
DReS-FL: Dropout-Resilient Secure Federated Learning for Non-IID Clients via Secret Data Sharing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated learning (FL) strives to enable collaborative training of machine learning
2 models without centrally collecting clients' private data. Different from centralized
3 training, the local datasets across clients in FL are non-independent and identically
4 distributed (non-IID). In addition, the data-owning clients may drop out of the
5 training process arbitrarily. These characteristics will significantly degrade the
6 training performance. This paper proposes a *Dropout-Resilient Secure Federated*
7 *Learning* (DReS-FL) framework based on Lagrange coded computing (LCC) to
8 tackle both the non-IID and dropout problems. The key idea is to utilize Lagrange
9 coding to secretly share the private datasets among clients so that the effects of
10 non-IID distribution and client dropouts can be compensated during local gradient
11 computations. To provide a strict privacy guarantee for local datasets and correctly
12 decode the gradient at the server, the gradient has to be a polynomial function in
13 a finite field, and thus we construct *polynomial integer neural networks* (PINNs)
14 to enable our framework. Theoretical analysis shows that DReS-FL is resilient to
15 client dropouts and provides privacy protection for the local datasets. Furthermore,
16 we experimentally demonstrate that DReS-FL consistently leads to significant
17 performance gains over baseline methods.

18 1 Introduction

19 Federated learning (FL) [1] is a machine learning framework in which a central server coordinates a
20 large number of clients to collaboratively train a shared model. The key idea of FL is to train the
21 model locally by individual clients and aggregate updates globally by the server. The main target is to
22 provide privacy protection for clients' local samples and solve the "data islands" problem. However,
23 it has been shown recently that local models may reveal substantial information about the local
24 datasets, and the private training data can be reconstructed through model inversion attacks [2, 3, 4].
25 Besides, as local data are typically non-independent and identically distributed (non-IID), the model
26 divergence during the local update may lead to unstable and slow convergence [5, 6, 7]. With many
27 clients involved in the training, some of the clients could drop out of the training process unexpectedly
28 (due to poor connectivity, battery level, etc), and it will cause detrimental model performance [8].
29 Thus, effective mechanisms are needed to tackle the non-IID data distribution and client dropouts,
30 while preserving the privacy of local datasets, which motivates this work.

31 To prevent information leakage from the local models, *secure aggregation* protocols [9, 10, 11, 12,
32 13, 14, 15] have been developed to allow for global aggregation without revealing the parameters of
33 clients' models. Even if some clients may drop out, these protocols can still recover the aggregated

34 results of the surviving clients. However, their training performance may degrade significantly in
 35 the non-IID setting due to the aggregation bias. To alleviate the non-IID problem, existing methods
 36 typically follow *algorithm-based approaches* [5, 16, 17, 18, 19] and add regularization terms to
 37 mitigate the model divergence. However, these methods are not dropout-resilient evidenced by the
 38 empirical results in [20]. This can be explained by the greatly varying data distributions among
 39 different rounds. Another fold of strategy for dealing with the non-IID problem is *data-centric*
 40 *approach* [21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31], which generates extra training samples to
 41 construct a more balanced data distribution for each client. The common practices are to share the
 42 synthesized samples [23, 24, 25, 26] or GAN-based augmented data [27, 28, 29, 30, 31]. However,
 43 these methods may leak private information about local datasets and violate the privacy criterion in
 44 FL.

45 **Contributions.** In this work, we develop a *Dropout-Resilient Secure Federated Learning* (DReS-FL)
 46 framework to address the above problems via Lagrange coded computing (LCC) [32]. The key
 47 idea of LCC is to encode the datasets using Lagrange polynomials, in order to create computational
 48 redundancy across the workers in a privacy-preserving way. At the beginning of the training, the
 49 clients secretly share their private datasets with each other via Lagrange coding. This allows clients to
 50 access an encoded version of the *global dataset*¹ for local gradient computations while guaranteeing
 51 privacy in an information-theoretic sense. After collecting the computation results from a certain
 52 number of surviving clients, the server performs polynomial interpolation to decode the gradient, and
 53 thus it is resilient to client dropouts. As the local gradients are computed on the mini-batches uniformly
 54 sampled from the global dataset, the server obtains *global gradient*² after decoding. Therefore, the
 55 training process in DReS-FL is made equivalent to the centralized training and eliminates the non-IID
 56 problem. To provide a privacy guarantee for the local datasets and correctly decode the gradient at the
 57 server, the gradient has to be a polynomial function in a finite field, which is a main design challenge
 58 of DreS-FL. Our main contributions are summarized as follows:

- 59 • The proposed DReS-FL framework provides a unified approach to tackle two critical
 60 problems of FL, namely, *non-IID data distribution* and *client dropouts*. Meanwhile, it
 61 maintains privacy and security guarantees such that no information about local datasets can
 62 be leaked beyond the global model parameters.
- 63 • We construct *polynomial integer neural networks* (PINNs) to ensure that the gradient is
 64 a polynomial, so that cryptographic primitives can be applied for secure computation. A
 65 PINN consists of affine transformation layers with parameters constrained in an integer set,
 66 and it adopts the quadratic function as the activation function. The convergence analysis of
 67 DReS-FL with PINNs is also provided.
- 68 • We conduct extensive experiments on FL benchmark datasets to demonstrate the effective-
 69 ness of DReS-FL. It is shown that DReS-FL outperforms baseline methods under the setting
 70 where local datasets are heterogeneous and clients may drop out of the training process
 71 arbitrarily.

72 2 Related Works

73 **Secure aggregation.** Secure model aggregation [9, 10, 11, 13, 12, 14, 15] is a key component of
 74 FL that protects the privacy of each client’s model while allowing their global aggregation amidst
 75 possible user dropouts. Existing protocols essentially rely on two main principles, including a
 76 pairwise random-seed agreement for mask cancellation and secret sharing of the random seeds to
 77 construct the dropped masks [9, 10, 11, 12, 14, 13, 15]. However, these approaches may suffer from

¹In the context of this paper, we use the global dataset to denote the concatenation of all the clients’ datasets.

²The global gradient corresponds to the stochastic gradient computed from mini-batches that are uniformly sampled from the global dataset. For simplicity, we consider that clients only perform one local stochastic gradient descent (SGD) step in each communication round. Note that the proposed DReS-FL framework, as stated in Remark 2, can be extended to more general cases in which clients can run multiple local SGD steps.

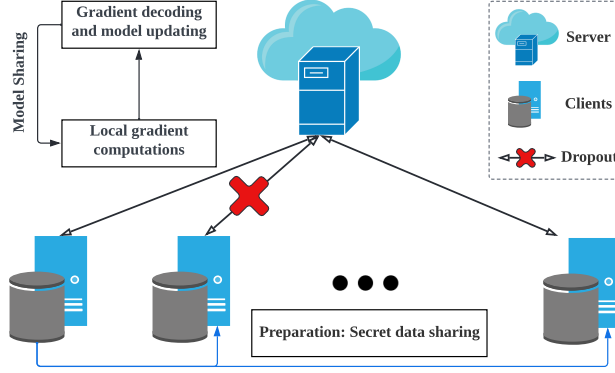


Figure 1: The DReS-FL system model. At the beginning of training, the clients secretly share the local datasets with each other. Then, the model parameters are iteratively trained by (1) local gradient computations and (2) gradient decoding and model updating until convergence.

78 severe performance degradation in non-IID FL, since the surviving clients in each round vary greatly,
 79 and thus the local gradients are biased towards different data distributions.

80 **Non-IID data and client dropouts.** Training with heterogeneous data is a unique challenge for FL
 81 [1], which significantly affects the convergence performance [8]. The client dropouts exacerbate
 82 the non-IID problem as the data distributions among different rounds could vary greatly. Many
 83 algorithm-based methods [5, 16, 17, 18, 19] attempt to mitigate the clients’ model divergence, but
 84 these methods cannot solve the essence of the non-IID problem due to the intrinsic difference between
 85 minimizing the local empirical loss and minimizing the global empirical loss. Another line of
 86 work adopts data-centric methods [27, 28, 29, 30, 31] to modify the local distributions. Ideally, a
 87 perfect data sharing mechanism should achieve that the local datasets have the same distribution
 88 as the global dataset while maintaining the privacy guarantee. Common practices include sharing
 89 raw datasets [21, 22], synthesized samples, [23, 24, 25, 26] or augmented data [27, 28, 29, 30, 31].
 90 However, these works cannot fully preserve local data privacy in an information-theoretic sense
 91 [33]. A special data-centric method is the secret coding scheme, which has been widely utilized
 92 in homomorphic encryption (HE) [34, 35, 36, 37, 38, 39, 40] and multiparty computation (MPC)
 93 techniques [41, 42, 43]. This coding scheme allows computations to be performed on encrypted
 94 data and has been used for privacy-preserving machine learning [36, 40, 43]. However, the HE
 95 methods often suffer from time-consuming cryptographic tools, and MPC techniques are difficult to
 96 generalize such primitives to a large number of clients. Recently, distributed secure machine learning
 97 frameworks [44, 45] have been proposed for logistic regression problems. They apply Lagrange
 98 coding for secret data sharing and approximate the Sigmoid function by a polynomial function. This
 99 paper proposes DReS-FL to further extend these works to train deep neural networks in the FL
 100 setting.

101 3 System Model

102 We consider a federated learning framework as shown in Fig. 1 that consists of one central server
 103 and N data-owning clients. Each client $i \in [N]$ holds a local dataset $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$ of size m_i , where
 104 $\mathbf{X}^{(i)} \in \mathbb{R}^{m_i \times d_x}$ represents the set of input features of dimension d_x and $\mathbf{Y}^{(i)} \in \mathbb{R}^{m_i \times d_y}$ corresponds
 105 to the output vector of dimension d_y . Accordingly, the size of the global dataset (\mathbf{X}, \mathbf{Y}) which
 106 concatenates all local datasets $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$, $\forall i \in [N]$ is denoted as $m \triangleq \sum_{i=1}^N m_i$. The clients aim
 107 to jointly train a neural network based on their local datasets without sharing private data samples.
 108 Particularly, the gradients are computed locally and aggregated globally. However, the local data may
 109 be highly heterogeneous, and the clients may drop out at any time unexpectedly, which makes the
 110 training process unstable. Our goal is to improve the convergence performance by secret data sharing
 111 while preserving the privacy of local datasets.

112 **Threat model and privacy requirements.** We consider a threat model where the clients are *honest-*
 113 *but-curious*. In particular, clients follow the protocol honestly but may collude amongst themselves to
 114 learn additional information. Besides, we assume that the server is also honest-but-curious, but does
 115 not collude with any other clients. To avoid leaking private information from the shared samples and
 116 protect the local updates against model inversion attacks [2, 3, 4], we impose two privacy requirements
 117 in the training process. First, the clients learn nothing about the private datasets of others from the
 118 shared samples even if up to T clients collude. Second, the server learns no information about the
 119 private datasets of clients from a single local computation result beyond the model parameters.

120 **Lagrange coded computing.** The LCC framework considers a scenario involving computations over
 121 massive datasets stored distributedly across multiple clients [32]. The key idea is to encode the data
 122 using Lagrange polynomial for redundant distributed computing, which fits nicely with federated
 123 learning due to its dropout-resiliency and privacy guarantees. DReS-FL applies LCC to secretly share
 124 the private datasets among clients for local gradient computations, and the global gradient can be
 125 decoded by the server for model updating. To provide a strong privacy guarantee for the datasets and
 126 correctly decode the gradient at the server, the gradient should be a polynomial function in a finite
 127 field. However, existing neural networks cannot satisfy this requirement, since the datasets are in the
 128 real field and the gradients are not polynomial functions due to the non-linear operations.

129 **Polynomial Integer Neural Networks.** We define a class of polynomial integer neural networks
 130 (PINNs) to ensure that the gradient is a polynomial function in a finite field \mathbb{F}_p with a prime number
 131 p . First, we transform the dataset (\mathbf{X}, \mathbf{Y}) from the real domain to the finite domain $(\bar{\mathbf{X}}, \bar{\mathbf{Y}})$. Besides,
 132 a PINN consists of affine transformation layers (e.g., fully connected layers and convolutional layers)
 133 and utilizes the quadratic function as the activation function. The model parameters of PINNs are
 134 defined in the integer set $\mathbb{Z}_p \triangleq \{-\lfloor \frac{p+1}{2} \rfloor, \dots, \lfloor \frac{p-1}{2} \rfloor\}$. Given a feed-forward function $f(\bar{\mathbf{X}}; \mathbf{w})$ and
 135 selecting the mean squared error (MSE) as the loss function, the gradient of the input samples is a
 136 multivariate polynomial with integer coefficients, i.e., $\mathbf{g}(\bar{\mathbf{X}}, \bar{\mathbf{Y}}; \mathbf{w}) \triangleq \nabla_{\mathbf{w}} \|\bar{\mathbf{Y}} - f(\bar{\mathbf{X}}; \mathbf{w})\|_2^2 \in \mathbb{Z}^{d_w}$,
 137 where d_w represents the number of model parameters. In particular, to avoid wrap-around when
 138 computing gradient in the finite field \mathbb{F}_p , we assume the prime number p is sufficiently large without
 139 leading to overflow errors in the integer set \mathbb{Z}_p .

140 4 The Proposed DReS-FL Framework

141 DReS-FL consists of two main phases, as shown in Fig. 1. In the first phase, the private datasets are
 142 transformed from the real domain to the finite field, and data-owning clients secretly share datasets
 143 by Lagrange coding. Then, the server and the clients train a PINN iteratively via (1) local gradient
 144 computations and (2) gradient decoding and model updating.

145 4.1 Data Transformation and Secret Sharing

146 To guarantee information-theoretic privacy, each client has to mask the datasets in a finite field \mathbb{F}_p
 147 using uniformly random matrices. Firstly, the local datasets $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$ are converted from the
 148 real domain to the finite field $(\bar{\mathbf{X}}^{(i)}, \bar{\mathbf{Y}}^{(i)})$. Considering an element-wise function $\phi(z) = z + c$
 149 that transforms a real value to a non-negative number by adding a proper scalar c^3 , we define
 150 $\bar{\mathbf{X}} \triangleq \text{Round}(2^l \cdot \phi(\mathbf{X}))$, where the rounding operation is element-wise that quantizes each entry to
 151 its closest integer, and $l \in \mathbb{Z}$ controls the quantization loss. We adopt the notation $\bar{\mathcal{D}}$ to represent the
 152 *global dataset*, which is the concatenation of all the local datasets $(\bar{\mathbf{X}}^{(i)}, \bar{\mathbf{Y}}^{(i)})$ for $i \in [N]$.

153 After converting the private datasets to the finite field, clients adopt T -private Lagrange coding [32]
 154 to secretly share their local data samples across other clients, where $T \geq 1$ is the privacy parameter
 155 in our system ensuring that the encoded datasets do not leak any information about the original
 156 datasets even if T clients collude. First, each client $i \in [N]$ partitions its local dataset to K shards
 157 as $\bar{\mathbf{X}}^{(i)} \triangleq [\bar{\mathbf{X}}_1^{(i)T}, \dots, \bar{\mathbf{X}}_K^{(i)T}]^T$ and $\bar{\mathbf{Y}}^{(i)} \triangleq [\bar{\mathbf{Y}}_1^{(i)T}, \dots, \bar{\mathbf{Y}}_K^{(i)T}]^T$. Assuming that m_i is divisible

³The scalar c could be the absolute value of the minimum entry in dataset, which is set to 0 in the experiments.

Table 1: Primary notations and descriptions

Notation	Description	Notation	Description
$(\bar{\mathbf{X}}^{(i)}, \bar{\mathbf{Y}}^{(i)})$	Transformed dataset at client i over finite field \mathbb{F}_p	$(\tilde{\mathbf{X}}_j^{(i)}, \tilde{\mathbf{Y}}_j^{(i)})$	Encoded dataset sent from client i to client j
$\bar{\mathcal{D}}$	Concatenation of datasets $(\bar{\mathbf{X}}^{(i)}, \bar{\mathbf{Y}}^{(i)})$ for $i \in [N]$	$(\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j)$	Concatenation of all the received encoded datasets at client j
$(\bar{\mathbf{X}}_k^{(i)}, \bar{\mathbf{Y}}_k^{(i)})$	k -th data shard at client i	$\mathbf{C}^{(t)}$	row selection matrix for data sampling in round t
$\bar{\mathcal{D}}_k^{(t)} \triangleq (\bar{\mathbf{X}}_k^{(t)}, \bar{\mathbf{Y}}_k^{(t)})$	k -th global mini-batch sampled from $\bar{\mathcal{D}}$ in round t	$(\tilde{\mathbf{X}}_j^{(t)}, \tilde{\mathbf{Y}}_j^{(t)})$	mini-batch sampled from $(\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j)$ at client j in round t

Algorithm 1 DReS-FL

Input: Local datasets $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$ for $i \in [N]$, batch size b , initialized parameters $\mathbf{w}^{(0)} \in \mathbb{Z}_p^{d_w}$, distinct elements $\{\alpha_j\}_{j \in [N]}$ and $\{\beta_k\}_{k \in [K+T]}$, prime number p , training round τ , learning rate η .

Output: Model parameter $\mathbf{w}^{(\tau)}$.

- 1: Clients encode the local datasets according to (1) and (2) and deliver them to other clients.
- 2: **for** $t = 1, 2, \dots, \tau$ **do**
- 3: Server sends the model parameters $\mathbf{w}^{(t)}$ to the clients.
- 4: **for** $j = 1, \dots, N$ **do**
- 5: Client j performs gradient computation on mini-batches $(\tilde{\mathbf{X}}_j^{(t)}, \tilde{\mathbf{Y}}_j^{(t)})$.
- 6: Upload local computation results $\tilde{\mathbf{g}}(\tilde{\mathbf{X}}_j^{(t)}, \tilde{\mathbf{Y}}_j^{(t)}; \mathbf{w}^{(t)})$ to the server.
- 7: **end for**
- 8: **if** Server receives at least $\deg(\mathbf{g})(K+T-1)+1$ uploads **then**
- 9: Decode K stochastic gradients $\tilde{\mathbf{g}}(\bar{\mathcal{D}}_k^{(t)}; \mathbf{w}^{(t)})$ for $k \in [K]$ by polynomial interpolation.
- 10: Convert gradients from the finite field to the integral domain $\mathbf{g}(\bar{\mathcal{D}}_k^{(t)}; \mathbf{w}^{(t)})$ by (5).
- 11: Update the global model by $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - Q(\frac{\eta}{bK} \sum_{j=1}^K \mathbf{g}(\bar{\mathcal{D}}_k^{(t)}; \mathbf{w}^{(t)}))$ based on (6).
- 12: **end if**
- 13: **end for**

158 by K , we have $\bar{\mathbf{X}}_k^{(i)} \in \mathbb{F}_p^{\frac{m_i}{K} \times d_x}$ and $\bar{\mathbf{Y}}_k^{(i)} \in \mathbb{F}_p^{\frac{m_i}{K} \times d_y}$ for $k \in [K]$. A large value of K helps to
 159 reduce the communication overhead in secret data sharing and computation costs in local gradient
 160 computations⁴. Then, the clients add padding from T uniform random masks to the data samples
 161 for privacy protection. Each client $i \in [N]$ forms the following polynomials $\mathbf{u}_i : \mathbb{F}_p \rightarrow \mathbb{F}_p^{\frac{m_i}{K} \times d_x}$ and
 162 $\mathbf{v}_i : \mathbb{F}_p \rightarrow \mathbb{F}_p^{\frac{m_i}{K} \times d_y}$ of degree $K+T-1$ to encode the local dataset:

$$\mathbf{u}_i(z) \triangleq \sum_{k \in [K]} \bar{\mathbf{X}}_k^{(i)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j} + \sum_{k=K+1}^{K+T} \mathbf{U}_k^{(i)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j}, \quad (1)$$

$$\mathbf{v}_i(z) \triangleq \sum_{k \in [K]} \bar{\mathbf{Y}}_k^{(i)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j} + \sum_{k=K+1}^{K+T} \mathbf{V}_k^{(i)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j}, \quad (2)$$

163 where $\{\mathbf{U}_k^{(i)}\}$'s and $\{\mathbf{V}_k^{(i)}\}$'s are random noise matrices uniformly sampled from $\mathbb{F}_p^{\frac{m_i}{K} \times d_x}$ and $\mathbb{F}_p^{\frac{m_i}{K} \times d_y}$,
 164 respectively. These matrices mask the local datasets and provide a privacy guarantee against up to T
 165 colluding workers. The clients and the server agree on $K+T$ distinct elements $\{\beta_1, \dots, \beta_{K+T}\}$
 166 from the finite field \mathbb{F}_p in advance. Particularly, setting $z = \beta_k$ for $k \in [K]$, we obtain $\mathbf{u}_i(\beta_k) = \bar{\mathbf{X}}_k^{(i)}$
 167 and $\mathbf{v}_i(\beta_k) = \bar{\mathbf{Y}}_k^{(i)}$. All the clients use the same N distinct elements $\{\alpha_1, \dots, \alpha_N\}$ selected from
 168 \mathbb{F}_p to encode the private datasets, where $\{\alpha_j\}_{j \in [N]} \cap \{\beta_k\}_{k \in [K+T]} = \emptyset$. Each client i obtains N
 169 encoded datasets $(\tilde{\mathbf{X}}_j^{(i)}, \tilde{\mathbf{Y}}_j^{(i)}) \triangleq (\mathbf{u}_i(\alpha_j), \mathbf{v}_i(\alpha_j))$ for $j \in [N]$, where $(\tilde{\mathbf{X}}_j^{(i)}, \tilde{\mathbf{Y}}_j^{(i)})$ is sent to client
 170 j from client i . All the received encoded datasets at client j are represented as $(\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j)$, where
 171 $\tilde{\mathbf{X}}_j \triangleq [\tilde{\mathbf{X}}_j^{(1)T}, \dots, \tilde{\mathbf{X}}_j^{(N)T}]^T \in \mathbb{F}_p^{\tilde{m} \times d_x}$ and $\tilde{\mathbf{Y}}_j \triangleq [\tilde{\mathbf{Y}}_j^{(1)T}, \dots, \tilde{\mathbf{Y}}_j^{(N)T}]^T \in \mathbb{F}_p^{\tilde{m} \times d_y}$ for $j \in [N]$.
 172 Accordingly, the number of samples in the encoded dataset is $\tilde{m} \triangleq \frac{1}{K} \sum_{i=1}^n m_i$.

⁴The complexity analysis of the proposed method has been deferred to Section C in Appendix.

173 **4.2 Federated Training**

174 **Local Gradient Computation.** The server randomly initializes a PINN at the beginning of the
 175 training process, and the model parameters are constrained to an integer set \mathbb{Z}_p in the training
 176 process. In each communication round, the server sends the model parameters to the clients, and
 177 they compute the stochastic gradient over the mini-batches with size b . Particularly, we assume
 178 that all the clients use the same row selection matrix $\mathbf{C}^{(t)} \in \{0, 1\}^{b \times \tilde{m}}$ for data sampling in each
 179 round t^5 , and the mini-batch at each client $j \in [N]$ is determined by $[\tilde{\mathbf{X}}_j^{(\mathcal{I}_t)}, \tilde{\mathbf{Y}}_j^{(\mathcal{I}_t)}] = \mathbf{C}^{(t)}[\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j]$.
 180 Here, $\mathcal{I}_t = \{l_1^{(t)}, \dots, l_b^{(t)}\} \subseteq [\tilde{m}]$ is a randomly selected index set in the t -th round with $l_i \in [\tilde{m}]$
 181 for $i \in [b]$. The entries of $\mathbf{C}^{(t)}$ satisfy $\mathbf{C}_{i, l_i}^{(t)} = 1$ for $i \in [b]$, and other entries are set to zero. Each
 182 client j computes the stochastic gradient $\tilde{\mathbf{g}}(\tilde{\mathbf{X}}_j^{(\mathcal{I}_t)}, \tilde{\mathbf{Y}}_j^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) \equiv \mathbf{g}(\tilde{\mathbf{X}}_j^{(\mathcal{I}_t)}, \tilde{\mathbf{Y}}_j^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) \pmod{p}$
 183 in the finite field, and uploads the result to the server. Particularly, each $\tilde{\mathbf{g}}(\tilde{\mathbf{X}}_j^{(\mathcal{I}_t)}, \tilde{\mathbf{Y}}_j^{(\mathcal{I}_t)}; \mathbf{w}^{(t)})$
 184 amounts to an evaluation of the polynomial $\tilde{\mathbf{g}}(\mathbf{u}_{\mathcal{I}_t}(z), \mathbf{v}_{\mathcal{I}_t}(z); \mathbf{w}^{(t)})$ at the point $z = \alpha_j$, where two
 185 $(K + T - 1)$ -degree polynomial functions $\mathbf{u}_{\mathcal{I}_t} : \mathbb{F}_p \rightarrow \mathbb{F}_p^{b \times d_x}$ and $\mathbf{v}_{\mathcal{I}_t} : \mathbb{F}_p \rightarrow \mathbb{F}_p^{b \times d_y}$ are defined as
 186 follows:

$$\mathbf{u}_{\mathcal{I}_t}(z) \triangleq \sum_{k \in [K]} \overline{\mathbf{X}}_k^{(\mathcal{I}_t)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j} + \sum_{k=K+1}^{K+T} \mathbf{U}_k^{(\mathcal{I}_t)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j}, \quad (3)$$

$$\mathbf{v}_{\mathcal{I}_t}(z) \triangleq \sum_{k \in [K]} \overline{\mathbf{Y}}_k^{(\mathcal{I}_t)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j} + \sum_{k=K+1}^{K+T} \mathbf{V}_k^{(\mathcal{I}_t)} \cdot \prod_{j \in [K+T] \setminus \{k\}} \frac{z - \beta_j}{\beta_k - \beta_j}, \quad (4)$$

187 with

$$\begin{aligned} \overline{\mathbf{X}}_k^{(\mathcal{I}_t)} &= \mathbf{C}^{(t)}[\overline{\mathbf{X}}_k^{(1)T}, \dots, \overline{\mathbf{X}}_k^{(N)T}]^T \in \mathbb{F}_p^{b \times d_x}, & \overline{\mathbf{U}}_k^{(\mathcal{I}_t)} &= \mathbf{C}^{(t)}[\mathbf{U}_k^{(1)T}, \dots, \mathbf{U}_k^{(N)T}]^T \in \mathbb{F}_p^{b \times d_x}, \\ \overline{\mathbf{Y}}_k^{(\mathcal{I}_t)} &= \mathbf{C}^{(t)}[\overline{\mathbf{Y}}_k^{(1)T}, \dots, \overline{\mathbf{Y}}_k^{(N)T}]^T \in \mathbb{F}_p^{b \times d_y}, & \overline{\mathbf{V}}_k^{(\mathcal{I}_t)} &= \mathbf{C}^{(t)}[\mathbf{V}_k^{(1)T}, \dots, \mathbf{V}_k^{(N)T}]^T \in \mathbb{F}_p^{b \times d_y}. \end{aligned}$$

188 Every $(\overline{\mathbf{X}}_k^{(\mathcal{I}_t)}, \overline{\mathbf{Y}}_k^{(\mathcal{I}_t)}) = (\mathbf{u}_{\mathcal{I}_t}(\beta_k), \mathbf{v}_{\mathcal{I}_t}(\beta_k))$ for $k \in [K]$ is a *global mini-batch* selected from the
 189 global dataset $\overline{\mathcal{D}}$. For notational simplicity, we define $\overline{\mathcal{D}}_k^{(\mathcal{I}_t)} \triangleq (\overline{\mathbf{X}}_k^{(\mathcal{I}_t)}, \overline{\mathbf{Y}}_k^{(\mathcal{I}_t)})$.

190 **Gradient Decoding and Model Updating.** After receiving any $\deg(\mathbf{g})(K + T - 1) + 1$ local com-
 191 putation results⁶, the server decodes the stochastic gradients $\tilde{\mathbf{g}}(\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) \triangleq \tilde{\mathbf{g}}(\overline{\mathbf{X}}_k^{(\mathcal{I}_t)}, \overline{\mathbf{Y}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)})$
 192 for $k \in [K]$ by interpolating the polynomial $\tilde{\mathbf{g}}(\mathbf{u}_{\mathcal{I}_t}(z), \mathbf{v}_{\mathcal{I}_t}(z); \mathbf{w}^{(t)})$. As the polynomial function
 193 $\tilde{\mathbf{g}}(\mathbf{u}_{\mathcal{I}_t}(z), \mathbf{v}_{\mathcal{I}_t}(z); \mathbf{w}^{(t)})$ is a composition of the encoding polynomials and the gradient function, its
 194 degree is $\deg(\mathbf{g})(K + T - 1)$. Therefore, the server needs at least $\deg(\mathbf{g})(K + T - 1) + 1$ evaluations
 195 to interpolate it. After recovering the coefficients of the polynomial $\tilde{\mathbf{g}}(\mathbf{u}_{\mathcal{I}_t}(z), \mathbf{v}_{\mathcal{I}_t}(z); \mathbf{w}^{(t)})$, the
 196 server decodes K stochastic gradients $\tilde{\mathbf{g}}(\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) = \tilde{\mathbf{g}}(\mathbf{u}_{\mathcal{I}_t}(\beta_k), \mathbf{v}_{\mathcal{I}_t}(\beta_k); \mathbf{w}^{(t)})$ on the global
 197 mini-batches $\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}$ by letting $z = \beta_k$ for $k \in [K]$. Then, the server converts the gradient from the fi-
 198 nite field to the integer set \mathbb{Z}_p by $\mathbf{g}(\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) = \psi(\tilde{\mathbf{g}}(\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}))$, where $\psi(z)$ is an element-wise
 199 function defined as follows:

$$\psi(z) = \begin{cases} z & \text{if } 0 \leq z < \frac{p-1}{2}, \\ z - p & \text{if } \frac{p-1}{2} \leq z < p. \end{cases} \quad (5)$$

200 As we assume that the prime number p is sufficiently large, the converted gradients do not
 201 have overflow errors. Thus, the central sever updates the global model by $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} -$
 202 $Q(\frac{\eta}{bK} \sum_{k=1}^K \mathbf{g}(\overline{\mathcal{D}}_k^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}))$, where η denotes the learning rate and bK represents the *global batch*

⁵This can be achieved by setting the same random seed across all the clients. The weighted sampling method has been adopted in this work, where the number of sampled data from $\tilde{\mathbf{X}}_j^{(i)}$ is proportional to m_i for $i \in [N]$. Note that other sampling schemes can also be applied in DReS-FL.

⁶Note that if the server cannot receive enough results due to the client dropouts, the training protocol continues to the next epoch without gradient decoding and model updating.

203 $size^7$. $Q(z)$ is a stochastic quantization function to ensure the model parameters are in the integer set
 204 \mathbb{Z}_p after updating, which is defined as follows:

$$Q(z) = \begin{cases} \lfloor z \rfloor & \text{with probability } 1 - (z - \lfloor z \rfloor) \\ \lfloor z \rfloor + 1 & \text{with probability } z - \lfloor z \rfloor. \end{cases} \quad (6)$$

205 Besides, the probability of rounding z to $\lfloor z \rfloor$ is proportional to the proximity of z to $\lfloor z \rfloor$ so that the
 206 stochastic rounding is unbiased. The overall procedure is summarized in Algorithm 1.

207 5 Theoretical Analysis

208 5.1 Privacy and Security Guarantees

209 Before training starts, client j receives an encoded version of the global dataset $(\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j)$ from other
 210 clients. Lagrange coding in DReS-FL provides a strong privacy guarantee that the clients cannot infer
 211 anything about the private datasets based on the received datasets even if up to T clients collude. The
 212 following theorem shows that our DReS-FL satisfies the first privacy requirement in Section 3, and
 213 the proof is available in Section IV of [32].

214 **Theorem 1.** (*T-private coding scheme [32]*) *Employing Lagrange coding with the privacy parameter*
 215 *T in DReS-FL, we have that for every subset of clients $\mathcal{T} \subseteq [N]$ of size at most T, the mutual*
 216 *information $I(\overline{\mathcal{D}}; \{\tilde{\mathbf{X}}_j, \tilde{\mathbf{Y}}_j\}_{j \in \mathcal{T}}) = 0$.*

217 Besides, DReS-FL can provide a security guarantee such that the server learns no information from
 218 the local gradients beyond the global model \mathbf{w} . This property corresponds to the second privacy
 219 requirement in Section 3. The following theorem provides a rigorous statement and its proof is
 220 deferred to Appendix B.1.

221 **Theorem 2.** *For any $i, j \in [N]$ and index set \mathcal{I}_t , the conditional mutual information*
 222 *$I(\overline{\mathbf{X}}^{(i)}, \overline{\mathbf{Y}}^{(i)}; \tilde{\mathbf{g}}(\tilde{\mathbf{X}}_j^{(\mathcal{I}_t)}, \tilde{\mathbf{Y}}_j^{(\mathcal{I}_t)}; \mathbf{w}^{(t)}) | \mathbf{w}^{(t)})$ equals to zero.*

223 5.2 Dropout-resiliency and Convergence

224 In the FL setting, it is common for clients to drop out at any time during protocol execution, which
 225 leads to model divergence especially when the clients' datasets are highly heterogeneous. The
 226 following theorem shows that DReS-FL is resilient to a certain number of client dropouts.

227 **Theorem 3.** (*Dropout-resiliency*) *Consider N clients in the federated learning system that*
 228 *use a T-private coding scheme to secretly share the local samples $[\overline{\mathbf{X}}_1^{(i)T}, \dots, \overline{\mathbf{X}}_K^{(i)T}]^T$ and*
 229 *$[\overline{\mathbf{Y}}_1^{(i)T}, \dots, \overline{\mathbf{Y}}_K^{(i)T}]^T$ for $i \in [N]$ for local gradient computations. DReS-FL guarantees that the*
 230 *server can decode the global gradient when there are no more than $D = N - \deg(\mathbf{g})(K + T - 1) - 1$*
 231 *client dropouts.*

232 *Proof.* Theorem 1 of [32] shows that given a number of N computing nodes and a K -shard dataset,
 233 the LCC framework provides a T -private coding scheme for computing any polynomial \mathbf{g} , as long as
 234 $\deg(\mathbf{g})(K + T - 1) + 1 \leq N$. Thus, it can tolerate at most $N - \deg(\mathbf{g})(K + T - 1) - 1$ dropouts. \square

235 **Remark 1.** There is a tradeoff among privacy guarantee (T), gradient computation cost ($1/K$),
 236 and dropout-resilience (D). Parameter T reflects the privacy threshold of Lagrange coding, and
 237 parameter K accounts for the computation load reduction. In particular, the local batch size of
 238 each client is $1/K$ of the global batch size. DReS-FL can achieve any T , K , and D as long as
 239 $D \deg(\mathbf{g})(K + T - 1) + 1 \leq N$. As T and K increase, DReS-FL can tolerate fewer client dropouts.

240 **Remark 2.** Our DReS-FL framework can be extended to more general cases in which clients can
 241 run s ($s \geq 1$) local SGD steps each round. More discussion has been deferred to Appendix D.

⁷The global batch size corresponds to the number of samples used for gradient computations in each round.

Table 2: Test accuracy (%) of different methods. Each experiment is repeated five times.

Dataset	MNIST	Fashion-MNIST	EMNIST	CIFAR-10	CIFAR-100	SVHN
FedAvg	96.17 ± 0.05	81.20 ± 0.07	71.50 ± 0.28	89.54 ± 0.09	67.71 ± 0.26	83.82 ± 0.20
FedAvg-IS	97.06 ± 0.10	85.94 ± 0.16	77.09 ± 0.34	89.83 ± 0.07	68.92 ± 0.14	85.27 ± 0.09
SCAFFOLD	71.89 ± 3.92	55.22 ± 1.83	55.15 ± 5.95	54.17 ± 9.13	29.97 ± 1.73	51.27 ± 3.43
DReS-FL (Ours)	97.38 ± 0.08	86.60 ± 0.32	78.04 ± 0.29	90.31 ± 0.19	69.15 ± 0.27	86.04 ± 0.15
Centralized	97.99 ± 0.04	89.02 ± 0.11	82.45 ± 0.23	90.37 ± 0.12	71.12 ± 0.09	86.18 ± 0.03

Next, we characterize the convergence performance of PINNs, which relies on the fact that the global gradients in the training process are unbiased. Define the empirical risk as $\ell(\mathbf{w}) \triangleq \mathbb{E}_{(\bar{\mathbf{X}}, \bar{\mathbf{Y}}) \sim \bar{\mathcal{D}}} \|\bar{\mathbf{Y}} - \mathbf{f}(\bar{\mathbf{X}}; \mathbf{w})\|_2^2$ and the corresponding gradient as $\mathbf{g}_e(\mathbf{w}) \triangleq \mathbb{E}_{(\bar{\mathbf{X}}, \bar{\mathbf{Y}}) \sim \bar{\mathcal{D}}} [\mathbf{g}(\bar{\mathbf{X}}, \bar{\mathbf{Y}}; \mathbf{w})]$. The variables $(\bar{\mathbf{X}}, \bar{\mathbf{Y}})$ are drawn from the distribution of the global dataset $\bar{\mathcal{D}}$. To prove that DReS-FL guarantees convergence to the optimal model parameters, we first present the following amputations to facilitate the analysis.

Assumption 1. (*L-smoothness*) There exists a constant $L > 0$ such that for all $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{Z}_p^{d_w}$, we have $\|\mathbf{g}_e(\mathbf{w}_1) - \mathbf{g}_e(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2$.

Assumption 2. (*Unbiased and variance-bounded stochastic gradient*) There exists a constant $\sigma > 0$ such that any stochastic gradient $\mathbf{g}(\bar{\mathcal{D}}_k; \mathbf{w}^{(t)})$ satisfies $\mathbb{E}[\frac{1}{b}\mathbf{g}(\bar{\mathcal{D}}_k; \mathbf{w}^{(t)})] = \mathbf{g}_e(\mathbf{w}^{(t)})$ and $\mathbb{E}[\|\frac{1}{b}\mathbf{g}(\bar{\mathcal{D}}_k; \mathbf{w}^{(t)}) - \mathbf{g}_e(\mathbf{w}^{(t)})\|_2^2] \leq \sigma^2$.

Assumption 3. (*Unbiased and variance-bounded rounding operation*) There exists a constant $\gamma > 0$ such that for any $z \in \mathbb{R}$, the stochastic quantization operation $Q(\cdot)$ satisfies $\mathbb{E}[Q(z)] = z$ and $\mathbb{E}[\|Q(z) - z\|_2^2] \leq \gamma^2 z^2$.

With the above preparations, we have the following theorem which ensures the convergence. The proof is deferred to Appendix B.2.

Theorem 4. (*Convergence*) Denote \mathbf{w}^* as the first-order optimal solution. With Assumption 1-3, selecting the learning rate as $\eta = \mathcal{O}\left(\frac{1}{\sqrt{\tau'}}\right)$ such that $\Psi \triangleq 1 - \eta L/2 - \eta \gamma^2 L/2 > 0$, after τ' times of model updates, we have:

$$\frac{1}{\tau'} \sum_{t=1}^{\tau'} \mathbb{E}[\|\mathbf{g}_e(\mathbf{w}^{(t)})\|_2^2] \leq \frac{\ell(\mathbf{w}^{(0)}) - \ell(\mathbf{w}^*)}{\eta \tau' \Psi} + \frac{\eta^2 L \sigma^2}{2bK\Psi} (\gamma^2 + 1). \quad (7)$$

6 Experiments

6.1 Experimental Setup

Dataset. We evaluate our proposed algorithm on several benchmark datasets: MNIST [46], Fashion-MNIST [47], EMNIST (Balanced) [48], CIFAR-10 [49], CIFAR-100 [49], and SVHN [50]. To simulate the non-IID data distribution, we assume there are $N = 20$ clients in the learning system and adopt the skewed label partition [51] to shuffle the datasets. Specifically, we sort a dataset by the labels, divide it into N shards, and assign one shard to each client. To simulate the client dropouts in the training process, we consider an extreme scenario, where the dropout rate of each client is set to 0.99 with a probability of 0.5 or is uniformly sampled from $[0, 0.1]$ otherwise.

Model structures. We adopt a multi-layer perception (MLP) with two hidden layers for the image classification tasks on MNIST, Fashion-MNIST, and EMNIST datasets. Besides, we utilize an ImageNet pretrained VGG19 model [52, 53] for CIFAR-10, CIFAR-100, and SVHN datasets. Specifically, the parameters in convolutional layers of VGG19 are fixed, and we utilize an MLP with two hidden layers as a classifier. The baseline methods train the neural networks on the real field and select the rectified linear unit (ReLU) function as the activation function. In each communication round, clients perform one SGD step for the local model update.

DReS-FL. Our method adopts the same size PINNs to replace MLPs in the federated training, and the degree of gradient is $\deg(\mathbf{g}) = 8$. Particularly, the extracted features from the last convolutional layer of VGG19 are secretly shared with other clients. We set the parameters $K = 1$ and $T = 1$ in the Lagrange coding, and the minimum number of clients needed to decode the global gradient is 9.

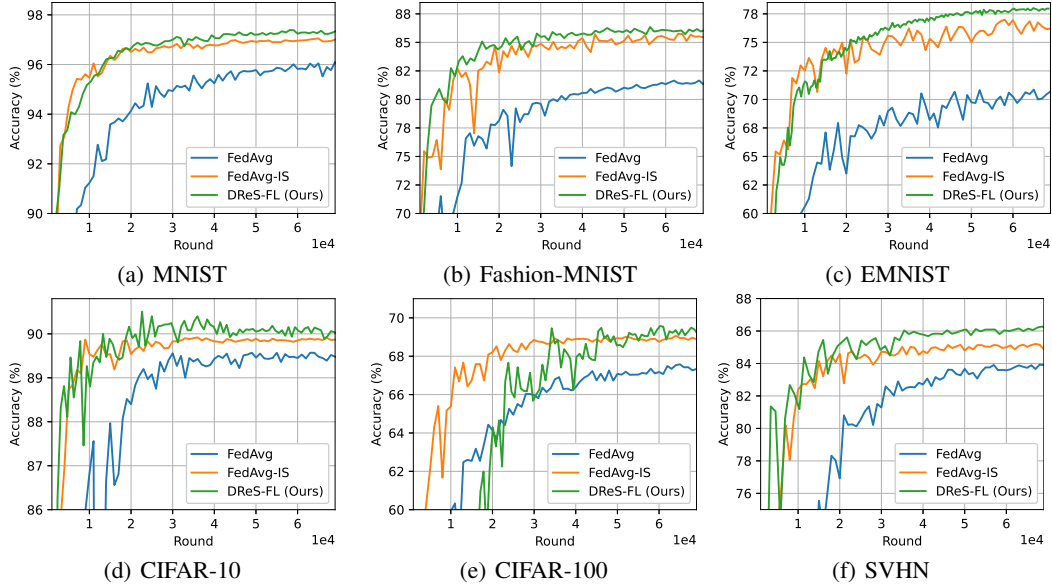


Figure 2: Test accuracies on different datasets.

281 **Baselines.** We select algorithm-based methods as baselines, including FedAvg [1], FedAvg with
 282 importance sampling (FedAvg-IS) [54, 55], and SCAFFOLD [5]. Specifically, we assume that the
 283 FedAvg-IS method knows the dropout distribution, and the local updates are weighted by the dropout
 284 probabilities to mitigate bias. Besides, we adopt the centralized training scheme as a performance
 285 upper bound in the comparison. More details of the experimental settings are deferred to Appendix
 286 A.

287 6.2 Performance Evaluation

288 We compare the performance of our DReS-FL method with base-
 289 lines. The experimental results are shown in Table 2 and Fig. 2.
 290 **The FedAvg method achieves worse performance than the central-
 291 ized training scheme. This is attributed to the non-IID data and
 292 client dropouts.** The FedAvg-IS method improves the test accuracy
 293 compared with FedAvg, but there is still a noticeable performance
 294 gap with the centralized training scheme. It shows that using the
 295 knowledge of dropout distribution can partially compensate for the
 296 biases in the aggregated models, but the local data distributions are
 297 still heterogeneous and degrade the performance. Besides, SCAF-
 298 FOLD has a low accuracy on all the settings. As the frequency of
 299 updating local control variates is low, the estimation of the update direction is highly inaccurate such
 300 that the model does not converge as shown in Fig. 3. These results are consistent with the findings in
 301 [20]. Our DReS-FL method is superior to all the baseline methods as the server can obtain global
 302 gradients after polynomial interpolation. In addition, DReS-FL achieves comparable performance
 303 to the centralized training scheme on some datasets, which demonstrates the effectiveness of our
 304 proposed framework in solving the non-IID and dropout problems.

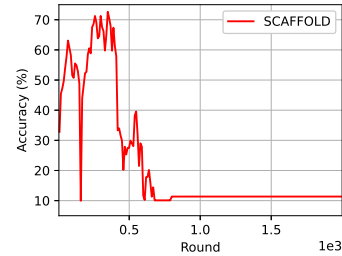


Figure 3: Test accuracy of SCAFFOLD on MNIST.

305 7 Conclusions

306 This paper proposed a Dropout-Resilient Secure Federated Learning (DReS-FL) framework via
 307 Lagrange coded computing (LCC) to simultaneously solve the data heterogeneity and dropout
 308 problems of FL, while providing privacy guarantees for the local datasets. The polynomial integer
 309 neural networks (PINNs) have been constructed to ensure that the server can correctly decode the
 310 global gradient without privacy leakage. Extensive experimental results validated the effectiveness of
 311 the proposed method. Potential limitations of our method have been deferred to Appendix E.

312 References

- 313 [1] McMahan, B., E. Moore, D. Ramage, et al. Communication-efficient learning of deep networks
314 from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- 315 [2] Shokri, R., M. Stronati, C. Song, et al. Membership inference attacks against machine learning
316 models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017.
- 317 [3] Nasr, M., R. Shokri, A. Houmansadr. Comprehensive privacy analysis of deep learning. In
318 *Proceedings of the 2019 IEEE Symposium on Security and Privacy (SP)*, pages 1–15. 2018.
- 319 [4] Geiping, J., H. Bauermeister, H. Dröge, et al. Inverting gradients-how easy is it to break privacy
320 in federated learning? *Advances in Neural Information Processing Systems*, 33:16937–16947,
321 2020.
- 322 [5] Karimireddy, S. P., S. Kale, M. Mohri, et al. Scaffold: Stochastic controlled averaging for
323 federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR,
324 2020.
- 325 [6] Kairouz, P., H. B. McMahan, B. Avent, et al. Advances and open problems in federated learning.
326 *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- 327 [7] Hsieh, K., A. Phanishayee, O. Mutlu, et al. The non-iid data quagmire of decentralized machine
328 learning. In *International Conference on Machine Learning*, pages 4387–4398. PMLR, 2020.
- 329 [8] Luo, M., F. Chen, D. Hu, et al. No fear of heterogeneity: Classifier calibration for federated
330 learning with non-iid data. *Advances in Neural Information Processing Systems*, 34, 2021.
- 331 [9] Bonawitz, K. and Ivanov, Vladimir and Kreuter, Ben and Marcedone, Antonio and McMahan,
332 H Brendan and Patel, Sarvar and Ramage, Daniel and Segal, Aaron and Seth, Karn. Practical
333 secure aggregation for privacy-preserving machine learning. In *proceedings of the 2017 ACM*
334 *SIGSAC Conference on Computer and Communications Security*, pages 1175–1191. 2017.
- 335 [10] Bonawitz, K., V. Ivanov, B. Kreuter, et al. Practical secure aggregation for privacy-preserving
336 machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer and*
337 *Communications Security*, pages 1175–1191. 2017.
- 338 [11] So, J., B. Güler, A. S. Avestimehr. Turbo-aggregate: Breaking the quadratic aggregation
339 barrier in secure federated learning. *IEEE Journal on Selected Areas in Information Theory*,
340 2(1):479–489, 2021.
- 341 [12] Kadhe, S., N. Rajaraman, O. O. Koyluoglu, et al. Fastsecagg: Scalable secure aggregation for
342 privacy-preserving federated learning. *arXiv preprint arXiv:2009.11248*, 2020.
- 343 [13] Bell, J. H., K. A. Bonawitz, A. Gascón, et al. Secure single-server aggregation with (poly)
344 logarithmic overhead. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and*
345 *Communications Security*, pages 1253–1269. 2020.
- 346 [14] Yang, C.-S., J. So, C. He, et al. Lightsecagg: Rethinking secure aggregation in federated
347 learning. *arXiv preprint arXiv:2109.14236*, 2021.
- 348 [15] Jahani-Nezhad, T., M. A. Maddah-Ali, S. Li, et al. Swiftagg+: Achieving asymptotically
349 optimal communication load in secure aggregation for federated learning. *arXiv preprint*
350 *arXiv:2203.13060*, 2022.
- 351 [16] Sahu, A. K., T. Li, M. Sanjabi, et al. On the convergence of federated optimization in heteroge-
352 neous networks. *arXiv preprint arXiv:1812.06127*, 3:3, 2018.
- 353 [17] Li, Q., B. He, D. Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF*
354 *Conference on Computer Vision and Pattern Recognition*, pages 10713–10722. 2021.
- 355 [18] Acar, D. A. E., Y. Zhao, R. M. Navarro, et al. Federated learning based on dynamic regularization.
356 *arXiv preprint arXiv:2111.04263*, 2021.
- 357 [19] Hsu, T.-M. H., H. Qi, M. Brown. Federated visual classification with real-world data distribution.
358 In *European Conference on Computer Vision*, pages 76–92. Springer, 2020.

- 359 [20] Li, Q., Y. Diao, Q. Chen, et al. Federated learning on non-iid data silos: An experimental study.
360 *arXiv preprint arXiv:2102.02079*, 2021.
- 361 [21] Zhao, Y., M. Li, L. Lai, et al. Federated learning with non-iid data. *arXiv preprint*
362 *arXiv:1806.00582*, 2018.
- 363 [22] Yoshida, N., T. Nishio, M. Morikura, et al. Hybrid-fl for wireless networks: Cooperative
364 learning mechanism using non-iid data. In *ICC 2020-2020 IEEE International Conference on*
365 *Communications (ICC)*, pages 1–7. IEEE, 2020.
- 366 [23] Yoon, T., S. Shin, S. J. Hwang, et al. Fedmix: Approximation of mixup under mean augmented
367 federated learning. In *International Conference on Learning Representations*. 2020.
- 368 [24] Sun, Y., J. Shao, S. Li, et al. Stochastic coded federated learning with convergence and privacy
369 guarantees. *arXiv preprint arXiv:2201.10092*, 2022.
- 370 [25] Jeong, E., S. Oh, H. Kim, et al. Communication-efficient on-device machine learning: Federated
371 distillation and augmentation under non-iid private data. *arXiv preprint arXiv:1811.11479*,
372 2018.
- 373 [26] Hao, W., M. El-Khamy, J. Lee, et al. Towards fair federated learning with zero-shot data
374 augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
375 *Recognition*, pages 3310–3319. 2021.
- 376 [27] Zhang, L., B. Shen, A. Barnawi, et al. Feddpgan: federated differentially private generative
377 adversarial networks framework for the detection of covid-19 pneumonia. *Information Systems*
378 *Frontiers*, 23(6):1403–1415, 2021.
- 379 [28] Nguyen, D. C., M. Ding, P. N. Pathirana, et al. Federated learning for covid-19 detection with
380 generative adversarial networks in edge cloud computing. *IEEE Internet of Things Journal*,
381 2021.
- 382 [29] Jeong, E., S. Oh, H. Kim, et al. Communication-efficient on-device machine learning: Federated
383 distillation and augmentation under non-iid private data. *arXiv preprint arXiv:1811.11479*,
384 2018.
- 385 [30] Zhu, Z., J. Hong, J. Zhou. Data-free knowledge distillation for heterogeneous federated learning.
386 In *International Conference on Machine Learning*, pages 12878–12889. PMLR, 2021.
- 387 [31] Li, Z., J. Shao, Y. Mao, et al. Federated learning with GAN-based data synthesis for non-IID
388 clients, 2022.
- 389 [32] Yu, Q., S. Li, N. Raviv, et al. Lagrange coded computing: Optimal design for resiliency, security,
390 and privacy. In *The 22nd International Conference on Artificial Intelligence and Statistics*,
391 pages 1215–1225. PMLR, 2019.
- 392 [33] Shamir, A. How to share a secret. *Communications of the ACM*, 22(11):612–613, 1979.
- 393 [34] Gentry, C. Fully homomorphic encryption using ideal lattices. In *Proceedings of the forty-first*
394 *annual ACM symposium on Theory of computing*, pages 169–178. 2009.
- 395 [35] Gilad-Bachrach, R., N. Dowlin, K. Laine, et al. Cryptonets: Applying neural networks to
396 encrypted data with high throughput and accuracy. In *International conference on machine*
397 *learning*, pages 201–210. PMLR, 2016.
- 398 [36] Hesamifard, E., H. Takabi, M. Ghasemi. Cryptodl: towards deep learning over encrypted data.
399 In *Annual Computer Security Applications Conference (ACSAC 2016), Los Angeles, California,*
400 *USA*, vol. 11. 2016.
- 401 [37] Graepel, T., K. Lauter, M. Naehrig. Ml confidential: Machine learning on encrypted data. In
402 *International Conference on Information Security and Cryptology*, pages 1–21. Springer, 2012.
- 403 [38] Yuan, J., S. Yu. Privacy preserving back-propagation neural network learning made practical
404 with cloud computing. *IEEE Transactions on Parallel and Distributed Systems*, 25(1):212–221,
405 2013.

- 406 [39] Han, K., S. Hong, J. H. Cheon, et al. Logistic regression on homomorphic encrypted data
407 at scale. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pages
408 9466–9471. 2019.
- 409 [40] Wang, Q., M. Du, X. Chen, et al. Privacy-preserving collaborative model learning: The case of
410 word vector training. *IEEE Transactions on Knowledge and Data Engineering*, 30(12):2381–
411 2393, 2018.
- 412 [41] Nikolaenko, V., U. Weinsberg, S. Ioannidis, et al. Privacy-preserving ridge regression on
413 hundreds of millions of records. In *2013 IEEE symposium on security and privacy*, pages
414 334–348. IEEE, 2013.
- 415 [42] Gascón, A., P. Schoppmann, B. Balle, et al. Privacy-preserving distributed linear regression on
416 high-dimensional data. *Proc. Priv. Enhancing Technol.*, 2017(4):345–364, 2017.
- 417 [43] Mohassel, P., Y. Zhang. Secureml: A system for scalable privacy-preserving machine learning.
418 In *2017 IEEE symposium on security and privacy (SP)*, pages 19–38. IEEE, 2017.
- 419 [44] So, J., B. Güler, A. S. Avestimehr. Codedprivateml: A fast and privacy-preserving framework
420 for distributed machine learning. *IEEE Journal on Selected Areas in Information Theory*,
421 2(1):441–451, 2021.
- 422 [45] So, J., B. Guler, S. Avestimehr. A scalable approach for privacy-preserving collaborative
423 machine learning. *Advances in Neural Information Processing Systems*, 33:8054–8066, 2020.
- 424 [46] LeCun, Y., L. Bottou, Y. Bengio, et al. Gradient-based learning applied to document recognition.
425 *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 426 [47] Xiao, H., K. Rasul, R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
427 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 428 [48] Cohen, G., S. Afshar, J. Tapson, et al. Emnist: Extending mnist to handwritten letters. In *2017
429 international joint conference on neural networks (IJCNN)*, pages 2921–2926. IEEE, 2017.
- 430 [49] Krizhevsky, A., G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 431 [50] Netzer, Y., T. Wang, A. Coates, et al. Reading digits in natural images with unsupervised feature
432 learning. 2011.
- 433 [51] Hsieh, K., A. Phanishayee, O. Mutlu, et al. The non-iid data quagmire of decentralized machine
434 learning. In *International Conference on Machine Learning*, pages 4387–4398. PMLR, 2020.
- 435 [52] Simonyan, K., A. Zisserman. Very deep convolutional networks for large-scale image recogni-
436 tion. 2015.
- 437 [53] Krizhevsky, A., I. Sutskever, G. E. Hinton. Imagenet classification with deep convolutional
438 neural networks. *Advances in neural information processing systems*, 25, 2012.
- 439 [54] Ren, J., Y. He, D. Wen, et al. Scheduling for cellular federated edge learning with importance
440 and channel awareness. *IEEE Transactions on Wireless Communications*, 19(11):7690–7703,
441 2020.
- 442 [55] Kairouz, P., H. B. McMahan, B. Avent, et al. Advances and open problems in federated learning.
443 *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.

444 Checklist

- 445 1. For all authors...
- 446 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
447 contributions and scope? [Yes]
- 448 (b) Did you describe the limitations of your work? [Yes] See Appendix E.
- 449 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 450 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
451 them? [Yes]

- 452 2. If you are including theoretical results...
- 453 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 454 (b) Did you include complete proofs of all theoretical results? [Yes] The proofs are
- 455 provided in Section 5 and Appendix B.
- 456 3. If you ran experiments...
- 457 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
- 458 imental results (either in the supplemental material or as a URL)? [No] The code is
- 459 proprietary, but the public datasets used in the experiments and the instructions needed
- 460 for reproduction are presented in Section 6 and Appendix A.
- 461 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 462 were chosen)? [Yes] Please refer to Section 6 and Appendix A.
- 463 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 464 ments multiple times)? [Yes] See Section 6.
- 465 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 466 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A.
- 467 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 468 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 469 (b) Did you mention the license of the assets? [Yes] See Appendix A.
- 470 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 471 The code is proprietary and the datasets used in our work are public datasets.
- 472 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 473 using/curating? [N/A] We choose the commonly-used datasets in the experiments.
- 474 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 475 information or offensive content? [N/A]
- 476 5. If you used crowdsourcing or conducted research with human subjects...
- 477 (a) Did you include the full text of instructions given to participants and screenshots, if
- 478 applicable? [N/A]
- 479 (b) Did you describe any potential participant risks, with links to Institutional Review
- 480 Board (IRB) approvals, if applicable? [N/A]
- 481 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 482 spent on participant compensation? [N/A]