}})] - f(\mathbf{w}^*) \leq \left(2G\|\mathbf{w}_0 - \mathbf{w}^*\| + \frac{6}{5}E \left\{ \frac{1}{n} \sum_{i=1}^n \Delta_i^* - \frac{11G\|\mathbf{w}_0 - \mathbf{w}^*\|\rho}{48} \left(\frac{\lambda^2}{L^2} \right) \right\} \right) \rho.$$

- So if $\frac{1}{n} \sum_{i=1}^n \Delta_i^* < \mathcal{O}\left(G\left(\frac{\lambda^2}{L^2}\right)\right)\|\mathbf{w}_0 - \mathbf{w}^*\|\rho$, then having a large value of E is beneficial; in particular, setting the maximum permissible value of E , which is $\frac{1}{2\rho}$, 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

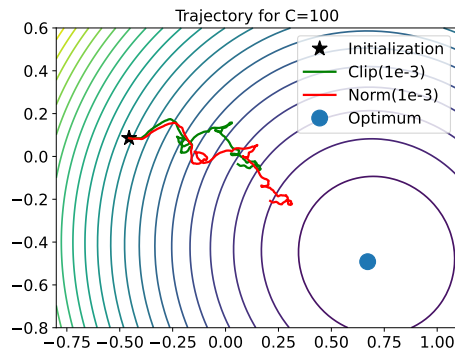
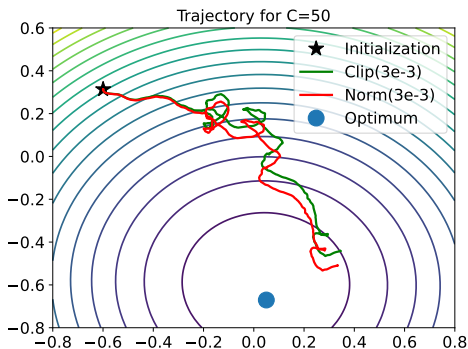
- Convergence of DP-FedAvg with clipping with and without Lipschitzness
- Role of E in DP-FedAvg with clipping

But clipping has a potential issue!

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

- As our FL algorithm converges, norm of model update $\mathbf{u}_k^{(i)} = \frac{\mathbf{w}_k - \mathbf{w}_{k,E}^{(i)}}{\eta}$ becomes small \rightarrow no clipping occurs
- But $\sigma \propto C$ regardless (can we adaptively reduce C ?)
- Hence, we enter a low SNR regime where added noise is dominant and hurts the convergence

- DP-FL on synthetic convex quadratic functions



- How can we ensure clients' update have higher SNR values?

Normalized Updates in DP-FL



Client-Update Normalization (Instead of Clipping)

- We propose to use

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

in DP-FedAvg, i.e.,

Normalized Updates:
$$\frac{C\mathbf{u}_k^{(i)}}{\|\mathbf{u}_k^{(i)}\|} + \zeta_k^{(i)}, \quad \zeta_k^{(i)} \sim \mathcal{N}(0_d, r\sigma^2\mathbf{I}_d)$$

- 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.

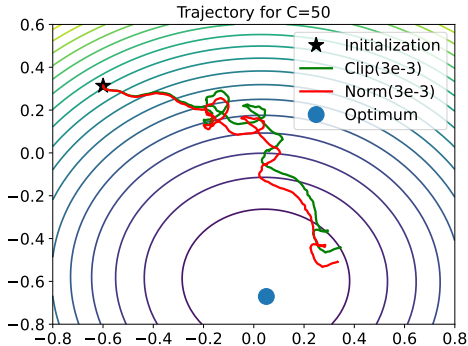
See Section 5.1 and Remark 1 in the paper for a precise comparison. In summary:

- 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.**

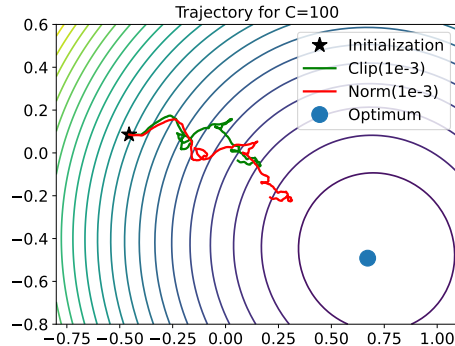
Theory-guided Recommendation

Do normalization if we can afford training for large K .

- We set $(\varepsilon, \delta) = (5, 10^{-6})$, $K = 500$, $E = 20$, and $n = 100$.
- $f_i(\mathbf{w}) = \frac{1}{2}(\mathbf{w} - \mathbf{w}_i^*)^T \mathbf{Q}_i(\mathbf{w} - \mathbf{w}_i^*)$
- \mathbf{w}_i^* and $\mathbf{Q}_i = \mathbf{A}_i \mathbf{A}_i^T$, where \mathbf{A}_i is 200×20 , are formed randomly.
- Two different initialization
 - I1: $\mathbf{w}_0 = \mathbf{w}^* + \mathbf{z}$, and
 - I2: $\mathbf{w}_0 = \mathbf{w}^* + \frac{\mathbf{z}}{5}$,
- Finally, we consider full-device participation and vary η and C .



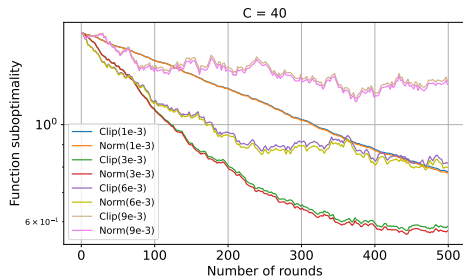
(a) I1: $C = 50$ and $\eta = 0.003$



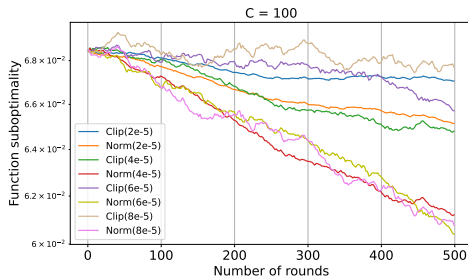
(b) I1: $C = 100$ and $\eta = 0.001$

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

Synthetic Problems: Convergence Curves



(a) I1: $C = 40$

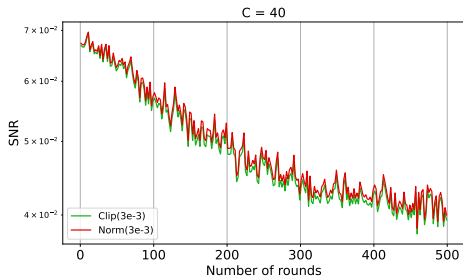


(b) I2: $C = 100$

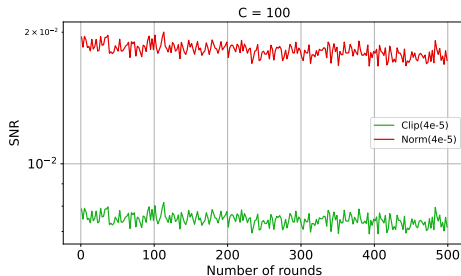
- Normalization does significantly better than clipping for large C .
- For smaller C normalization and clipping are nearly equivalent (as expected).

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.

- 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 $\frac{1}{n} \sum_{i=1}^n \Delta_i^* < \mathcal{O}\left(G\left(\frac{\lambda^2}{L^2}\right)\right) \|\mathbf{w}_0 - \mathbf{w}^*\|_\rho$ 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)