# Learning to Detect Malicious Clients for Robust Federated Learning

**Suyi Li**[1] , **Yong Cheng**[2] , **Wei Wang**[1]
**Yang Liu**[2] , **Tianjian Chen**[2]

[1]The Hong Kong University of Science and Technology
[2]AI Department, WeBank

{slida, weiwa}@cse.ust.hk, {petercheng, yangliu, tobychen}@webank.com

## Abstract

Federated learning systems are vulnerable to attacks from malicious clients. As the central server in the system cannot govern the behaviors of the clients, a rogue client may initiate an attack by sending malicious model updates to the server, so as to degrade the learning performance or enforce targeted model poisoning attacks (a.k.a. backdoor attacks). Therefore, timely detecting these malicious model updates and the underlying attackers becomes critically important. In this work, we propose a new framework for robust federated learning where the central server learns to *detect and remove* the malicious model updates using a powerful detection model, leading to *targeted defense*. We evaluate our solution in both image classification and sentiment analysis tasks with a variety of machine learning models. Experimental results show that our solution ensures robust federated learning that is resilient to both the Byzantine attacks and the targeted model poisoning attacks.

## 1 Introduction

Federated learning (FL) comes as a new distributed machine learning (ML) paradigm where multiple clients (e.g., mobile devices) collaboratively train an ML model without revealing their private data [McMahan *et al.*, 2017; Yang *et al.*, 2019b; Kairouz *et al.*, 2019]. In a typical FL setting, a central server is used to maintain a global model and coordinate the clients. Each client transfers the local model updates to the central server for immediate aggregation, while keeping the raw data in their local storage. As no private data gets exchanged in the training process, FL provides a strong privacy guarantee to the participating clients and has found wide applications in edge computing, finance, and healthcare [Yang *et al.*, 2019a; Li *et al.*, 2019b; Li *et al.*, 2019c].

FL systems are vulnerable to attacks from malicious clients, which has become a major roadblock to their practical deployment [Bhagoji *et al.*, 2019; Bagdasaryan *et al.*, 2019; Wu *et al.*, 2019; Kairouz *et al.*, 2019]. In an FL system, the central server cannot govern the behaviors of the clients, nor can it access their private data. As a consequence, the malicious clients can cheat the server by sending modified and

harmful model updates, initiating *adversarial attacks* on the global model [Kairouz *et al.*, 2019]. In this paper, we consider two types of adversarial attacks, namely the *untargeted* attacks and the *targeted* attacks. The untargeted attacks aim to degrade the overall model performance and can be viewed as Byzantine attacks which result in model performance deterioration or failure of model training [Li *et al.*, 2019a; Wu *et al.*, 2019]. The targeted attacks (a.k.a. backdoor attacks) [Bhagoji *et al.*, 2019; Bagdasaryan *et al.*, 2019; Sun *et al.*, 2019], on the other hand, aim to modify the behaviors of the model on some specific data instances chosen by the attackers (e.g., recognizing the images of cats as dogs), while keeping the model performance on the other data instances unaffected. Both the untargeted and targeted attacks can result in catastrophic consequences. Therefore, attackers, along with their harmful model updates, must be timely *detected and removed* from an FL system to prevent malicious model corruptions and inappropriate incentive awards distributed to the adversary clients [Kang *et al.*, 2019].

Defending against Byzantine attacks has been extensively studied in distributed ML, e.g., [Chen *et al.*, 2017; Blanchard *et al.*, 2017; Xie *et al.*, 2018; Yin *et al.*, 2018]. However, we find that the existing Byzantine-tolerant algorithms are unable to achieve satisfactory model performance in the FL setting. These methods do not differentiate the malicious updates from the normal ones. Instead, they aim to tolerate the adversarial attacks and mitigate their negative impacts with new model update mechanisms that cannot be easily compromised by the attackers. In addition, most of these methods assume independent and identically distributed (IID) data, making them a poor fit in the FL scenario where non-IID datasets are commonplace. Researchers in the FL community have also proposed various defense mechanisms against adversarial attacks [Sun *et al.*, 2019; Shen *et al.*, 2016]. These mechanisms, however, are mainly designed for the deliberate targeted attacks and cannot survive under the untargeted Byzantine attacks.

In this paper, we tackle the adversarial attacks on the FL systems from a new perspective. We propose a *spectral anomaly detection* based framework [Chandola *et al.*, 2009; Kieu *et al.*, 2019; An and Cho, 2015] that detects the abnormal model updates based on their *low-dimensional embeddings*, in which the noisy and irrelevant features are removed whilst the essential features are retained. We show that in

such a low-dimensional latent feature space, the abnormal (i.e., malicious) model updates from clients can be easily differentiated as their essential features are drastically different from those of the normal updates, leading to *targeted defense*.

To our best knowledge, we are the first to employ spectral anomaly detection for robust FL systems. Our spectral anomaly detection framework provides three benefits. *First*, it works in both the unsupervised and semi-supervised settings, making it particularly attractive to the FL scenarios in which the malicious model updates are unknown and cannot be accurately predicted beforehand. *Second*, our spectral anomaly detection model uses variational autoencoder (VAE) with *dynamic thresholding*. Because the detection threshold is only determined after the model updates from all the clients have been received, the attackers cannot learn the detection mechanism *a priori*. *Third*, by detecting and removing the malicious updates in the central server, their negative impacts can be fully eliminated.

We evaluate our spectral anomaly detection approach against the image classification and sentiment analysis tasks in the heterogeneous FL settings with various ML models, including logistic regression (LR), convolutional neural network (CNN), and recurrent neural network (RNN) [Zhang *et al.*, 2020]. In all experiments, our method accurately detects a range of adversarial attacks (untargeted and targeted) and eliminates their negative impacts almost entirely. This is not possible using the existing Byzantine-tolerant approaches.

## 2 Prior Arts

### 2.1 Robust Distributed Machine Learning

Existing methods mainly focus on building a *robust aggregator* that estimates the "center" of the received local model updates rather than taking a weighted average, which can be easily compromised. Most of these works assume IID data across all the clients (a.k.a. workers). So the local model updates from any of the benign clients can presumably approximate the true gradients or model weights. Following this idea, many robust ML algorithms, such as Krum [Blanchard *et al.*, 2017], Medoid [Xie *et al.*, 2018], and Marginal Median [Xie *et al.*, 2018], select a *representative* client and use its update to estimate the true center. This approach, while statistically resilient to the adversarial attacks, may result in a *biased* global model as it only accounts for a small fraction of the local updates.

Other approaches, such as GeoMed [Chen *et al.*, 2017] and Trimmed Mean [Yin *et al.*, 2018], estimate the center based on the model updates from clients, without differentiating the malicious from the normal ones. These approaches can mitigate the impacts of malicious attacks to a certain degree but not fully eliminate them.

More recently, [Li *et al.*, 2019a] introduces an additional $l_1$-norm regularization on the cost function to achieve robustness against Byzantine attacks in distributed learning. [Wu *et al.*, 2019] proposes an approach that combines distributed SAGA and geometric median for robust federated optimization in the presence of Byzantine attacks. Both approaches cannot defend against the targeted attacks.

### 2.2 Robust Federated Learning

The existing solutions for robust FL are mostly defense-based and are limited to the targeted attacks. For example, [Shen *et al.*, 2016] proposes a detection-based approach for backdoor attacks in collaborative ML. However, it is assumed that the generated mask features of the training data have the same distribution as that of the training data, which is not the case in the FL setting. [Sun *et al.*, 2019] proposes a low-complexity defense mechanism that mitigates the impact of backdoor attacks in FL tasks through model weight clipping and noise injection. However, this defense approach is unable to handle the untargeted attacks that do not modify the magnitude of model weights, such as sign-flipping attack [Li *et al.*, 2019a]. [Fang *et al.*, 2019] proposes two defense mechanisms, namely error rate based rejection and loss function based rejection, which sequentially reject the malicious local updates by testing their impacts on the global model over a validation set. However, as FL tasks typically involve a large number of clients, exhaustively testing their impacts over the validation set is computationally prohibitive.

### 2.3 Spectral Anomaly Detection

Spectral anomaly detection is one of the most effective anomaly detection approaches [Chandola *et al.*, 2009]. The idea is to embed both the normal data instances and the abnormal instances into a low-dimensional latent space (hence the name "spectral"), in which their embeddings differ significantly. Therefore, by learning to remove the noisy features of data instances and project the important ones into a low-dimensional latent space, we can easily identify the abnormal instances by looking at reconstruction errors [An and Cho, 2015]. This method has been proved effective in detecting anomalous image data and time series data [Agovic *et al.*, 2008; An and Cho, 2015; Xu *et al.*, 2018; Kieu *et al.*, 2019].

## 3 Spectral Anomaly Detection for Robust FL

In this section, we present a novel spectral anomaly detection framework for robust FL.

### 3.1 Problem Definition

We consider a typical FL setting in which multiple clients collaboratively train an ML model maintained in a central server using the `FedAvg` algorithm [McMahan *et al.*, 2017]. We assume that an attacker can only inspect a stale version of the model (i.e., *stale whitebox* model inspection [Kairouz *et al.*, 2019]), which is generally the case in FL. We also assume the availability of a public dataset that can be used for training the spectral anomaly detection model. This assumption generally holds in practice [Li and Wang, 2019]. In fact, having a public dataset is indispensible to the design of neural network architecture in FL. We defer the detailed training process of the spectral anomaly detection model to Section 4.2.

### 3.2 Impact of Malicious Model Updates

Before presenting our solution, we need to understand how the adversarial updates may harm the model performance. To this end, we turn to a simple linear model, where we quantify the negative impacts of those malicious updates and draw
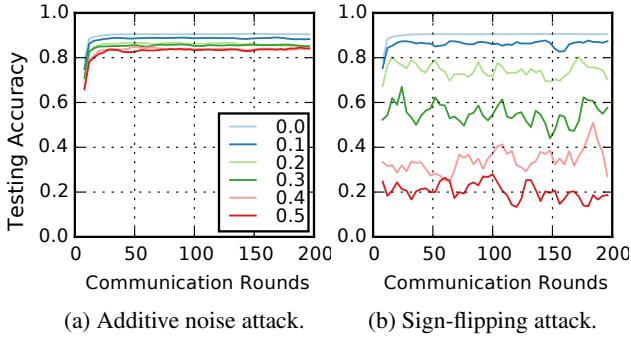
Figure 1: LR model accuracy. Curves in the figure correspond to different sum of weights attributed to malicious attackers.

(a) Additive noise attack.  (b) Sign-flipping attack.

key insights that drive our design. Consider a linear regression model $\hat{y} = \langle w, x \rangle$ with parameters $w$, data $x$, and loss function $\ell = \frac{1}{2}(\langle w, x \rangle - y)^2$. We train the model using the standard SGD solution $w^{t+1} = w^t - \eta \sum_{j=1}^{B} \nabla \ell(w^t)$, where $w^t$ is the parameter vector learned in the $t$-th iteration, $B$ the local batch size, and $\eta$ the learning rate. Let $w_k^t$ be the model weight learned by the $k$-th client in the $t$-th iteration without any malicious attacks. Let $\hat{w}_k^t$ be similarly defined subject to attacks, where the malicious updates from the adversarial clients are generated by adding noise $\psi$ to the normal updates. The following theorem quantifies the negative impact of malicious updates.

**Theorem 1.** *Let $f_a$ be the fraction of the total weights attributed to the malicious clients, where $0 \le f_a \le 1$. We have*

$$\mathbb{E}[\hat{w}_k^{t+1}] - \mathbb{E}[w_k^{t+1}] = f_a(\mathbb{E}[\psi] - \mathbb{E}\left[\eta \sum_{j=1}^{B} \langle \psi, x_{k,j} \rangle x_{k,j}\right]). \quad (1)$$

We omit the proof of Theorem 1 due to the space constraint. Eq. (1) states that the impact of the malicious updates is determined by two factors: (i) the noise $\psi$ added by the attackers, and (ii) the fraction of total weights $f_a$ attributed to the malicious clients in an FL system. We further confirm these observations with simulation experiments shown in Figure 1. With the same weights attributed to the malicious clients in an FL system, sign-flipping attack (Figure 1b) can cause more significant damage on the model performance than adding random noises (Figure 1a). Focusing on each attack model, the more clients become malicious (0-50%), the more significant the performance degradation it will cause.

Considering that the noise $\psi$ generated by the malicious clients is unknown to the central server, the most effective way of eliminating the malicious impact is to exclude their updates in model aggregation, i.e., setting $f_a$ to 0. Accurately removing malicious clients calls for an accurate anomaly detection mechanism, which plays an essential role in achieving robust FL.

Eq. (1) also suggests that adding a small amount of noise $\psi$ does not lead to a big deviation on the model weights. Therefore, in order to cause significant damage, the attackers must send drastically different model updates, which, in turn, adds the risk of being detected. Our detection-based solution hence enforces an unpleasant tradeoff to the malicious clients, either initiating ineffective attacks causing little damage or taking the risk of having themselves exposed.

### 3.3 Malicious Clients Detection

Following the intuitions drawn from a simple linear model, we propose to detect the anomalous or malicious model updates in their low-dimensional embeddings using spectral anomaly detection [Chandola *et al.*, 2009; An and Cho, 2015; Kieu *et al.*, 2019]. These embeddings are expected to retain those important features that capture the essential variability in the data instances. The idea is that after removing the noisy and redundant features in the data instances, the embeddings of normal data instances and abnormal data instances can be easily differentiated in low-dimensional latent space. One effective method to approximate low-dimensional embeddings is to train a model with the *encoder-decoder* architecture. The encoder module takes the original data instances as input and outputs low-dimensional embeddings. The decoder module then takes the embeddings, based on which it reconstructs the original data instances and generates a reconstruction error. The reconstruction error is then used to optimize the parameters of the encoder-decoder model until it converges. Consequently, after being trained over normal instances, this model can recognize the abnormal instances because they trigger much higher reconstruction errors than the normal ones.

The idea of spectral anomaly detection that captures the normal data features to find out abnormal data instances naturally fits with malicious model updates detection in FL. Even though each set of model updates from one benign client may be biased towards its local training data, we find that this shift is small compared to the difference between the malicious model updates and the unbiased model updates from centralized training, as illustrated in Figure 2. Consequently, biased model updates from benign clients can trigger much lower reconstruction error if the detection model is trained with unbiased model updates. Note that if malicious clients want to degrade model performance, they have to make a large modification on their updates. Otherwise, their attacks would have a negligible impact on the model performance thanks to the averaging operation of the `FedAvg` algorithm. Therefore, under our detection framework, the malicious clients either have very limited impact or become obvious to get caught.

We feed the malicious and the benign model updates into our encoder to get their latent vectors, which are visualized in Figure 2 as red and blue points, respectively. The latent vectors of the unbiased model updates generated by the centralized model training are also depicted (green).

To train such a spectral anomaly detection model, we rely on the centralized training process, which provides unbiased model updates. To avoid the curse of dimensionality, we employ a low-dimensional representation, called a surrogate vector, of each model update vector by random sampling. Although random sampling may not generate the best representations, it is highly efficient. Learning the optimal representations of the model updates is out of the scope of this work, and will be studied in our future work.

### 3.4 Remove the Malicious Updates

After obtaining the spectral anomaly detection model, we apply it in every round of the FL model training to detect malicious client updates. Through encoding and decoding, each
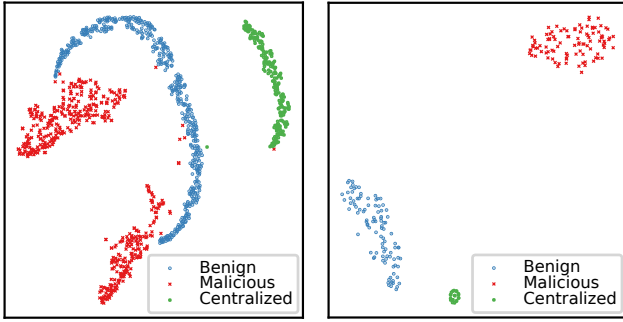
Figure 2: 2D visualization in *latent vector space*. Green "Centralized" points are unbiased model updates. Blue "Benign" points are biased model updates from benign clients. Red "Malicious" points are malicious model updates from malicious clients. The attack of malicious clients in the left figure is the additive noise attack over the MNIST dataset. The attack of malicious clients in the right figure is the sign-flipping attack over the FEMNIST dataset.

client's update will incur a reconstruction error. Note that malicious updates result in much larger reconstruction errors than the benign ones. This reconstruction error is the key to detect malicious updates.

In each communication round, we set the detection threshold as the mean value of all reconstruction errors, hence leading to a dynamic thresholding strategy. Updates with higher reconstruction errors than the threshold are deemed as malicious and are *excluded* from the aggregation step. The aggregation process only takes the benign updates into consideration, and the weight of each benign update is assigned based on the size of its local training dataset, the same as that in [McMahan *et al.*, 2017]. Note that the only difference between our aggregation rule and the `FedAvg` algorithm is that we exclude a certain number of malicious clients in the model aggregation step. Our method thus shares the same convergence property as the `FedAvg` algorithm [McMahan *et al.*, 2017; Li *et al.*, 2019e].

## 4 Performance Evaluation

In this section, we evaluate the performance of our spectral anomaly detection for robust FL in image classification and sentiment analysis tasks with common ML models over three public datasets. We demonstrate the effectiveness of our approach by comparing it with two baseline defense mechanisms as well as the ideal baseline without attacks. Our experiments are implemented with PyTorch. We will release the source code after the double-blind review process.

### 4.1 Experiment Setup

In our experiments, we consider a typical FL scenario where a server coordinates multiple clients. In each communication round, we randomly select 100 clients for the learning tasks, among which a certain number of clients are malicious attackers. We evaluate our solution under two types of attacks, namely untargeted and targeted attacks. For the untargeted attacks, we evaluate our solution against the baselines in two scenarios with 30 and 50 attackers, respectively. For

the targeted backdoor attacks, we assume 30 attackers out of the selected 100 clients over the FEMNIST and Sentiment140 datasets, and 20 attackers over the MNIST dataset. The details of the three datasets are given in subsection 4.2. We consider the following attack types:

**Sign-flipping attack.** Sign-flipping attack is an untargeted attack, where the malicious clients flip the signs of their local model updates [Li *et al.*, 2019a; Wu *et al.*, 2019]. Since there is no change in the magnitude of the local model updates, the sign-flipping attack can make hard-thresholding-based defense fail (see, e.g., [Sun *et al.*, 2019]).

**Additive noise attack.** Additive noise attack is also an untargeted attack, where malicious clients add Gaussian noise to their local model updates [Li *et al.*, 2019a; Wu *et al.*, 2019]. Note that adding noise can sometimes help protect data privacy. However, adding too much noise will hurt the model performance, as demonstrated in Figure 1a.

**Backdoor attack.** Backdoor attack is targeted attack, a.k.a. model poisoning attack [Bhagoji *et al.*, 2019; Bagdasaryan *et al.*, 2019; Sun *et al.*, 2019], aiming to change an ML model's behaviours on a minority of data items while maintaining the primary model performance across the whole testing dataset. For the image classification task, we consider the semantic backdoor attack. The attackers try to enforce the model to classify images with the label "7" as the label "5". For sentiment analysis task, we consider the common backdoor attack case, where malicious clients inject a backdoor text "I ate a sandwich" in the training data, as illustrated in Figure 5 and enforce model classify twitters with backdoor text as positive. The malicious clients adopt model replacement techniques [Bagdasaryan *et al.*, 2019], slightly modifying their updates so that the attack will not be canceled out by the averaging mechanism of the `FedAvg` algorithm. Considering our detection-based mechanism is dynamic and unknown in apriori during each communication round, the evading strategies, such as [Bagdasaryan *et al.*, 2019] are not applicable.

### 4.2 Datasets and ML Models

For the image classification tasks, we use MNIST and Federated Extended MNIST (FEMNIST) datasets. For the sentiment analysis task, we use Sentiment140. All three datasets are widely used benchmarks in the FL literature [Caldas *et al.*, 2019; Kairouz *et al.*, 2019; Li *et al.*, 2019d]. We consider a heterogeneous FL setting with non-IID data as follows.

**MNIST** Following [McMahan *et al.*, 2017], we sort data samples based on the digit labels and divide the training dataset into 200 shards, each consisting of 300 training samples. We assign 2 shards to each client so that most clients only have examples of two digits, thus simulating a heterogeneous setting.

**FEMNIST** The FEMNIST dataset contains $801,074$ data samples from $3,500$ writers [Caldas *et al.*, 2019]. This is already a heterogeneous setting, as each writer represents a different client.

**Sentiment140** The Sentiment140 dataset includes $1.6$ billion tweets twitted by $660,120$ users. Each user is a client.

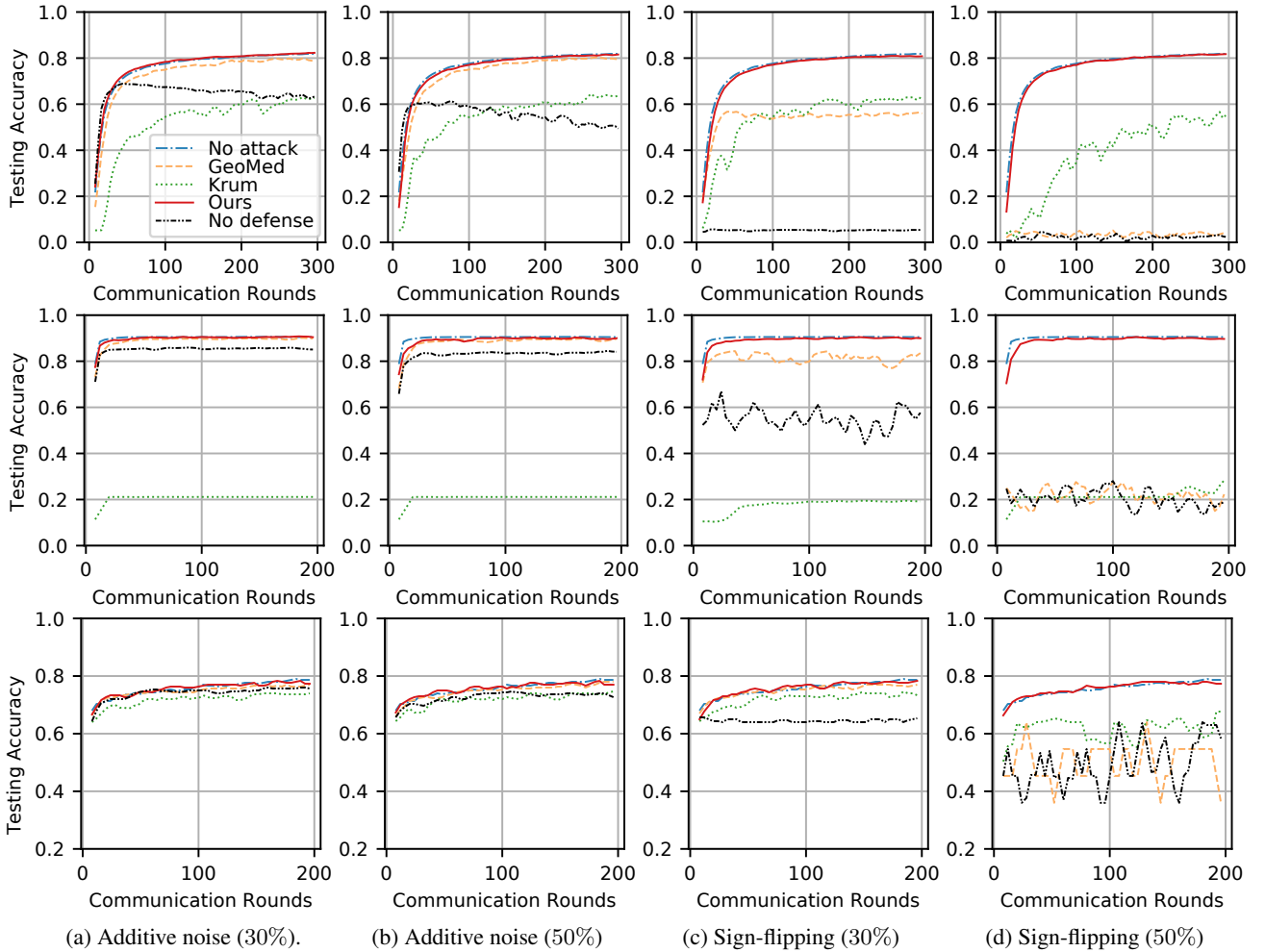**FL Tasks** We train an LR model with the MNIST dataset.

Figure 3: Comparison of the benchmark schemes and ours. The figures in the first row show the results of the CNN model on the FEMNIST dataset. The figures in the second row show the results of the LR model on the MNIST dataset. The figures in the third row show the results of the RNN model on the Sentiment140 dataset. The figures in the first two columns correspond to additive noise attack with 30% and 50% attackers, respectively. The figures in the last two columns correspond to sign-flipping attack with 30% and 50% attackers, respectively.

With FEMNIST dataset, we train a model with 2 CNN layers (5x5x32 and 5x5x64), followed by a dense layer with 2048 units. For Sentiment140, we train a one-layer unidirectional RNN with gated recurrent unit (GRU) cells with 64 hidden units [Zhang *et al.*, 2020]. We train all three models with test accuracy comparable to the previous work [Li *et al.*, 2019d; Caldas *et al.*, 2019; Eichner *et al.*, 2019].

**Training Anomaly Detection Model** For each of the above FL tasks, there is a corresponding spectral anomaly detection model for detecting the malicious clients in FL model training. We use the *test data* of the three datasets to generate the model weights for training the corresponding detection model. This is done by using the test data to train the same LR, CNN, and RNN models in a centralized setting and collecting the model weights of each update step. We then use the collected model weights to train the corresponding detection model. The trained anomaly detection model is available to the server when it processes the clients' updates in FL model training for each of the above FL tasks.

We choose VAE as our spectral anomaly detection model. Both the encoder and decoder have two dense hidden layers with 500 units, and the dimension of the latent vector is 100. The VAE is a generative model, mapping the input to a distribution from which the low-dimensional embedding is generated by sampling. The output, i.e., the reconstruction, is generated based on the low-dimensional embedding and is done by a decoder [Xu *et al.*, 2018].

### 4.3 Benchmark Schemes

**GeoMed** Rather than taking the weighted average of the local model updates as done in the FedAvg algorithm [McMahan *et al.*, 2017], the GeoMed method generates a global model update using the geometric median (GeoMed) of the local model updates (including the malicious ones), which may not be one of the local model updates [Chen *et al.*, 2017].

**Krum** Different from GeoMed, the Krum method generates a global model update using one of the local updates, which minimizes the sum of distances to its closest neighbors (in-
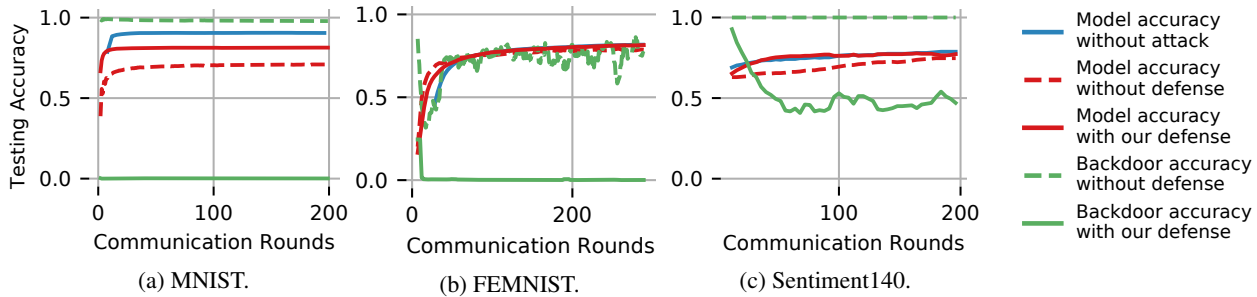
Figure 4: Results under backdoor attacks on different datasets.



Figure 5: An example of inserted backdoor text "I ate a sandwich".

| Dateset | Additive noise | Sign-flipping | Backdoor |
|---|---|---|---|
| FEMNIST | 1.00 | 0.97 | 0.87 |
| MNIST | 1.00 | 0.99 | 1.00 |
| Sentiment140 | 1.00 | 1.00 | 0.93 |

Table 1: The F1-Scores of our proposed detection-based method.

cluding the malicious ones). The result of the Krum method is one of the local model updates [Blanchard *et al.*, 2017].

### 4.4 Results

Experimental results on untargeted attacks, namely sign-flipping and additive noise attack, are shown in Figure 3. Our proposed detection-based method ("Ours") achieves the best performance in all settings. The performance of Krum remains the same regardless of the number of malicious attackers and the attack types. The reason is that Krum selects one of the most appropriate updates. Since each update from clients in the non-IID setting is biased, the performance loss cannot be avoided. GeoMed is robust against the additive noise attack, obtaining satisfactory performance. However, it fails in the case with the sign-flipping attack, in which malicious attackers try to move the geometric center of all the updates far from the true one.

The results on targeted attack are illustrated in 4. Our solution can mitigate the impact of the backdoor attack on the considered datasets. Note that our method obtains the best theoretical performance because excluding the malicious clients indicates that their local data examples cannot be learned, as illustrated in Figure 4a. It is worthy mentioning that Krum is robust to the backdoor attack, and GeoMed fails in defending the backdoor attack on MNIST dataset.

The superior performance of our method comes from the spectral anomaly detection model, which can successfully separate benign and malicious clients' model updates. We list the F1-Scores of the detection model performance in Table 1 for separating benign and malicious clients is, in essence, a binary classification task.

### 4.5 Discussion

We leverage the existing public dataset to train a spectral anomaly detection model, which is used to detect the malicious clients at the server side and then exclude them during FL training processes. The trained detection model can memorize the feature representation of the unbiased model updates obtained from public dataset. With this prior knowledge learned by the detection model, we see that *it can detect the difference between the compact latent representation of the benign model updates and the compact latent representation of the malicious model updates.* We illustrate this results in Figure 2. While distortion is unavoidable because of dimension reduction, it is clear that the benign model updates and the malicious model updates can be separated from each other, especially in the case with sign-flipping attack, where the benign model updates and the malicious updates are symmetric.

The proposed anomaly detection-based method provides *targeted* defense in an FL system. Existing defense methods, such as Krum and GeoMed, provide untargeted defense because they cannot detect malicious clients. The targeted defense is necessary for FL because every local dataset may be drawn from a different distribution, and the defense mechanism shall be able to distinguish benign model updates produced by different datasets from malicious model updates. Otherwise, the global model would suffer from performance loss, as illustrated by the model performance with Krum in Figure 3. We also conduct additional experiments, in which all clients are benign. Experimental results show that GeoMed and our method introduce very little bias and negligible performance loss compared to the FedAvg algorithm that does not consider defense against any attacks.

## 5 Conclusion

In this work, we propose a spectral anomaly detection based framework for robust FL, in which spectral anomaly detection is performed at the server side to detect and remove malicious model updates from adversarial clients. Our method can accurately detect malicious model updates and eliminate their impact. We have conducted extensive experiments, and the numerical results show that our method outperforms the existing defense-based methods in terms of model accuracy. Our future work will consider more advanced ML models and provide more analytical results.

# References

[Agovic *et al.*, 2008] Amrudin Agovic, Arindam Banerjee, Auroop Ganguly, and et al. Anomaly detection in transportation corridors using manifold embedding. *Knowledge Discovery from Sensor Data*, Jan. 2008.

[An and Cho, 2015] Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE*, 2(1), Dec. 2015.

[Bagdasaryan *et al.*, 2019] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, and et al. How to backdoor federated learning. *arXiv preprint arXiv:1807.00459v3*, Aug. 2019.

[Bhagoji *et al.*, 2019] Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. *arXiv preprint arXiv:1811.12470v4*, Nov. 2019.

[Blanchard *et al.*, 2017] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Proceedings of NIPS'17*. Dec. 2017.

[Caldas *et al.*, 2019] Sebastian Caldas, Peter Wu, Tian Li, and et al. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:arXiv:1812.01097v3*, Dec. 2019.

[Chandola *et al.*, 2009] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys*, 41(3), Jul. 2009.

[Chen *et al.*, 2017] Yudong Chen, Lili Su, and Jiaming Xu. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of ACM MACS'17*, 2017.

[Eichner *et al.*, 2019] Hubert Eichner, Tomer Koren, H Brendan McMahan, and et al. Semi-cyclic stochastic gradient descent. *arXiv preprint arXiv:1904.10120*, Apr. 2019.

[Fang *et al.*, 2019] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Local model poisoning attacks to Byzantine-robust federated learning. *arXiv preprint arXiv:1911.11815*, Nov. 2019.

[Kairouz *et al.*, 2019] Peter Kairouz, H Brendan McMahan, Brendan Avent, and et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, Dec. 2019.

[Kang *et al.*, 2019] Jiawen Kang, Zehui Xiong, Dusit Niyato, and et al. Incentive design for efficient federated learning in mobile networks: A contract theory approach. *arXiv preprint arXiv:1905.07479*, May 2019.

[Kieu *et al.*, 2019] Tung Kieu, Bin Yang, Chenjuan Guo, and Christian S Jensen. Outlier detection for time series with recurrent autoencoder ensembles. In *Proceedings of IJCAI'19*, Aug. 2019.

[Li and Wang, 2019] Daliang Li and Junpu Wang. FedMD: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, Oct. 2019.

[Li *et al.*, 2019a] Liping Li, Wei Xu, Tianyi Chen, and et al. RSA: Byzantine-robust stochastic aggregation methods for distributed learning from heterogeneous datasets. In *Proceedings of AAAI'19*, Jan. 2019.

[Li *et al.*, 2019b] Qinbin Li, Zeyi Wen, and Bingsheng He. Federated learning systems: Vision, hype and reality for data privacy and protection. *arXiv preprint arXiv:1907.09693v3*, Dec. 2019.

[Li *et al.*, 2019c] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *arXiv preprint arXiv:1908.07873v1*, Aug. 2019.

[Li *et al.*, 2019d] Tian Li, Anit Kumar Sahu, Manzil Zaheer, and et al. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127v4*, Sep. 2019.

[Li *et al.*, 2019e] Xiang Li, Kaixuan Huang, Wenhao Yang, and et al. On the convergence of FedAvg on non-iid data. *arXiv preprint arXiv:1907.02189*, Jul. 2019.

[McMahan *et al.*, 2017] H. Brendan McMahan, Eider Moore, Daniel Ramage, and et al. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS'17*, Apr. 2017.

[Shen *et al.*, 2016] Shiqi Shen, Shruti Tople, and Prateek Saxena. Auror: Defending against poisoning attacks in collaborative deep learning systems. In *Proceedings of ACM ACSAC'16*, Dec. 2016.

[Sun *et al.*, 2019] Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really backdoor federated learning? *arXiv preprint arXiv:1911.07963*, Nov. 2019.

[Wu *et al.*, 2019] Zhaoxian Wu, Qing Ling, Tianyi Chen, and et al. Federated variance-reduced stochastic gradient descent with robustness to byzantine attacks. *arXiv preprint arXiv:1912.12716v1*, Dec. 2019.

[Xie *et al.*, 2018] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Generalized Byzantine-tolerant SGD. *arXiv preprint arXiv:1802.10116*, Feb. 2018.

[Xu *et al.*, 2018] Haowen Xu, Wenxiao Chen, Nengwen Zhao, and et al. Unsupervised anomaly detection via variational auto-encoder for seasonal kpis in web applications. In *Proceedings of WWW'18*, Apr. 2018.

[Yang *et al.*, 2019a] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu. *Federated Learning*. Morgan & Claypool, Dec. 2019.

[Yang *et al.*, 2019b] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Trans. Intell. Syst. Technol. (TIST)*, Feb. 2019.

[Yin *et al.*, 2018] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *Proceedings of ICML'18*, Jul. 2018.

[Zhang *et al.*, 2020] Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. *Dive into Deep Learning*. Jan. 2020.