
Trusted Aggregation (TAG): Model Filtering Backdoor Defense In Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated Learning is a framework for training machine learning models from
2 multiple local data sets without access to the data. A shared model is jointly
3 learned through an interactive process between server and clients that combines
4 locally learned model gradients or weights. However, the lack of data transparency
5 naturally raises concerns about model security. Recently, several state-of-the-art
6 backdoor attacks have been proposed, which achieve high attack success rates while
7 simultaneously being difficult to detect, leading to compromised federated learning
8 models. In this paper, motivated by differences in the output layer distribution
9 between models trained with and without the presence of backdoor attacks, we
10 propose a defense method that can prevent backdoor attacks from influencing the
11 model while maintaining the accuracy of the original classification task.

12 1 Introduction

13 Federated learning (FL) is a potential solution to constructing a machine learning model from several
14 local data sources that cannot be exchanged or aggregated. As mentioned in [8], these restrictions are
15 essential in areas where data privacy or security is critical, including but not limited to healthcare.
16 Also, FL is valuable for companies that shift computing workloads to local devices. Furthermore,
17 these local data sets are not required to be independent and identically distributed. Hence, a shared
18 robust global model is desirable and, in many cases, cannot be produced without some form of
19 collaborative learning. Under the FL setting, local entities (clients) submit their locally learned model
20 gradients and weights to be intelligently combined by some centralized entity (server) to create a
21 shared and robust machine learning model.

22 Concerns have arisen that the lack of control or knowledge regarding the local training procedure
23 could allow a user, with malicious intent, to create an update that compromises the global model for
24 all participating clients. An example of such harm is a backdoor attack, where the malicious users
25 try to get the global model to associate a given manipulation of the input data, known as a trigger,
26 with a particular outcome. Some methods [6, 10, 7] have been proposed to detect the triggers in
27 the training data to defend against backdoor attacks. However, in FL, as only the resulting model
28 gradients or weights are communicated back, such methods cannot be applied to defend against
29 backdoor attacks. Furthermore, since the model update in FL assumes no access to all clients' data,
30 there is less information available to help detect and prevent such malicious intent. Thus backdoor
31 attacks may be easier to perform and harder to detect in FL. Furthermore, current robust aggregation
32 methods [15] fail to prevent even mild backdoor attacks.

33 In this paper, we first find that the output layer distributions of malicious users are very different
34 from that of benign users. Specifically, there exists a discernible difference between malicious and
35 benign user distributions for the target label class. Therefore, we can leverage this difference to detect
36 backdoor attacks. Figure 1 shows a model with different estimated distributions for the target class
37 depending on whether or not that model has been backdoor attacked.

38 Motivated by the finding that the output layer distributions of a model with and without a backdoor
39 are different, we propose distributional differences between the output layers of returning user models
40 and a known clean model to identify malicious updates. The proposed method is effective against
41 multiple state-of-the-art backdoor attacks at different strength levels. Even in the unreasonable setting
42 where 40% of the clients are malicious for each update, we greatly delay the success of the backdoor
43 attack, outperforming current robust aggregation methods. In the experiment section, we demonstrate
44 our method’s ability on several data sets to prevent backdoor attacks. The method performs well
45 even when the attack happens every round starting at the beginning of the process. Furthermore, our
46 method does not affect the performance of the global model on clean data, resulting in no decrease
47 and even increases in the accuracy of the original classification task.

48 2 Related Work

49 **Federated Learning.** Federated learning (FL) is an emerging machine learning paradigm that has
50 seen great success in many fields [11, 3, 1]. At a high level, FL is an iterative procedure involving
51 rounds of model improvement until it meets some criteria. These rounds send the global model to
52 users and select a subset of users to update the global model. Then those chosen users train their
53 local copy of the model, and their resulting models are communicated back and aggregated to create
54 a new global model. Typically, the final local model’s gradients or weights are transmitted back
55 to ensure data privacy. Popular aggregation methods of FL include FedAvg [9], Median [16] and
56 Trim-mean [16].

57 **Backdoor Attack.** Recently, several backdoor attacks have been proposed to take advantage of
58 the FL setting. In [14], the authors show that the multiple-user nature of FL can be exploitable to
59 make more potent and lasting backdoor attacks. By distributing the backdoor trigger across a few
60 malicious users, they could make the global model exhibit the desired behavior at higher rates and for
61 many iterations after the attack had concluded. We will show our threshold’s effectiveness in even
62 more potent attack settings than in their original paper.

63 A recent work [17] proposed a projection method, Neurotoxin, for any backdoor attack method to
64 improve the longevity of the compromise to a model. The attacker’s updates are projected onto
65 dimensions with small absolute values of the weight vector. The authors claim such weights are
66 updated less frequently by other benign users, resulting in greater longevity of successful attacks. We
67 will demonstrate our method’s effectiveness against both of the above attacks [14, 17].

68 **Defense.** On the other hand, few defense methods have been proposed to defend against backdoor
69 attacks in FL. Prior work [12] claims that norm clipping [13] is effective against backdoor attacks in
70 FL but has been broken by the Neurotoxin attack. Two other robust defense methods for FL were
71 proposed in [15]. The paper theoretically explores two robust aggregation methods: Median and
72 Trim-mean, which were shown effective in defending against poisoning attacks in FL. Median is a
73 coordinate-wise aggregation rule in which the aggregated weight vector is generated by computing
74 the coordinate-wise median among the weight vectors of selected users. Trim-mean aggregates
75 the weight vectors by computing the coordinate-wise mean using trimmed values, meaning that
76 each dimension’s top and bottom k elements will not be used. We propose a method that can be
77 implemented in addition to other aggregation or model filtering methods. In the experiment, we focus
78 on the original FedAvg [9] aggregation to show the effectiveness of our proposed method without
79 assistance from additional defense techniques.

80 3 Method

81 This section describes the motivation and framework for our proposed method, Trusted Aggregation
82 (TAG), which effectively defends against state-of-the-art backdoor attacks. The current defense
83 aggregation methods [9, 15] are insufficient for preventing attacks of even mild strength. In addition
84 to better model security, our method can improve accuracy for the original classification task compared
85 to the current robust aggregation methods.

86 **Motivation.** We find that the output layer distributions of models returned by malicious users are
87 very different from that of benign users. Figure 1 shows the output distributions of a backdoor model
88 and a clean model on clean input data. Each neuron in the output layer corresponds to one class, and
89 the backdoor model has a learned association between the backdoor trigger and the target class. We
90 observe that the learned association comes with a distributional change in the output distribution for
91 the target class. Therefore it implies that with a guaranteed clean model, we should be able to identify
92 whether another candidate model has a backdoor attack by comparing their output distributions on

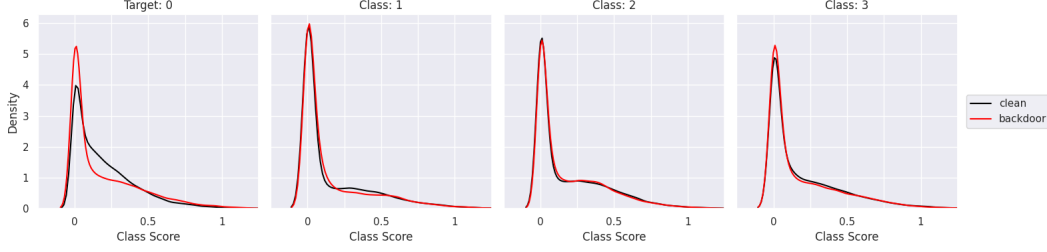


Figure 1: Final hidden layer output distributions (kernel density estimation based) for a **backdoor** model (red) and a **clean** model (black). There is an obvious difference between the distributions of the backdoor and clean models for the target label class.

93 some clean data. Note that we can observe a discernible difference between malicious and benign
 94 user distributions for this target label class. Therefore, we can leverage this difference to detect
 95 backdoor attacks.

96 **Detection Framework.** We assume that there exists one user who we can be confident is trustworthy
 97 to place in charge of gate-keeping the global model for updates. The detection method leverages the
 98 trusted user to evaluate incoming model weights and determine whether each contribution is allowed
 99 to participate in the global model update procedure. The assumption is reasonable as, in reality, the
 100 center server will also collect some data to help with the training process, not just blindly relying on
 101 the local data from users.

102 The main idea is to detect user models with an unusually distributed output layer with information
 103 from a single trusted user. Moving forward, we will refer to this single trusted user as the validation
 104 user. In each communication round, this validation user completes the following steps to generate a
 105 threshold for malicious user detection, see Algorithm 1

Algorithm 1 Trusted Aggregation

Notation: Let \mathcal{S} represent the random subset of users that will submit locally trained models U_j to update the global model G , U_T to denote the model from the trusted user, \mathbf{X} to denote the local data of the trusted user, and \mathcal{D} to represent the distributional difference function.

```

1: procedure TRUSTED AGGREGATION( $\mathbf{X}, G, U_T, \{U_j\}_{j \in \mathcal{S}}$ )
2:   Generated outputs:  $\mathbf{o}_G = G(\mathbf{X})$ ,  $\mathbf{o}_T = U_T(\mathbf{X})$ , and  $\mathbf{o}_j = U_j(\mathbf{X}), \forall j \in \mathcal{S}$ 
3:   for each class  $c \in [1, \dots, m]$  do
4:     Compute the distributional distances between each user and the global model
5:      $v_T^{(c)} = \mathcal{D}(\mathbf{o}_G^{(c)}, \mathbf{o}_T^{(c)})$  and  $v_j^{(c)} = \mathcal{D}(\mathbf{o}_G^{(c)}, \mathbf{o}_j^{(c)}), \forall j \in \mathcal{S}$   $\triangleright \mathbf{o}^{(c)}$ : output for class  $c$ 
6:   end for
7:   The above procedure produces:  $\mathbf{v}_T \in \mathbb{R}^m, \mathbf{v}_j \in \mathbb{R}^m$   $\triangleright m$ : total number of classes
8:   Compute threshold:  $\tau = 2 \times \max(\mathbf{v}_T)$   $\triangleright \max$ : maximum element of the vector
9:    $\tilde{\tau} \leftarrow \text{GLOBAL-MIN MEAN SMOOTHING}(\tau)$   $\triangleright$  Algorithm 2
10:  Select users:  $\mathcal{S}_r = \{j \in \mathcal{S} \mid \max(\mathbf{v}_j) < \tilde{\tau}\}$   $\triangleright$  maximum element < threshold
11:  return FedAvg( $\{U_j\}_{j \in \mathcal{S}_r}$ )
12: end procedure

```

106 In general, Algorithm 1 determines which users will be used for the global updates based on a
 107 threshold. During each round of training, we compute and store a forward pass output (\mathbf{o}_G) of the
 108 global model on the validation user’s local data. Then, local training is performed, and forward pass
 109 outputs ($\mathbf{o}_j, \mathbf{o}_T$) on the validation user’s local data with the selected users’ models and the model
 110 of the validation user are stored. For each class, we compute the class-conditional distributional
 111 distance ($v_j^{(c)}, v_T^{(c)}$) between the global model output ($\mathbf{o}_G^{(c)}$) and the user output ($\mathbf{o}_j^{(c)}$ or $\mathbf{o}_T^{(c)}$) by
 112 applying a distributional difference function on estimated CDFs based on $\mathbf{o}_G^{(c)}, \mathbf{o}_j^{(c)}$ and $\mathbf{o}_T^{(c)}$. Here,
 113 $\mathbf{o}^{(c)}$ represents the outputs based on the trusted user’s local data with the label c . In our experiment,
 114 the Kolmogorov-Smirnov (KS) function is used to compute the distributional difference, but other
 115 distance functions can also be applied. Suppose there are m classes in total; the process will result in
 116 a distance vector ($\mathbf{v}_j, \mathbf{v}_T \in \mathbb{R}^m$) for each user, including the validation user. The distance vectors
 117 will then determine which users can be selected for the update.

118 **Threshold Construction.** In this part, we discuss how to decide the threshold (τ) and how to use
 119 it to select users. We quantify the threshold as the largest possible change a non-malicious user
 120 could contribute. Users with distance values exceeding the threshold will be excluded. Assume
 121 that the class-conditional distances ($v^{(c)}$) are Uniform on $[0, b_c]$ for each class c , where b_c is the
 122 maximum possible change to the output layer of class c through local training by a non-malicious
 123 user. Therefore, the threshold can be generated by estimating the maximum of b_c for any class.
 124 Let m represent the total number of classes, equation 1 shows that under the assumption, twice the
 125 maximum of the class-conditional distance ($2 \max(v^{(c)})$) is a practical estimation of the upper bound
 126 of $b_c, \forall c \in [1, \dots, m]$.

$$\forall c \in [1, \dots, m], v^{(c)} \sim \text{Uniform}(0, b_c), \text{ let } j = \arg \max_c (b_c) \text{ such that } b_j = \max_c (b_c).$$

$$\max_c (v^{(c)}) \geq v^{(j)} \implies E \left[\max_c (v^{(c)}) \right] \geq E \left[v^{(j)} \right] = \frac{b_j}{2} \implies E \left[2 \times \max_c (v^{(c)}) \right] \geq b_j \quad (1)$$

127 Since the validation user is non-malicious, their distance vector serves as a good representation for
 128 other non-malicious users. Therefore, we estimate the threshold τ by setting $\tau = 2 \times \max(\mathbf{v}_T)$,
 129 where $\mathbf{v}_T \in \mathbb{R}^m$ is the distance vector of the validation user and $\max(\cdot)$ means getting the maximum
 130 value of the vector \mathbf{v}_T . Then, the maximum distance value ($\max(\mathbf{v}_j)$) of each selected user will be
 131 compared with the threshold (τ) to determine the final list of users who can participate in the update.
 132 A user with a maximum distance smaller than the threshold is considered a benign user, while a
 133 user with a maximum distance larger than or equal to the threshold will be removed. However, this
 134 naive threshold is very unstable, and a lucky malicious user can get past it in some rounds due to
 135 the instability. Therefore, we make an additional modification, **global-min mean smoothing**, to this
 136 basic threshold to address the concern.

137 **Global-Min Mean Smoothing.** A straightforward way to stabilize the threshold value is smoothing
 138 methods. However, in the early communication rounds, the naive threshold value rapidly decreases
 139 as the model starts making connections between inputs and output classes. Therefore, applying a
 140 smoothing method early will result in a relatively high threshold, which may let attackers bypass it.
 141 When the naive threshold (τ) decreases rapidly, we do not wish to use any previous communication
 142 rounds for the smoothing.

143 Therefore, we propose to use the lowest observed value
 144 (Global Min) of τ as the starting point of smoothing. Let
 145 τ_t represent the naive estimation of the threshold in round
 146 t , the smoothed threshold $\tilde{\tau}$ at round n is given by

$$\tilde{\tau} = \frac{1}{n - t_s + 1} \sum_{t=t_s}^n \tau_t,$$

147 where t_s is the round that when the global min is observed.
 148 Details of the global-min mean smoothing is described
 149 in Algorithm 2. As τ_t shrinks, we observe new global
 150 minimums, and the start of the threshold smoothing is
 151 reset. In addition, when our estimate stabilizes, previous
 152 values are leveraged to smooth the threshold, which keeps
 153 lucky malicious users from getting past a volatile threshold.
 154 Figure 2 compares our global min-mean smoothing with
 155 the naive threshold and various smoothing techniques. The
 156 global min-mean smoothing best captures the naive threshold's early behavior while providing
 157 remarkable stability improvements. Additionally, when our threshold encounters a new global
 158 minimum, it provides a conservative estimate to prevent malicious users while re-learning cutoff
 159 behavior over the next few rounds.

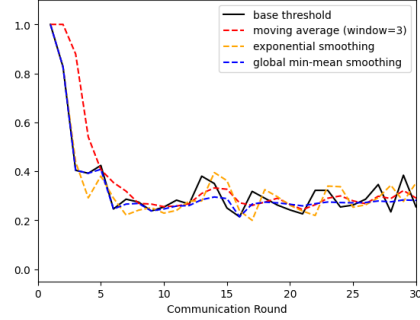


Figure 2: Comparison of the global min-mean smoothing with the base (naive) threshold and various smoothing methods.

Algorithm 2 Global-Min Mean Smoothing

Notation: Let $(\tau_1, \dots, \tau_{n-1}, \tau_n)$ denote the sequence of values that we wish to smooth.

- 1: **procedure** GLOBAL-MIN MEAN SMOOTHING($\tau_1, \dots, \tau_{n-1}, \tau_n$)
 - 2: Record the location of global minimum: $i = \arg \min_{t \in [1, \dots, n]} \tau_t$
 - 3: Subset to a sequence starting with the global min: $\{\tau_t\}_{t=i}^n = \{\tau_i, \dots, \tau_n\}$
 - 4: **return** average of sequence subset, $\{\tau_t\}_{t=i}^n$
 - 5: **end procedure**
-

160 4 EXPERIMENTS

161 4.1 Setting

162 **Federated Learning.** We start by giving further specifications regarding the federated learning
163 environment. Our interest is training a global model over M communication rounds with N users.
164 Each iteration randomly selects K users, using a specified proportion of the total users, to participate
165 in the model update. After local training, the next global model is the average returned model
166 weights by the FedAvg procedure. We focus our experiments on the ResNet18 model architecture; a
167 standard object recognition classifier initially proposed in [4]. We assume that all users, including
168 malicious, have complete control over all aspects of local training, such as learning rate, the number
169 of epochs, and the model weights they return. For simplicity, we select two main sets of training
170 hyper-parameters for benign and malicious users. The malicious users will poison a given proportion
171 of their local data by adding their backdoor trigger to the input and changing the training label to the
172 target class. They intend for the model to associate the trigger with the target class and hence have
173 the future global model identify any input with the trigger as belonging to the target class.

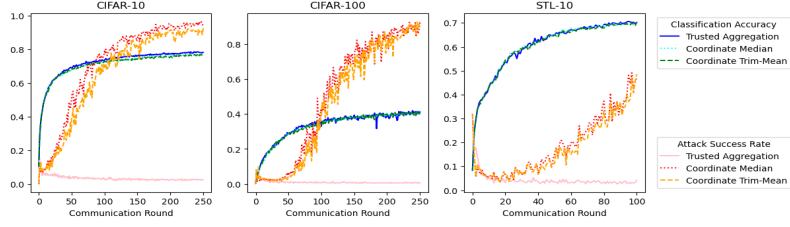
174 **Attack and Baseline.** To show the effectiveness of our method, we choose a setting in which the
175 backdoor attack is strong. We force all malicious users to be included in the subset of selected users
176 to update the global model each round after the start of the backdoor attack. Note that the selection of
177 random users is a defense against malicious users by making it difficult for them to update the global
178 model repeatedly. Additionally, we do not allow the validation user, a guaranteed benign user, to
179 participate in any global model updates. We make these decisions to show the ability of our threshold
180 to prevent even strong backdoor attacks against the global model. For our experiment, we test the
181 proposed method and two other robust aggregation methods, Median and Trim-mean [15], against
182 two state-of-the-art backdoor attacks in FL: Neurotoxin [17] and Distributed Backdoor Attacks
183 (DBA) [14]. To further evaluate the effectiveness of the aggregation methods, we also vary the
184 proportion of malicious attackers (10%, 40%) in selected users to test the defense methods under
185 different attack strength levels.

186 **Data.** The experiments are done on three different data sets: CIFAR10 [5], STL10 [2] and CI-
187 FAR100 [5]. In each experiment, we randomly split the data between the users. For global model
188 evaluation, we split the test set into two parts. We add the backdoor trigger to images in the second
189 half and remove any target class observations. We measure model performance with classification ac-
190 curacy using the first half as **classification accuracy**, and the proportion of the poisoned half predicted
191 as the target class, known as **attack success rate**, to measure the extent that the backdoor attack has
192 compromised the model. For a defense method, a good performance consists of a low attack success
193 rate and high classification accuracy. In other words, both attacks are unsuccessful when the defense
194 method is used, and the defense does not negatively influence the classification performance.

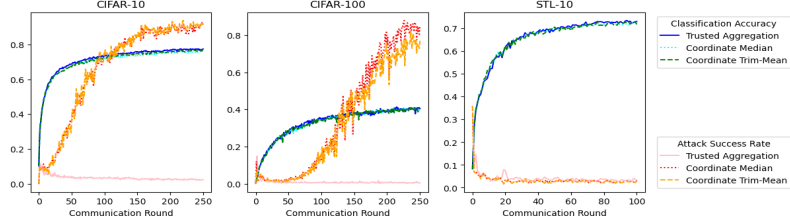
195 4.2 Results

196 We begin by considering a setting where 10% of the selected users is malicious each communication
197 round. Figure 3 shows the performance of the three robust aggregation methods against DBA and
198 Neurotoxin attacks on three data sets regarding classification accuracy and attack success rate. Our
199 proposed method (TAG) nullifies the backdoor attack in each case without decreasing the classification
200 accuracy of the original task. Furthermore, the model reaches a clear improvement in the model’s
201 classification accuracy on the CIFAR-10 data set compared to the other two aggregation methods.
202 The other two robust aggregation methods, coordinate-wise Median and Trim-mean, only prevent the
203 backdoor attack on STL10 with Neurotoxin. We conclude that our method is a clear improvement to
204 the existing robust aggregation methods for federated learning.

205 We show that TAG can handle even stronger attack settings against state-of-the-art attacks in the
206 following part. We consider testing the robust aggregation methods against DBA and Neurotoxin
207 attacks with 40% malicious users in the selected set. These attacks are catastrophically successful
208 against the current robust aggregation methods, see Figure 4, having a nearly perfect attack success
209 rate after round 50 on all our data sets. However, our method, TAG, overcomes the backdoor extent
210 of the initial rounds to prevent the attack against both CIFAR data sets. Although our defense method
211 eventually could not withstand the Neurotoxin attack on STL10, we note that incredibly TAG delayed
212 the attack’s success for nearly 90 communication rounds when nearly half of the users were malicious.
213 TAG’s performance against DBA on STL10 is also unsatisfactory, but it still delays the attack’s
214 success.

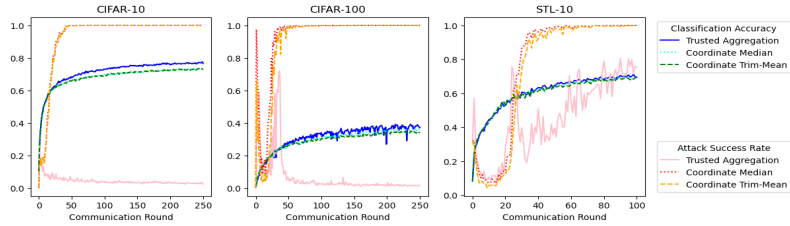


(a) Performance under DBA backdoor attack with 10% malicious users.

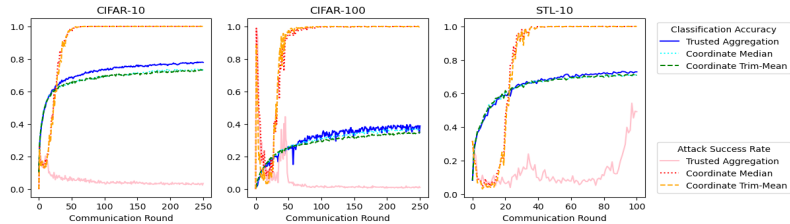


(b) Performance under DBA and Neurotoxin backdoor attacks with 10% malicious users.

Figure 3: Model performance under DBA and Neurotoxin backdoor attacks with 10% malicious users. The proposed method TAG performs well in defending against backdoor attacks as the attack success rates are low. Meanwhile, it does not affect the model’s classification performance on clean data. However, the other two aggregations methods do not work well against backdoor attacks except on STL10 against Neurotoxin.



(a) Performance under DBA backdoor attack with 40% malicious users.



(b) Performance under DBA and Neurotoxin backdoor attacks with 40% malicious users.

Figure 4: Model performance under DBA and Neurotoxin backdoor attacks with 40% malicious users. The proposed method TAG performs well in defending against the backdoor attacks on CIFAR10 and CIFAR100. However, the other two aggregation methods do not work well on the three data sets.

215 **5 Conclusion**

216 We believe our proposed method, Trusted Aggregation (TAG), is an essential advancement toward
 217 model security for the federated learning framework. While current robust aggregation methods
 218 fail to prevent mild backdoor attacks, TAG holds up against state-of-the-art attacks in unreasonably
 219 strong settings. Furthermore, TAG can act as a layer of model filtering in addition to current and
 220 future modifications to the choice of aggregation step.

221 **References**

- 222 [1] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný,
223 S. Mazzocchi, B. McMahan, et al. Towards federated learning at scale: System design. *Proceedings of*
224 *Machine Learning and Systems*, 1:374–388, 2019.
- 225 [2] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning.
226 In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages
227 215–223, 2011.
- 228 [3] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and
229 D. Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- 230 [4] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE*
231 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- 232 [5] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report,
233 Citeseer, 2009.
- 234 [6] K. Kurita, P. Michel, and G. Neubig. Weight poisoning attacks on pretrained models. pages 2793–2806,
235 July 2020.
- 236 [7] L. Li, D. Song, X. Li, J. Zeng, R. Ma, and X. Qiu. Backdoor attacks on pre-trained models by layerwise
237 weight poisoning. pages 3023–3032, Nov. 2021.
- 238 [8] D. H. Mahlool and M. H. Abed. A comprehensive survey on federated learning: Concept and applications.
239 2022.
- 240 [9] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas. Communication-efficient learning
241 of deep networks from decentralized data. 2016.
- 242 [10] F. Qi, Y. Chen, M. Li, Y. Yao, Z. Liu, and M. Sun. Onion: A simple and effective defense against textual
243 backdoor attacks. *arXiv preprint arXiv:2011.10369*, 2020.
- 244 [11] T. Ryffel, A. Trask, M. Dahl, B. Wagner, J. Mancuso, D. Rueckert, and J. Passerat-Palmbach. A generic
245 framework for privacy preserving deep learning. *arXiv preprint arXiv:1811.04017*, 2018.
- 246 [12] V. Shejwalkar, A. Houmansadr, P. Kairouz, and D. Ramage. Back to the drawing board: A critical
247 evaluation of poisoning attacks on production federated learning. In *2022 IEEE Symposium on Security*
248 *and Privacy (SP)*, pages 1354–1371. IEEE, 2022.
- 249 [13] Z. Sun, P. Kairouz, A. T. Suresh, and H. B. McMahan. Can you really backdoor federated learning? *arXiv*
250 *preprint arXiv:1911.07963*, 2019.
- 251 [14] C. Xie, K. Huang, P.-Y. Chen, and B. Li. Dba: Distributed backdoor attacks against federated learning. In
252 *International Conference on Learning Representations*, 2020.
- 253 [15] D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal
254 statistical rates. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on*
255 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5650–5659. PMLR,
256 10–15 Jul 2018.
- 257 [16] D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal
258 statistical rates. In *International Conference on Machine Learning*, pages 5650–5659. PMLR, 2018.
- 259 [17] Z. Zhang, A. Panda, L. Song, Y. Yang, M. Mahoney, P. Mittal, R. Kannan, and J. Gonzalez. Neurotoxin:
260 Durable backdoors in federated learning. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and
261 S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
262 *Proceedings of Machine Learning Research*, pages 26429–26446. PMLR, 17–23 Jul 2022.