



Abstract

- Considered the problem of Federated Learning (FL) under non-i.i.d data setting
- Provided an improved estimate of the empirical loss at each node by using a weighted average of losses across nodes with a penalty term
- Assigned uneven weights to different nodes by taking a Bayesian approach to the problem where learning for each node is cast as maximizing the likelihood of a joint distribution of losses for a given neural network of a node, by using data across nodes
- Provided a PAC learning guarantee on the objective function which revealed that the true average risk is no more than the proposed objective and the error term
- Leveraged this guarantee to propose an algorithm called Omni-Fedge
- Using MNIST and Fashion MNIST data-sets, we showed that the performance of the proposed algorithm is significantly better than existing algorithms

Index Terms – Federated Learning, Neural Network, Bayesian Approach, Distributed Machine Learning, PAC Learning.

Introduction and Problem Setting

- We address the problem of improving FL performance with non-i.i.d data
- We consider a federated system with N edge-devices that communicate with one federating server (FS)
- We assume that the data points are independent but not necessarily identically distributed across edge-devices
- Further, we assume that the data at edge-device $i \in \{1, \dots, N\}$ is sampled from a distribution D_i
- Neural network weights are divided into two parts, viz, shared ($\theta^{(sh)}$) and task-specific ($\theta^{(i)}$)

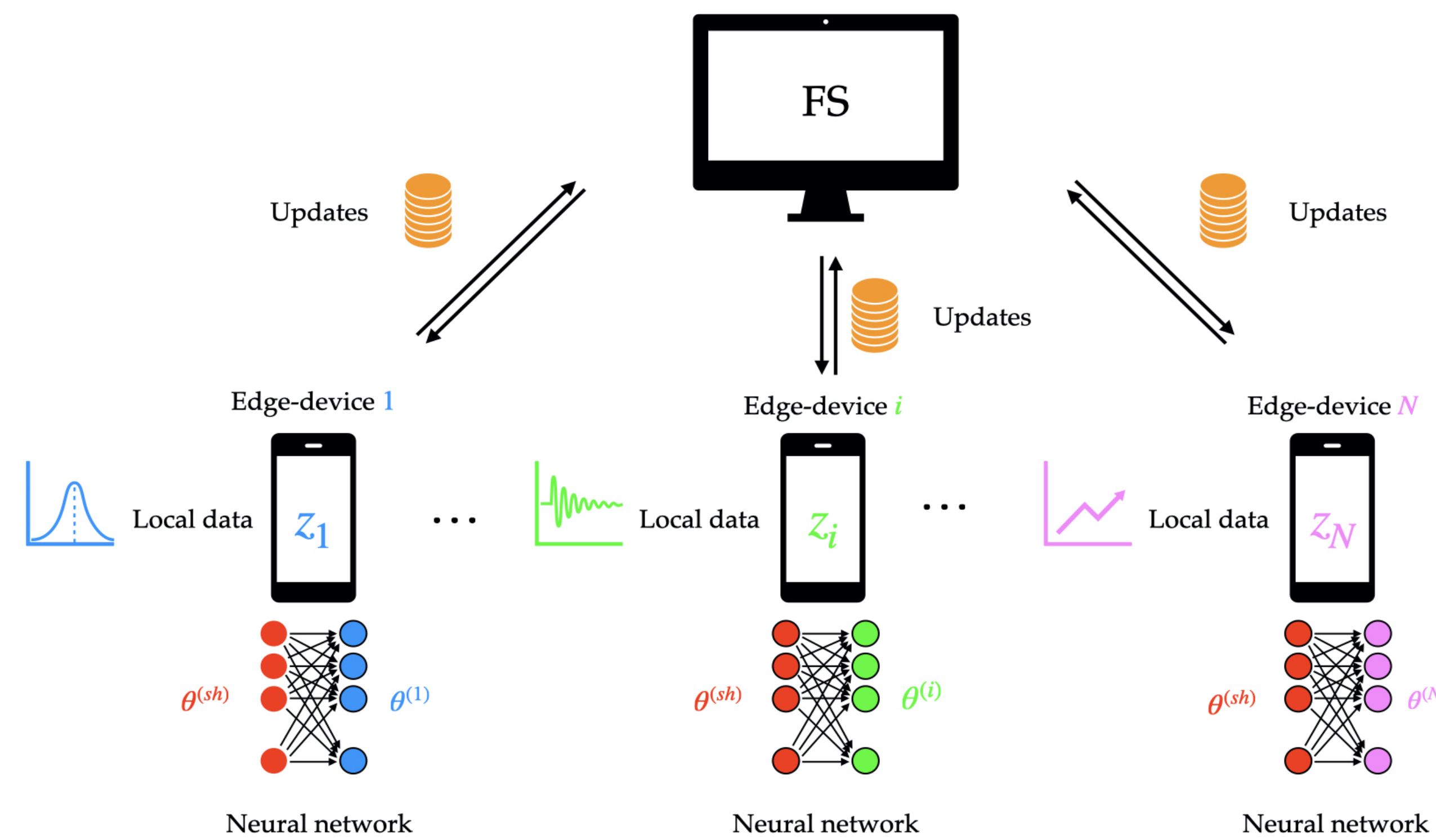


Figure: Federated Setup

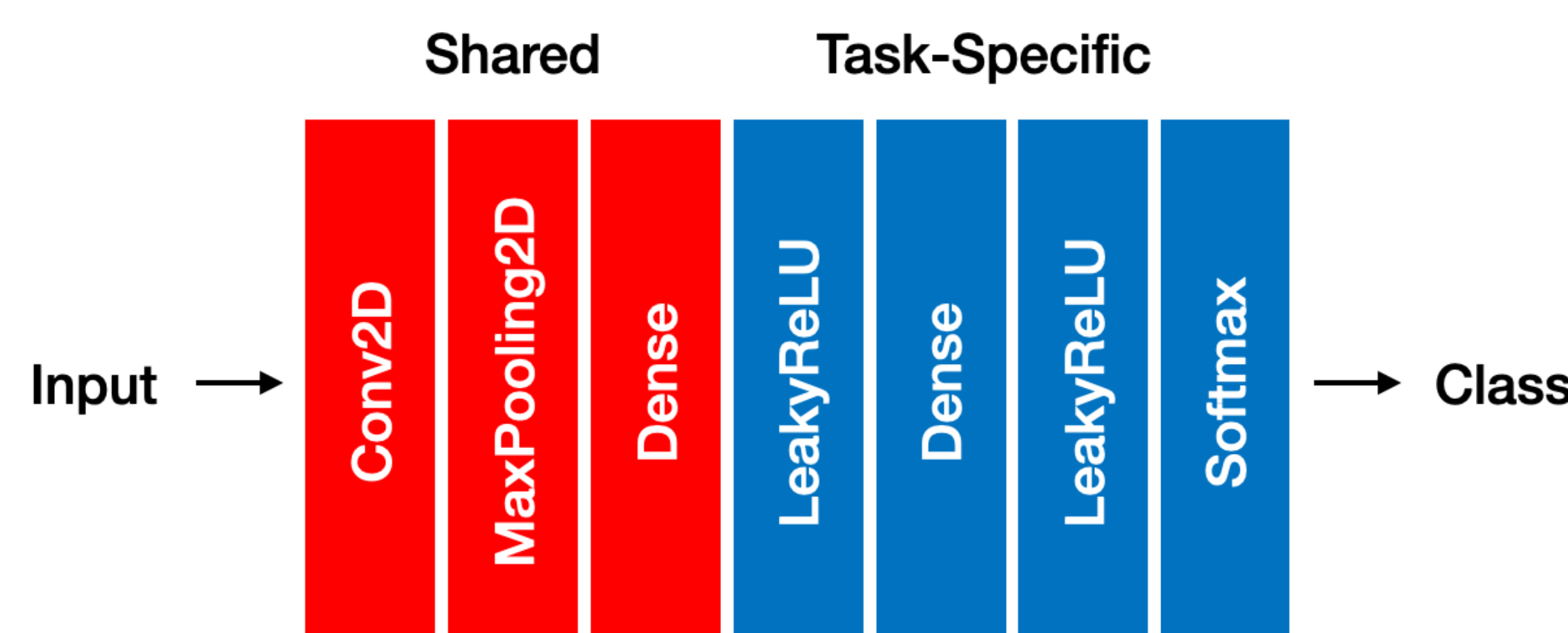
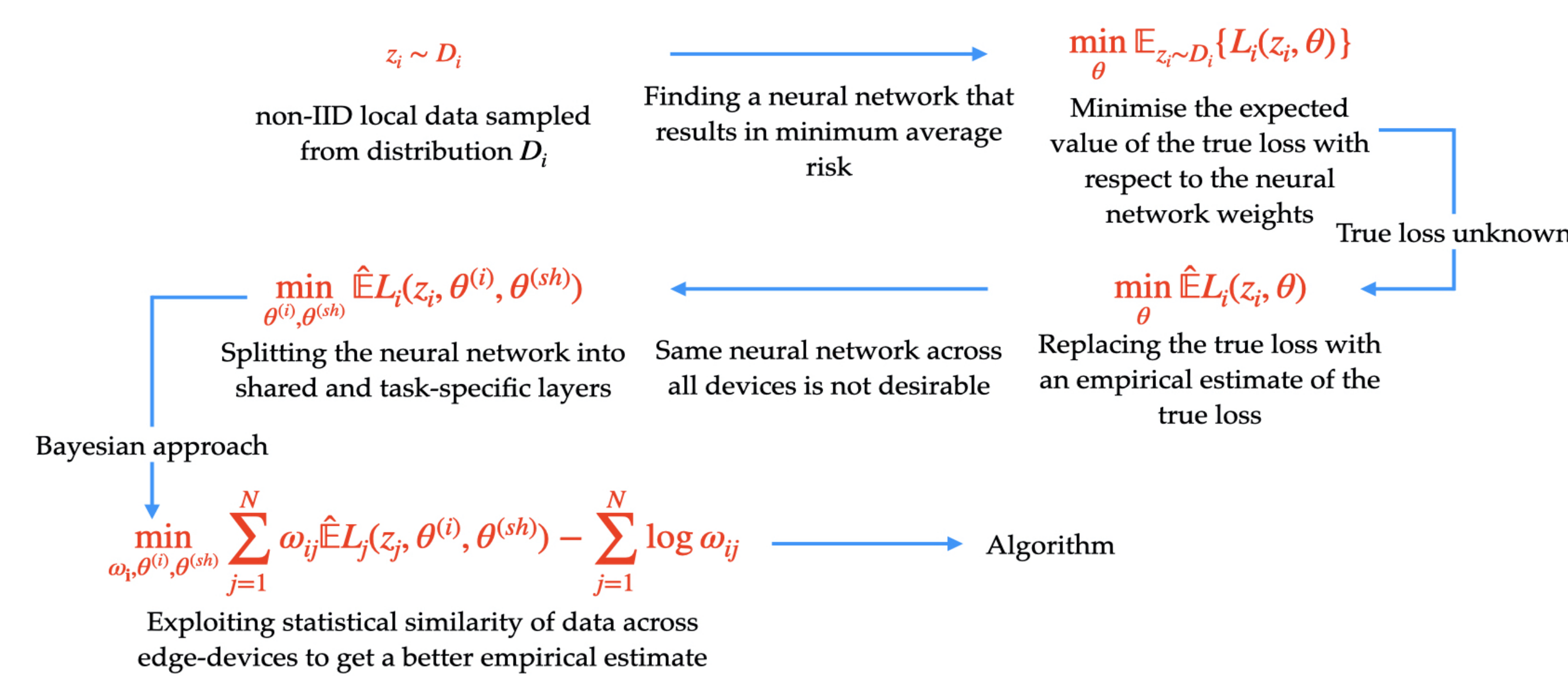


Figure: Neural Network

Motivation



Bayesian Approach



$$Q(\theta^{(i)}, \theta^{(sh)}) = \left(\prod_{j=1}^N \omega_{ij} \right) e^{-\sum_{j=1}^N \omega_{ij} \hat{E}L_j(z_j, \theta^{(i)}, \theta^{(sh)})}$$

Joint distribution of losses as assumed by edge-device i

$$P(\theta^{(i)}, \theta^{(sh)}) \quad (\neq Q(\theta^{(i)}, \theta^{(sh)}))$$

True joint distribution

$$\min_{\omega_i, \theta^{(i)}, \theta^{(sh)}} \sum_{j=1}^N \omega_{ij} \hat{E}L_j(z_j, \theta^{(i)}, \theta^{(sh)}) - \sum_{j=1}^N \log \omega_{ij}$$

Maximum likelihood estimate

Theoretical Guarantees

Definition: log – exp Complexity

Let $\theta^{(i)}$ and $\theta^{(sh)}$ be a family of weights corresponding to task/edge specific and shared neural networks, respectively. The log – exp complexity of the neural network with respect to the distribution $Q_{\theta^{(i)}, \theta^{(sh)}}$ (Q for short) for $i = 1, 2, \dots, N$ is defined as

$$\mathcal{R}_i(\theta) := \log \mathbb{E}_Q \sup_{\theta^{(i)}, \theta^{(sh)}} \frac{\exp \left\{ \mathbb{E}_{z \sim D_i} L_i(z, \theta^{(i)}, \theta^{(sh)}) \right\}}{\prod_{j=1}^N \hat{E}L_j(z_j, \theta^{(i)}, \theta^{(sh)})}. \quad (1)$$

Theorem: PAC bound

For a given neural network θ , and the log – exp complexity, the following bound holds with a probability of at least $1 - \delta$, ($\delta > 0$)

$$\inf_{\theta} \mathbb{E}_{z_i \sim D_i} \{L_i(z_i, \theta)\} \leq \inf_{\theta^{(sh)}} \left[\text{Obj}_i(\theta^{(sh)}) + \mathcal{R}_i(\theta) + \sup_{\theta^{(i)}, \theta^{(sh)}, \omega_i} \text{KL}(Q||P) + l_{max} \sqrt{\sum_{j=1}^N \frac{\omega_{ij}^2}{2n_j^2} \log \left(\frac{1}{\delta} \right)} \right],$$

where $\text{KL}(Q||P)$ is the KL-divergence between two joint distributions Q and P ,

$$\text{Obj}_i(\theta^{(sh)}) := \inf_{\omega_i} \sum_{j=1}^N \left[\omega_{ij} \inf_{\theta^{(i)}} \hat{E}L_j(z_j, \theta^{(i)}, \theta^{(sh)}) - \log \omega_{ij} \right]. \quad (2)$$

Algorithm

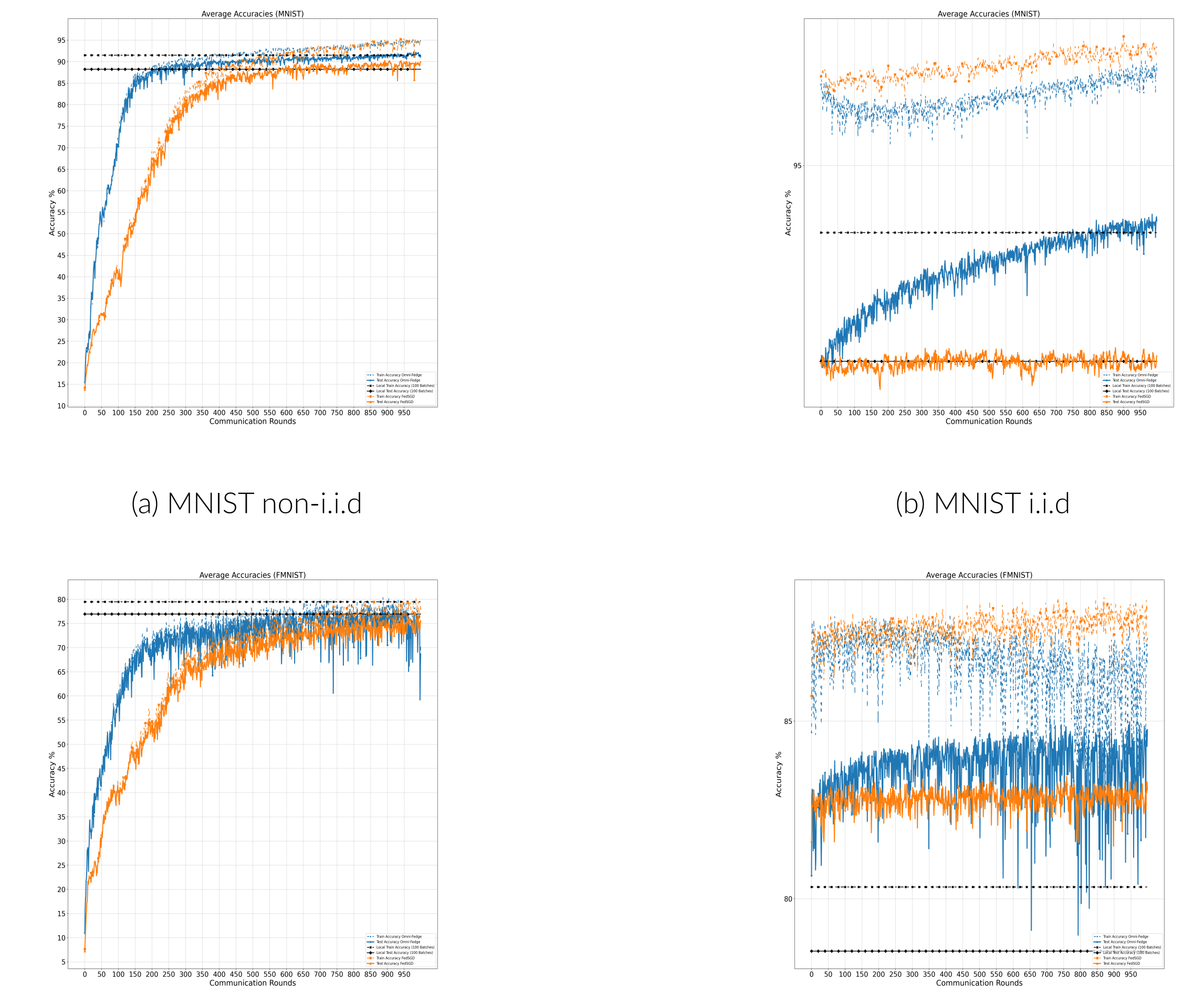
Algorithm 1: Omni-Fedge

```

1 Omni-Fedge () :
2 INITIALIZE  $\theta^{sh}$  and BROADCAST (BC) to all nodes
3 for  $t \in \{1, 2, \dots\}$  do
4   for  $i = 1, 2, \dots, N$  do
5      $\theta^{(i)} = \arg \min_{\theta^{(i)}} \hat{E}L_i(z_i, \theta^{(i)}, \theta^{sh})$ 
6     Each device  $i$  BCs  $\theta^{(i)}$  to all other nodes
7     COMPUTE AND SEND  $\hat{E}L_i(z_i, \theta^{(i)}, \theta^{sh})$ 
8     to all nodes.
9     Minimize-Objective ()
10    to get  $\omega_i$  for all  $i$ .
11    At each node, COMPUTE
12     $\sum_{j=1}^N \omega_j \nabla_{\theta^{sh}} \hat{E}L_j(z_j, \theta^{(j)}, \theta^{sh})$  and
13    BC it to all nodes through FS.
14    Perform GRADIENT UPDATE
15     $\theta_{t+1}^{sh} := \theta_t^{sh} - \eta^{comm} \gamma_t^{sh}$ , where  $\gamma_t^{sh} :=$ 
16     $\frac{1}{N} \left( \sum_{i=1}^N \sum_{j=1}^N \omega_{ij} \nabla_{\theta^{sh}} \hat{E}L_j(z_j, \theta^{(j)}, \theta^{sh}) \right)$ 
17    GO TO step 3.
18 Minimize-Objective () :
19 COMPUTE  $\omega_i^* =$ 
20  $\arg \min_{\omega_i} \left( \sum_{j=1}^N \omega_{ij} \hat{E}L_j(z_j, \theta^{(j)}, \theta^{sh}) - \log \prod_{j=1}^N \omega_{ij} \right)$ 

```

Experimental Results



(a) MNIST non-i.i.d

(b) MNIST i.i.d

(c) FMNIST non-i.i.d

(d) FMNIST i.i.d

Figure: Plots of Average Accuracies vs Communication Rounds for Omni-Fedge and FedSGD

References

- Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018.
- Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016.
- Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *Advances in Neural Information Processing Systems*, pages 527–538, 2018.