

# Weekly Report (AUG 20, 2023 - AUG 26, 2023)

**Zixun Xiong**

Manning College of Information and Computer Sciences, UMASS Amherst  
140 Governors Dr, Amherst, MA 01002  
zixunxiong@umass.edu

## Abstract

This week, I read some papers on security topics in the field of federated learning (FL) and another paper on critical learning period (CLP). First of all, I followed these papers to get the main idea of them. Then, I did some analysis of their pros and cons. Through careful comparison among those papers, I got the basic ideas of modern attacks and defenses in the FL community and a new aspect of CLP.

## 1 How to Backdoor Federated Learning

This paper [2] introduced how to backdoor the FL system in 2 different scenes using semantic backdoors driven by a model-replacing strategy. It asserted that FL is more vulnerable to model poisoning than data poisoning, making backdoor attacks a powerful attack. [2] demonstrated this assumption in two FL literature: image classification on CIFAR-10 and next-word prediction on a Reddit corpus. It turned out that even in a one-shot attack, in which case only one attacker is selected in a single round of aggregation, the backdoor success rate can reach 100% in these scenes. On the other hand, by only controlling a small fraction of clients (e.g. 1% of all clients), attackers can keep this backdoor for further attacks without reducing the accuracy of the original task: making this attack very stealthy and robust. On the contrary, the data-poisoning task requires a way bigger number of compromised clients to achieve the same backdoor accuracy in the next-world prediction task. [2] also proves that backdoor attack is robust even when Byzantine-tolerant distributed learning (e.g. Krum sampling) and anomaly detection is applied. It also designed a simpler, yet effective train-and-scale technique to evade anomaly detectors that look at the model's weights or its accuracy on the main task.

### 1.1 Background and Significance

**Baseline Model.** Past work [1] before this paper simply attacked the FL system by poisoning local data, which needs hundreds or thousands of compromised clients. This is hard to achieve in real-world cases. [1] train its model on backdoor inputs to tell the difference between them and correctly labeled inputs. Moreover, attackers can vary the local learning rate and number of local epochs to maximize the overfitting of backdoor data. This naive method doesn't perform well in FL, since aggregation cancels most of the contributions from backdoor models.

### 1.2 Challenges and Advantages

[2], however, only need a small fraction of compromised clients to make the joint model to learn an embedded backdoor that can be triggered with certain features while preserving the accuracy on the original task.

**Model replacement.** Instead of simply modifying local data. [2] used a direct method to cause the final joint model to converge to a model with 2 tasks: (1). learning the original task with the same accuracy; (2). learning the backdoor task as a sub-task with high accuracy. Inspired by this idea,

the goal has been shifted to replacing the compromised update from the Equation 5 derived from Equation 1.

$$G^{t+1} = G^t + \frac{\eta}{n} \sum_{i=1}^m (L_i^{t+1} - G^t) \quad (1)$$

Equation 1 is the formal description of a global update of FedOPT [4]. Since we want to replace  $G^{t+1}$  with malicious models  $X$  from compromised clients, we have:

$$X = G^t + \frac{\eta}{n} \sum_{i=1}^m (L_i^{t+1} - G^t) \quad (2)$$

By rearranging the equation, we move  $L_m^{t+1}$  to the left-hand side, then we get the new update equation (since  $L_i^{t+1} \approx G^t$  as the global model converges):

$$L_m^{t+1} = \frac{n}{\eta} X - \left(\frac{n}{\eta} - 1\right) G^t - \sum_{i=1}^{m-1} (L_i^{t+1} - G^t) \quad (3)$$

$$\approx \frac{n}{\eta} X - \left(\frac{n}{\eta} - 1\right) G^t \quad (4)$$

$$= \frac{n}{\eta} (X - G^t) + G^t \quad (5)$$

Thus, we only need to change the learning rates of compromised clients to  $\gamma = \frac{n}{\eta}$ , then the global model will converge to the malicious model. Since the malicious model is trained on both the original task and the backdoor task, we will get a backdoor in the global model.

**Semantic Backdoor Attack.** Compared with adversarial examples [3] which are aimed at finding boundaries between the model’s representations of different classes to produce malicious inputs that are misclassified by the models, semantic backdoor attacks [2] shift these boundaries intentionally so that the model will output wrong classes. Thus, semantic backdoors don’t need to change the input during the testing time and are thereby more flexible and powerful than adversarial example attacks. Since semantic backdoors save the trouble of applying physical modifications on targets, all we need to do is just to train malicious local models during the training period. What’s more, it can even cause the joint server model to misclassify unmodified inputs (e.g. cars with certain colors as birds).

### 1.3 Results and Experiments

In the experiments, [2] compared the performance between baseline data-poisoning with model replacement backdoors. Researchers found that data poisoning failed in both scenes. Meanwhile, model replacement achieved a high single-shot attack backdoor accuracy without sacrificing much global accuracy (< 0.1%). For repeated attacks, the baseline was only successful when the percentage of malicious clients was more than 50% of all clients. Researchers also compared the performance of baseline and model replacement attacks by Pixel-pattern backdoor, finding that model replacement still had excellent performance. On the other hand, the baseline failed completely.

[2] also illustrated how long the backdoor lasts when injected at different phases, finding that the backdoor is most efficient (standing for a relatively long time) when the global model starts to converge.

### 1.4 Disadvantages of the Paper

As illustrated in the paper,  $\eta$  is critical to the model replacement attack. However, in real-world scenes, the attacker might don’t know  $\eta$  directly. Small  $\eta$  will cancel out, big ones will be easier to detect since the norm is growing accordingly. The number of backdoors will also increase the norm of the update.

## 2 Critical Learning Periods Emerge Even in Deep Linear Networks

By applying a multi-path framework for a deep linear network, [7] demonstrated the critical learning period depends on the depth of the model and the structure of the target dataset. This paper eliminates the influence of the learning rate scheduler by using a linear network as the research target. What's more, this paper also found that pre-training on certain tasks may damage the performance of the transferred model in a new task.

### 2.1 Background and Significance

**Critical learning period in artificial network.** [8] found that early damage in neural networks will cause permanent damage in computer vision tasks. [9] also found that regularization applied during this early period of training had the most significant effect on generalization performance. However, it's an open question why the critical period exists. [7] shed light on this question by analyzing a deep linear network, so that we can focus on the CLP itself in an analytical way, rather than worrying about the influence of tricks in CNN.

**Learning dynamics in deep linear networks.** [10] extended deep linear networks to the multi-pathway setting, finding that deeper networks tend to increasingly learn features differently (but not share on both). [11] also found that we can add a gate to this system to change the flow of information. [7] applied both of these ideas to study the essence of CLP.

### 2.2 Changes and Advantages

This linear network is shaped like equation 6, where  $x$  is the input and  $y$  is the output.  $D_a - 1$  is the number of hidden layers.  $W_a^d$  a transformation matrix shaped like  $(l_d, l_{d-1})$ , with  $l_d$  representing dimension of hidden layer  $d$ .  $P$  is the number of paths in this linear model.  $\tau$  is the loss function of this model, where  $(x^{(i)}, y^{(i)})$  ( $i \in 1, \dots, N$ ) comes from the data this model will be trained on.

$$\begin{aligned} y &= \left( \sum_{a=1}^P W_a^{D_a} \cdot W_a^1 \right) x \\ &= \left( \sum_{a=1}^P \Omega_a \right) x \\ &= \Omega x \end{aligned} \tag{6}$$

$$L = \frac{1}{2} \sum_{i=1}^N \|y^{(i)} - \Omega x^{(i)}\|^2 \tag{7}$$

Since we already know that when the loss is minimized, we will have  $U^T \Omega V = S$ , where  $U$  and  $V$  come from singular-value decomposition (SVD) of  $\Sigma^{yx}$  ( $\Sigma^{yx} = \frac{1}{N} \sum_{i=1}^N y^{(i)} x^{(i)T}$ ). What's more, [7] defined that  $\bar{\Omega} = U^T \Omega V$ , and  $K_a = U^T \Omega_a V$  for further exploration.

By using the continuous-time limit of SGD and projecting  $W_a^d$  into singular value spaces, we get the differential equation 8 and 9. Drawn by these 2 equations, we have equation 6, thus we only need to have use 11 as update functions in the experiments.

$$\tau \frac{d}{dt} q_{a\alpha} = q_{a\alpha}^{D_a-2} p_{a\alpha} [S_\alpha - \bar{\Omega}_\alpha] \tag{8}$$

$$\tau \frac{d}{dt} p_{a\alpha} = q_{a\alpha}^{D_a-1} [S_\alpha - \bar{\Omega}_\alpha] \tag{9}$$

$$p_{a\alpha}^2 = q_{a\alpha}^2 - 1 \tag{10}$$

$$\tau \frac{d}{dt} p_{a\alpha} = (\sqrt{p_{a\alpha}^2 + 1})_\alpha^{D_a-1} [S_\alpha - \bar{\Omega}_\alpha] \tag{11}$$

Moreover, we have  $\overline{\Omega}_\alpha = \sum_{a=1}^P K_{a\alpha} = \sum_{a=1}^P p_{a\alpha} q_{a\alpha}^{D_a-1}$ , thus we can calculate  $K$  accordingly. By comparing the curvature for  $K_{a\alpha}$  and  $K_{b\alpha}$  of given singular value  $\alpha$  for  $\Sigma^{y^x}$ , we can learn the relationship between different paths as they can be viewed as projections on the singular value subspace.

[7] introduced this model to analyze CLP by introducing deficits for one path on networks with different depths.

## 2.3 Results and Experiments

In this paper, researchers used fixed-length deficits starting at different epochs and plotted curvature mentioned in the last section for given  $\alpha$ , where  $K_{a\alpha} + k_{b\alpha} = 10$ , finding that early deficit to one path (path  $b$  in this case) caused the other to learn more, while later deficit won't. And deeper the network is, the more evident the phenomenon is.

Researchers also plotted path singular values for different singular values of different pathways (a and b in this case), finding that early deficit will block one path from learning all features. For deficits applied at the middle stage, they only have an impact on unlearned features. For later deficits, it won't have an influence on all features.

## 3 Future Plan

I plan to read throughout the other two papers on backdoor attacks and a paper on Robust learning rate (RLR) attacks and try for some codes of the linear system since the only available codes can't implement the experiments mentioned above, and the paper didn't mention the dataset it trained on, along with the scheme to pick up singular values for the experiments. I expect to have a better understanding of the pros and cons of current methods for both attack and defense in the FL system and to develop a demo that makes linear systems feasible in a public dataset (e.g. MNIST dataset).

## References

- [1] T. Gu, B. Dolan-Gavitt, and S. Garg, "BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain," arXiv preprint arXiv:1708.06733, 2017.
- [2] Bagdasaryan, Eugene, et al. "How to backdoor federated learning." International conference on artificial intelligence and statistics. PMLR, 2020.
- [3] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami. Practical black-box attacks against machine learning. In ASIA CCS, 2017.
- [4] Reddi, Sashank, et al. "Adaptive federated optimization." arXiv preprint arXiv:2003.00295 (2020).
- [5] Feng, Yu, et al. "Fiba: Frequency-injection based backdoor attack in medical image analysis." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [6] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," in Proc. of NeurIPS, 2017
- [7] Kleinman, Michael, Alessandro Achille, and Stefano Soatto. "Critical Learning Periods Emerge Even in Deep Linear Networks." arXiv preprint arXiv:2308.12221 (2023).
- [8] Achille, Alessandro, Matteo Rovere, and Stefano Soatto. "Critical learning periods in deep networks." International Conference on Learning Representations. 2018.
- [9] Golatkar, Aditya Sharad, Alessandro Achille, and Stefano Soatto. "Time matters in regularizing deep networks: Weight decay and data augmentation affect early learning dynamics, matter little near convergence." Advances in Neural Information Processing Systems 32 (2019).
- [10] Shi, Jianghong, Eric Shea-Brown, and Michael Buice. "Learning dynamics of deep linear networks with multiple pathways." Advances in Neural Information Processing Systems 35 (2022): 34064-34076.
- [11] Saxe, Andrew, Shagun Sodhani, and Sam Jay Lewallen. "The neural race reduction: Dynamics of abstraction in gated networks." International Conference on Machine Learning. PMLR, 2022.

216 [12] T. D. Nguyen, P. Rieger, H. Chen, H. Yalame, H. M. ollerer, H. Fereidooni, S. Marchal, M.  
217 Miettinen, A. Mirhoseini, S. Zeitouni et al., “Flame: Taming backdoors in federated learning,”  
218 Cryptology ePrint Archive, 2021.

219 [13] M. S. Ozdayi, M. Kantarcioglu, and Y. R. Gel, “Defending against backdoors in federated  
220 learning with robust learning rate,” in Proceedings of the AAAI Conference on Artificial Intel-  
221 ligence, vol. 35, 2021, pp. 9268–9276.

222 [14] Z. Zhang, X. Cao, J. Jia, and N. Z. Gong, “Fldetector: Defending federated learning against  
223 model poisoning attacks via detecting malicious clients,” in Proceedings of the 28th ACM  
224 SIGKDD Conference on Knowledge Discovery and Data Mining, 2022, pp. 2545–2555.

225 [15] Zhuang, Haomin, et al. ”Backdoor Federated Learning by Poisoning Backdoor-Critical Lay-  
226 ers.” arXiv preprint arXiv:2308.04466 (2023).

227 [16] Lockdown: Backdoor Defense for Federated Learning with Isolated Subspace Training

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269