

Federated Learning - Summary

Willy Fitra Hendria

willyfitrahendria@gmail.com



Introduction

- Introduced in a paper entitled “Communication-Efficient Learning of Deep Networks from Decentralized Data”, H. Brendan McMahan, et al
- The client’s data doesn’t need to be shared in order to learn a global model. Instead, the global model is learned by aggregating the locally-computed updates from each client devices.
- Addresses the concerns of privacy and communication costs .

Ideal Problems for Federated Learning

Have the following properties:

- 1) Training on real-world data provides a distinct advantage
- 2) The data is privacy sensitive or large in size
- 3) Labels on the data can be inferred naturally from user interaction

Federated Optimization

Several key properties that differ from typical distributed optimization problem:

- **Non-IID**
 - Any particular client's dataset will not be representative of the population distribution
- **Unbalanced**
 - Varying amounts of local training data
- **Massively distributed**
 - The number of clients to be much larger than the average number of examples per client
- **Limited communication**
 - Clients are frequently offline or on slow or expensive connections

Algorithm (FederatedAveraging)

Algorithm 1 FederatedAveraging. The K clients are indexed by k ; B is the local minibatch size, E is the number of local epochs, and η is the learning rate.

Server executes:

```
initialize  $w_0$ 
for each round  $t = 1, 2, \dots$  do
   $m \leftarrow \max(C \cdot K, 1)$ 
   $S_t \leftarrow$  (random set of  $m$  clients)
  for each client  $k \in S_t$  in parallel do
     $w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$ 
   $w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ 
```

ClientUpdate(k, w): // Run on client k

```
 $\mathcal{B} \leftarrow$  (split  $\mathcal{P}_k$  into batches of size  $B$ )
for each local epoch  $i$  from 1 to  $E$  do
  for batch  $b \in \mathcal{B}$  do
     $w \leftarrow w - \eta \nabla \ell(w; b)$ 
return  $w$  to server
```

- The amount of computation is controlled by three key parameters: C (Client Fraction), E (Local Epoch), and B (Local Minibatch).
- Models are locally-trained on each 1 client:
$$w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$$
- Central server averages the resulting models
- $$w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$$
- - With $E=1$ and $B=\infty$, it corresponds to FedSGD

Experimental Results

2NN <i>C</i>	IID		Non-IID	
	$B = \infty$	$B = 10$	$B = \infty$	$B = 10$
0.0	1455	316	4278	3275
0.1	1474 (1.0×)	87 (3.6×)	1796 (2.4×)	664 (4.9×)
0.2	1658 (0.9×)	77 (4.1×)	1528 (2.8×)	619 (5.3×)
0.5	— (—)	75 (4.2×)	— (—)	443 (7.4×)
1.0	— (—)	70 (4.5×)	— (—)	380 (8.6×)
CNN, $E = 5$				
0.0	387	50	1181	956
0.1	339 (1.1×)	18 (2.8×)	1100 (1.1×)	206 (4.6×)
0.2	337 (1.1×)	18 (2.8×)	978 (1.2×)	200 (4.8×)
0.5	164 (2.4×)	18 (2.8×)	1067 (1.1×)	261 (3.7×)
1.0	246 (1.6×)	16 (3.1×)	— (—)	97 (9.9×)

With $B=\infty$, there is only a small advantage in increasing C . Using smaller $B=10$ shows a significant improvement in using $C \geq 0.1$, especially in the non-IID case

Experimental Results

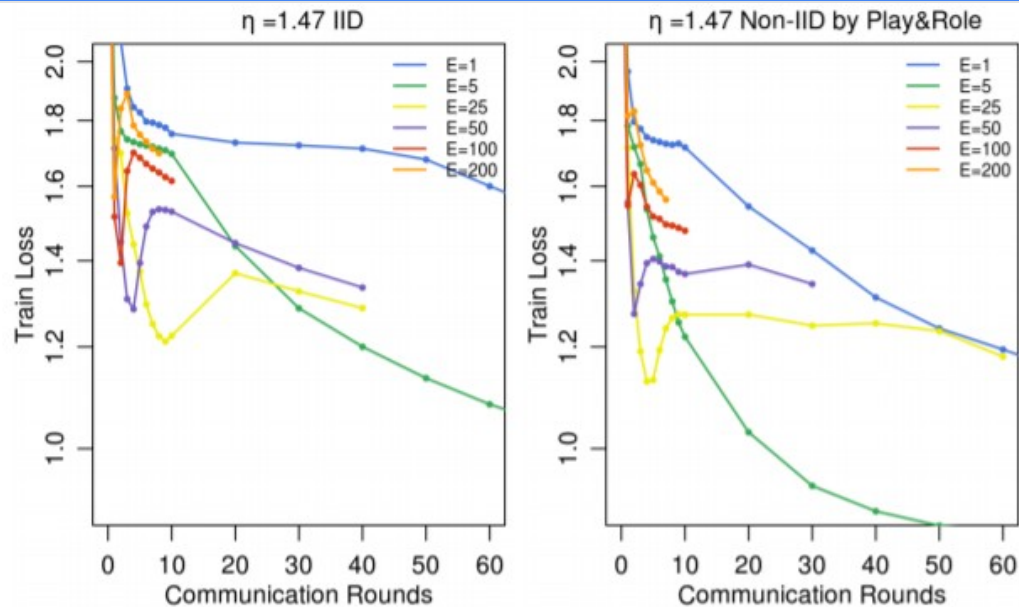
MNIST CNN, 99% ACCURACY					
CNN	E	B	u	IID	Non-IID
FEDSGD	1	∞	1	626	483
FEDAVG	5	∞	5	179 (3.5 \times)	1000 (0.5 \times)
FEDAVG	1	50	12	65 (9.6 \times)	600 (0.8 \times)
FEDAVG	20	∞	20	234 (2.7 \times)	672 (0.7 \times)
FEDAVG	1	10	60	34 (18.4 \times)	350 (1.4 \times)
FEDAVG	5	50	60	29 (21.6 \times)	334 (1.4 \times)
FEDAVG	20	50	240	32 (19.6 \times)	426 (1.1 \times)
FEDAVG	5	10	300	20 (31.3 \times)	229 (2.1 \times)
FEDAVG	20	10	1200	18 (34.8 \times)	173 (2.8 \times)

SHAKESPEARE LSTM, 54% ACCURACY					
LSTM	E	B	u	IID	Non-IID
FEDSGD	1	∞	1.0	2488	3906
FEDAVG	1	50	1.5	1635 (1.5 \times)	549 (7.1 \times)
FEDAVG	5	∞	5.0	613 (4.1 \times)	597 (6.5 \times)
FEDAVG	1	10	7.4	460 (5.4 \times)	164 (23.8 \times)
FEDAVG	5	50	7.4	401 (6.2 \times)	152 (25.7 \times)
FEDAVG	5	10	37.1	192 (13.0 \times)	41 (95.3 \times)

MNIST 2NN	E	B	u	IID	Non-IID
FEDSGD	1	∞	1	1468	1817
FEDAVG	10	∞	10	156 (9.4 \times)	1100 (1.7 \times)
FEDAVG	1	50	12	144 (10.2 \times)	1183 (1.5 \times)
FEDAVG	20	∞	20	92 (16.0 \times)	957 (1.9 \times)
FEDAVG	1	10	60	92 (16.0 \times)	831 (2.2 \times)
FEDAVG	10	50	120	45 (32.6 \times)	881 (2.1 \times)
FEDAVG	20	50	240	39 (37.6 \times)	835 (2.2 \times)
FEDAVG	10	10	600	34 (43.2 \times)	497 (3.7 \times)
FEDAVG	20	10	1200	32 (45.9 \times)	738 (2.5 \times)

With $C=0.1$, adding more local updates per round (increase E & decrease B) can produce a dramatic decrease in communication costs

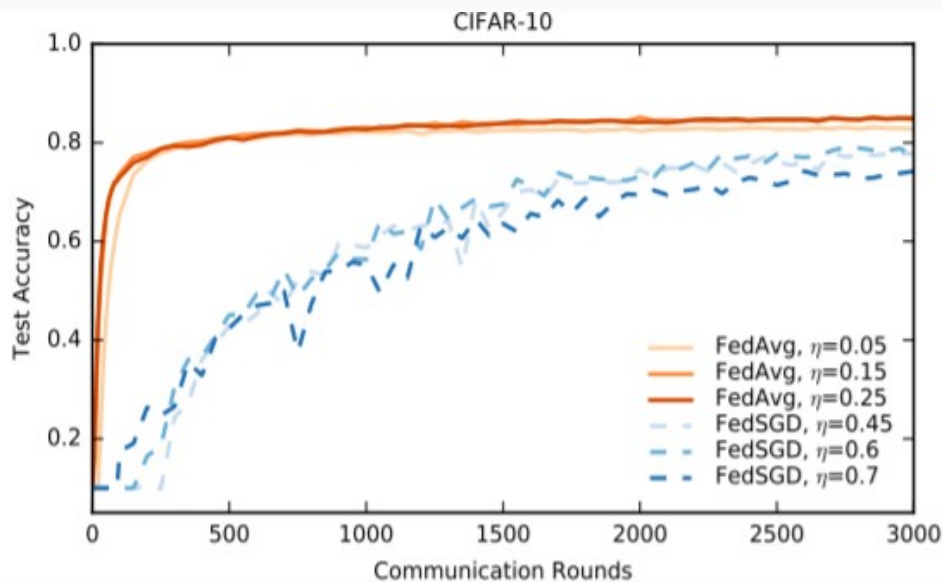
Experimental Results



For very large E , FedAvg can plateau or diverge.

It may be useful to decay the amount of local computation per round (moving to smaller E or larger B)

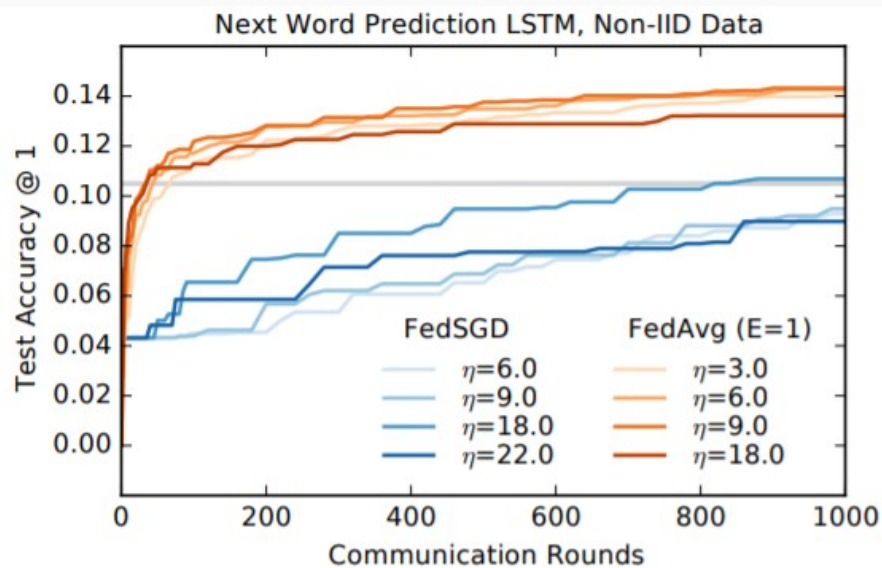
Experimental Results



Acc.	80%		82%		85%	
SGD	18000	(—)	31000	(—)	99000	(—)
FedSGD	3750	(4.8×)	6600	(4.7×)	N/A	(—)
FedAVG	280	(64.3×)	630	(49.2×)	2000	(49.5×)

On the CIFAR-10 dataset, FedAVG has less number of communication rounds compares to FedSGD and baseline SGD.

Experimental Results



On a large-scale LSTM experiment, FedSGD with $\eta = 18.0$ required 820 rounds to reach 10.5%, while FedAvg with $\eta=9.0$ reached an accuracy of 10.5% in only 35 rounds, which is 23X fewer than fedSGD.