# Communication-Efficient and Deferentially-Private Federated Learning via Sketching with Sharp Guarantees

**Abstract**

Fedearted learning...

## 1 Introduction

The main contributions of this paper are as follows:

*

## 2 Federated Learning with Sketching

---

**Algorithm 1** CS: Count Sketch to compress $\mathbf{g} \in \mathbb{R}^d$.

---

1: **Inputs:** $\mathbf{g} \in \mathbb{R}^d, t, k, \mathbf{S}_{t \times k}, h_i(1 \leq i \leq t), sign_i(1 \leq i \leq t)$
2: **Compress vector $\tilde{\mathbf{g}} \in \mathbb{R}^d$ into $\mathbf{S}(\tilde{\mathbf{g}})$:**
3: **for** $\mathbf{g}_i \in \mathbf{g}$ **do**
4:    **for** $j = 1, \cdots, t$ **do**
5:       $\mathbf{S}[j][h_j(i)] = \mathbf{S}[j-1][h_{j-1}(i)] + \mathrm{sign}_j(i).\mathbf{g}_i$
6:    **end for**
7: **end for**
8: **return** $\mathbf{S}_{t \times k}$
9: **Query $\mathbf{g}_S \in \mathbb{R}^d$ from $\mathbf{S}(\mathbf{g})$:**
10: **for** $i = 1, \ldots, d$ **do**
11:    $\mathbf{S_g} = \mathrm{Median}\{\mathrm{sign}_j(i).\mathbf{S}[j][h_j(i)] : 1 \leq j \leq t\}$
12: **end for**
13: **Output:** $\mathbf{S}(\mathbf{g})$

---

**Algorithm 2** HEAVYMIX [1]

---

1: **Inputs:** $\mathbf{S_g}$; parameter-$k$
2: **Compress vector $\tilde{\mathbf{g}} \in \mathbb{R}^d$ into $\mathbf{S}(\tilde{\mathbf{g}})$:**
3: Query $\hat{\ell}_2^2 = (1 \pm 0.5)\|\mathbf{g}\|^2$ from sketch $\mathbf{S_g}$
4: $\forall j$ query $\hat{\mathbf{g}}_j^2 = \hat{\mathbf{g}}_j^2 \pm \frac{1}{2k}\|\mathbf{g}\|^2$ from sketch $\mathbf{S_g}$
5: $H = \{j|\hat{\mathbf{g}}_j \geq \frac{\hat{\ell}_2^2}{k}\}$ and $NH = \{j|\hat{\mathbf{g}}_j < \frac{\hat{\ell}_2^2}{k}\}$
6: $\mathrm{Top}_k = H \cup rand_\ell(NH)$, where $\ell = k - |H|$
7: Second round of communication to get exact values of $\mathrm{Top}_k$
8: **Output:** $\mathbf{g}_S : \forall j \in \mathrm{Top}_k : \mathbf{g}_{Si} = \mathbf{g}_i$ and $\forall \notin \mathrm{Top}_k : \mathbf{g}_{Si} = 0$

---

---

**Algorithm 3** PFL$(R, \tau, \eta, \gamma)$: Private Federated Learning with Sketching for homogeneous setting.

---

1: **Inputs:** $\boldsymbol{w}^{(0)}$ as an initial model shared by all local devices, the number of communication rounds $R$, the the number of local updates $\tau$, and global and local learning rates $\gamma$ and $\eta$, respectively
2: **for** $r = 0, \dots, R-1$ **do**
3:     **parallel for device** $j = 1, \dots, n$ **do:**
4:         Set $\boldsymbol{w}_j^{(0,r)} = \boldsymbol{w}^{(r-1)} - \gamma \mathbf{S}^{(r-1)}$
5:         **for** $c = 0, \dots, \tau-1$ **do**
6:             Sample a mini-batch $\xi_j^{(\ell,r)}$ and compute $\tilde{\mathbf{g}}_j^{(\ell,r)} \triangleq \nabla f_j(\boldsymbol{w}_j^{(\ell,r)}, \xi_j^{(c,r)})$
7:             $\boldsymbol{w}_j^{(\ell+1,r)} = \boldsymbol{w}_j^{(\ell,r)} - \eta \, \tilde{\mathbf{g}}_j^{(c,r)}$
8:         **end for**
9:         Device $j$ sends $\mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$ back to the server.
10:     Server **computes**
11:         $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1}^{} \mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$ and **broadcasts** $\mathbf{S}^{(r)}$ to all devices.
12:     **end parallel for**
13: **end**
14: **Output:** $\boldsymbol{w}^{(R-1)}$

---

---

**Algorithm 4** PFLGT$(R, \tau, \eta, \gamma)$: Private Federated Learning with Sketching and gradient tracking.

---

1: **Inputs:** $\boldsymbol{w}^{(0)}$ as an initial model shared by all local devices, the number of communication rounds $R$, the the number of local updates $\tau$, and global and local learning rates $\gamma$ and $\eta$, respectively
2: **for** $r = 0, \dots, R-1$ **do**
3:     **parallel for device** $j = 1, \dots, n$ **do:**
4:         Set $\mathbf{c}_j^{(r)} = \mathbf{c}_j^{(r-1)} - \frac{1}{\tau}\left(\mathbf{S}^{(r-1)} - \mathbf{S}_j^{(r-1)}\right)$
5:         Set $\boldsymbol{w}_j^{(0,r)} = \boldsymbol{w}^{(r-1)} - \gamma \mathbf{S}^{(r-1)}$
6:         **for** $\ell = 0, \dots, \tau-1$ **do**
7:             Sample a minibatch $\xi_j^{(\ell,r)}$ and compute $\tilde{\mathbf{g}}_j^{(\ell,r)} \triangleq \nabla f_j(\boldsymbol{w}_j^{(\ell,r)}, \xi_j^{(\ell,r)})$
8:             $\boldsymbol{w}_j^{(\ell+1,r)} = \boldsymbol{w}_j^{(\ell,r)} - \eta \left(\tilde{\mathbf{g}}_j^{(\ell,r)} - \mathbf{c}_j^{(r)}\right)$
9:         **end for**
10:         Device $j$ sends $\mathbf{S}_j^{(r)} \triangleq \mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$ back to the server.
11:     Server **computes**
12:         $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1}^{} \mathbf{S}_j^{(r)}$ and **broadcasts** $\mathbf{S}^{(r)}$ to all devices.
13:     **end parallel for**
14: **end**
15: **Output:** $\boldsymbol{w}^{(R-1)}$

---

---

**Algorithm 5** CFL($R, \tau, \eta, \gamma$): Communication-efficient Federated Learning with Sketching for homogeneous setting.

---

1: **Inputs:** $\boldsymbol{w}^{(0)}$ as an initial model shared by all local devices, the number of communication rounds $R$, the the number of local updates $\tau$, and global and local learning rates $\gamma$ and $\eta$, respectively
2: **for** $r = 0, \ldots, R-1$ **do**
3:     **parallel for device** $j = 1, \ldots, n$ **do**:
4:         Set $\boldsymbol{w}_j^{(0,r)} = \boldsymbol{w}^{(r-1)} - \gamma \underline{\mathbf{S}}^{(r)}$
5:         **for** $c = 0, \ldots, \tau - 1$ **do**
6:             Sample a mini-batch $\xi_j^{(\ell,r)}$ and compute $\tilde{\mathbf{g}}_j^{(\ell,r)} \triangleq \nabla f_j(\boldsymbol{w}_j^{(\ell,r)}, \xi_j^{(c,r)})$
7:             $\boldsymbol{w}_j^{(\ell+1,r)} = \boldsymbol{w}_j^{(\ell,r)} - \eta\, \tilde{\mathbf{g}}_j^{(c,r)}$
8:         **end for**
9:         Device $j$ sends $\mathbf{S}_j^{(r)} = \mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$ back to the server.
10:    Server **computes**
11:        $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1}^n \mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$
12:        Sever runs $\underline{\mathbf{S}}^{(r)} = \texttt{HEAVYMIX}(\mathbf{S}^{(r)})$.
13:    **end parallel for**
14: **end**
15: **Output:** $\boldsymbol{w}^{(R-1)}$

---

---

**Algorithm 6** CFL($R, \tau, \eta, \gamma$): Communication-efficient Federated Learning with Sketching and gradient tracking.

---

1: **Inputs:** $\boldsymbol{w}^{(0)}$ as an initial model shared by all local devices, the number of communication rounds $R$, the the number of local updates $\tau$, and global and local learning rates $\gamma$ and $\eta$, respectively
2: **for** $r = 0, \ldots, R-1$ **do**
3:     **parallel for device** $j = 1, \ldots, n$ **do**:
4:         Set $\boldsymbol{w}_j^{(0,r)} = \boldsymbol{w}^{(r-1)} - \gamma \underline{\mathbf{S}}^{(r)}$
5:         Set $\mathbf{c}_j^{(r)} = \mathbf{c}_j^{(r-1)} - \frac{1}{\tau}\left(\underline{\mathbf{S}}^{(r)} - \underline{\mathbf{S}}_j^{(r)}\right)$
6:         **for** $c = 0, \ldots, \tau - 1$ **do**
7:             Sample a mini-batch $\xi_j^{(\ell,r)}$ and compute $\tilde{\mathbf{g}}_j^{(\ell,r)} \triangleq \nabla f_j(\boldsymbol{w}_j^{(\ell,r)}, \xi_j^{(c,r)})$
8:             $\boldsymbol{w}_j^{(\ell+1,r)} = \boldsymbol{w}_j^{(\ell,r)} - \eta\, \tilde{\mathbf{g}}_j^{(c,r)}$
9:         **end for**
10:        Device $j$ sends $\mathbf{S}_j^{(r)} = \mathbf{S}\left(\boldsymbol{w}_j^{(0,r)} - \boldsymbol{w}_j^{(\tau,r)}\right)$ back to the server.
11:    Server **runs**
12:        $\underline{\mathbf{S}}_j^{(r)} = \texttt{HEAVYMIX}(\mathbf{S}_j^{(r)})$ and returns $\underline{\mathbf{S}}_j^{(r)}$ to server $j$.
13:        $\underline{\mathbf{S}}^{(r)} = \frac{1}{p} \sum_{j=1}^n \underline{\mathbf{S}}_j^{(r)}$
14:        Sever broadcasts $\underline{\mathbf{S}}^{(r)}$.
15:    **end parallel for**
16: **end**
17: **Output:** $\boldsymbol{w}^{(R-1)}$

---

# 3 Convergence Analysis

## 3.1 Assumptions

# 4 Experiments

# 5 Conclusion

# References

[1] N. Ivkin, D. Rothchild, E. Ullah, I. Stoica, R. Arora *et al.*, "Communication-efficient distributed sgd with sketching," in *Advances in Neural Information Processing Systems*, 2019, pp. 13 144–13 154.