
Variance Reduced Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To be completed

2 1 SAGA-like Algorithms

3 1.1 Local SGD

Algorithm 1 SAGA Local SGD

- 1: **Input:** Local learning rate γ and global learning rate η and communication period p .
- 2: **Init:** $g^{(k)} = \frac{1}{n} \sum_{i=1}^n g^{(0)}$.
- 3: **for** $k = 0, 1, \dots, K$ **do**
- 4: Draw two independent and distinct indices i_k and j_k
- 5: **for** $\tau = k, \dots, k + p - 1$ **do**
- 6: Compute the following quantity

$$v_{i_k}^{(\tau)} = v_{i_k}^{(\tau-1)} - \gamma(\nabla f_{i_k}(x^{(k)}) - \nabla f_{i_k}(\alpha_{i_{(k)}}^t) + g^{(k)})$$

- 7: **end for**
- 8: $v_{i_k}^{(k)} \leftarrow v_{i_k}^{(k+p-1)}$
- 9: Update the global model

$$x^{(k+1)} = x^{(k)} - \eta v_{i_k}^{(k)}$$

- 10: Update $\alpha_{j_k}^{(k+1)} = x^{(k)}$ and $\alpha_j^{(k+1)} = \alpha_j^{(k)}$ for $j \neq j_k$
 - 11: Update $g^{(k+1)} = g^{(k)} - \frac{1}{n} \left(\nabla f_{j_k}(\alpha_{j_k}^{(k)}) - \nabla f_{j_k}(\alpha_{j_k}^{(k+1)}) \right)$
 - 12: **end for**
-

4 1.2 FedSVRG with Sketching

Algorithm 2 FedSVRG: SVRG Federated Learning algorithm with Sketching.

```

1: Inputs:  $\mathbf{x}^{(0)}$  initial common model, communication rounds  $R$ , the number of local updates  $K$ ,
   and global and local learning rates  $\gamma$  and  $\eta$ 
2: for  $r = 0, \dots, R - 1$  do
3:   parallel for device  $j = 1, \dots, n$  do:
4:     Computes  $\Phi^{(r)} \triangleq \mathbf{Q} [\mathbf{S}^{(r-1)}]$  ( $\mathbf{Q}$  is any query function)
5:     Set  $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)} - \gamma \Phi^{(r)}$ 
6:     Set  $\mathbf{x}_j^{(0,r)} = \mathbf{x}^{(r)}$ 
7:     Compute full gradient  $\nabla f_j(\mathbf{x}^{(\kappa(r))}) = \frac{1}{n_j} \sum_{i=1}^{n_j} \nabla f_j(\mathbf{x}^{(\kappa(r))}, \xi_i)$ 
8:     for  $\ell = 0, \dots, K - 1$  do
9:       Sample an index  $i_\ell$  uniformly on  $[n_j]$ 
10:       $\mathbf{x}_j^{(\ell+1,r)} = \mathbf{x}_j^{(\ell,r)} - \eta (\nabla f_j(\mathbf{x}_j^{(\ell,r)}, \xi_{j,i_\ell}) - \nabla f_j(\mathbf{x}^{(\kappa(r))}, \xi_{j,i_\ell}) + \nabla f_j(\mathbf{x}^{(\kappa(r))}))$ 
11:    end for
12:    Device  $j$  sends  $\mathbf{S}_j^{(r)} \triangleq \mathbf{S}_j (\mathbf{x}_j^{(0,r)} - \mathbf{x}_j^{(K,r)})$  back to the server.
13:  Server computes
14:     $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1} \mathbf{S}_j^{(r)}$  and broadcasts  $\mathbf{S}^{(r)}$  to all devices.
15:  end parallel for
16: end
17: Output:  $\mathbf{x}^{(R-1)}$ 

```

5 **Particularity:** $\kappa(r)$ is the epoch number at which we compute both control variate terms. We can
6 have $\kappa(r) = r$ or something else (to tune).

7 1.3 FedSAGA with Sketching

Algorithm 3 FedSAGA: SAGA Federated Learning algorithm with Sketching.

```

1: Inputs:  $\mathbf{x}^{(0)}$  initial common model, communication rounds  $R$ , the number of local updates  $K$ ,
   and global and local learning rates  $\gamma$  and  $\eta$ 
2: for  $r = 0, \dots, R - 1$  do
3:   parallel for device  $j = 1, \dots, n$  do:
4:     Computes  $\Phi^{(r)} \triangleq \mathbf{Q} [\mathbf{S}^{(r-1)}]$  ( $\mathbf{Q}$  is any query function)
5:     Set  $\mathbf{x}^{(r)} = \mathbf{x}^{(r-1)} - \gamma \Phi^{(r)}$ 
6:     Set  $\mathbf{x}_j^{(0,r)} = \mathbf{x}^{(r)}$ 
7:     Compute full gradient  $\nabla f_j(\mathbf{x}^{(r)}) = \frac{1}{n_j} \sum_{i=1}^{n_j} \nabla f_j(\mathbf{x}^{(r)}, \xi_i)$ 
8:     for  $\ell = 0, \dots, K - 1$  do
9:       Sample an indices  $(i_\ell, q_\ell)$  independently and uniformly on  $[n_j]$ 
10:       $\mathbf{x}_j^{(\ell+1,r)} = \mathbf{x}_j^{(\ell,r)} - \eta (\nabla f_j(\mathbf{x}_j^{(\ell,r)}, \xi_{j,i_\ell}) - \nabla f_j(\mathbf{x}_j^{(t_{i_\ell},r)}, \xi_{j,i_\ell}) + \bar{F}_j^{(r)})$ 
11:      where  $\bar{F}_j^{(r)} = \frac{1}{n_j} \sum_{i=1}^{n_j} \nabla f_j(\mathbf{x}^{(r)}, \xi_i) + \frac{1}{n_j} (\nabla f_j(\mathbf{x}_j^{(\ell,r)}, \xi_{j,q_\ell}) - \nabla f_j(\mathbf{x}_j^{(t_{q_\ell},r)}, \xi_{j,q_\ell}))$ 
12:    end for
13:    Device  $j$  sends  $\mathbf{S}_j^{(r)} \triangleq \mathbf{S}_j (\mathbf{x}_j^{(0,r)} - \mathbf{x}_j^{(K,r)})$  back to the server.
14:  Server computes
15:     $\mathbf{S}^{(r)} = \frac{1}{p} \sum_{j=1} \mathbf{S}_j^{(r)}$  and broadcasts  $\mathbf{S}^{(r)}$  to all devices.
16:  end parallel for
17: end
18: Output:  $\mathbf{x}^{(R-1)}$ 

```

8 **Particularity:** We need to store each local models on each device j in order to compute the control
9 variate terms. No epoch tuning though. Indeed $\mathbf{x}_j^{(t_{i_\ell}, r)}$ is the value of local model on device j when
10 index i_ℓ was drawn last.

11 **2 Numerical Examples**