

Dubhe: Towards <u>Data Unb</u>iasedness with <u>H</u>omomorphic <u>Encryption in Federated Learning Client Selection</u>

Shulai Zhang, Zirui Ii, Quan Chen, Wenli Zheng, Jingwen Leng, Minyi Guo Shanghai Jiao Tong University



Background



Federated Learning

Demand: Security promise

• Homomorphic Encryption

Problem: Statistical heterogeneity

- Local data skewness
- Data discrepancy among clients
- Global data skewness

Method: Client selection

- Random selection
- Greedy selection





The framework of Federated Learning



Background



Federated Learning

- **Demand:** Security promise
 - Homomorphic Encryption

Problem: Statistical heterogeneity

- Local data skewness
- Data discrepancy among clients
- Global data skewness

Method: Client selection

- Random selection
- Greedy selection





The framework of Federated Learning



Motivation



Shanghai Jiao Tong University

Federated Learning

- **Demand:** Security promise
 - Homomorphic Encryption
- Problem: Statistical heterogeneity
 - Local data skewness
 - Data discrepancy among clients
 - Global data skewness

Method: Client selection

- Random selection
- Greedy selection



CIFAR10

Discrepancy of clients in CIFAR10 classification



50th International Conference on Parallel Processing (ICPP) August 9-12, 2021 in Virtual Chicago, IL



8 9

mmm mmm

<u>പ്പം</u>

6 7

Motivation



Shanghai Jiao Tong University

Federated Learning

- **Demand:** Security promise
 - Homomorphic Encryption

Problem: Statistical heterogeneity

- Local data skewness
- Data discrepancy among clients
- Global data skewness

Method: Client selection

- Random selection
- Greedy selection



Global data skewness in CIFAR10 classification





Mathematical Demonstration



Shanghai Jiao Tong University

Federated Learning



• Greedy selection









The framework of Dubhe







- The *registry* encodes the client's data distribution information in a one-hot manner.
- In classification problems, the registry encodes each client's data distribution by its *dominating classes*.









Registration and **Probability Calculation**









Shanghai Jiao Tong University

Registration and Probability Calculation









MNIST **MNIST-2/0.5** MNIST-2/1.0 **MNIST-2/1.5** and the second second second 0.98 Test accuracy o 8.0 0.8 0.8 0.96 Lest accuracy 0.94 0.92 0.975 0.98 0.975 0.97 0.6 0.6 0.96 0.97 0.965 0.4 0.4 0.4 0.96 0 0.94 0.965 $\begin{array}{c} 0.5\\ E_{MD} \\ a_{v_g} \end{array}$ 200 180 160 160 180 200 180 200 160 2 0.2 0.2 0.2 5 Ø 200 200 200 0 150 0 150 1.5 10 50 100 150 0 50 100 50 100 Number of rounds CIFAR10 CIFAR10-10/0.5 CIFAR10-10/1.0 CIFAR10-10/1.5 0.6 0.6 0.6 0.6 0.5 0.5 0.5 Test accuracy 8.0 Test accuracy 7.0 Test accuracy Test accuracy 0.5 0.4 0.4 0.4 0.3 0.3 Random 0.3 0.2 0.2 Dubhe 0.2 0.1 0.1 0. Greedy 0 0.5 2 0 EMD 1.0 0 5ρ 0 200 400 600 800 1000 0 200 600 800 1000 0 200 400 600 800 1000 400 1.5 10 Number of rounds

Average accuracy





50th International Conference on Parallel Processing (ICPP)

Test accuracy curves on MNIST and CIFAR10

August 9-12, 2021 in Virtual Chicago, IL







The data balancing performance of Dubhe is approaching the performance of the greedy selection.

Results on FEMNIST







Multi-time selection









Shanghai Jiao Tong University

Multi-time selection







Multi-time selection



Shanghai Jiao Tong University

Results with multi-time selection. *M* refers to MNIST and *C* refers to CIFAR10.

Н	1	2	5	10	20	opt
EMD^*	0.2946	0.2588	0.2176	0.1971	0.1750	0.0144
Acc^{M}	0.9662	0.9668	0.9665	0.9684	0.9678	0.9694
eta^M	0.0%	17.6%	10.5%	69.5%	51.5%	100%
Acc^{C}	0.4300	0.4518	0.4486	0.4441	0.4577	0.5295
β^{C}	0.00%	14.8%	12.6%	9.5%	18.8%	100%

EMD^{*} decreases with larger *H*, thereby improving the model accuracy









Shanghai Jiao Tong University

In-Cooperation

Registry sparsity



Conclusions



- The impact of data skewness on the performance degradation in FL is mathematically demonstrated.
- We propose *Dubhe*, a proactive client selection system to balance skewed data
 - pluggable, adaptive and robust to various FL settings
 - with negligible encryption and communication overhead
 - improves the training performance without bringing security threats.









Thanks for listening!

Any questions? Comments? Please contact:

Shulai Zhang: zslzsl1998@sjtu.edu.cn