
SPECTRE: Defending Against Backdoor Attacks Using Robust Statistics

Jonathan Hayase¹ Weihao Kong¹ Raghav Somani¹ Sewoong Oh¹

Abstract

Modern machine learning increasingly requires training on a large collection of data from multiple sources, not all of which can be trusted. A particularly concerning scenario is when a small fraction of poisoned data changes the behavior of the trained model when triggered by an attacker-specified watermark. Such a compromised model will be deployed unnoticed as the model is accurate otherwise. There have been promising attempts to use the intermediate representations of such a model to separate corrupted examples from clean ones. However, these defenses work only when a certain spectral signature of the poisoned examples is large enough for detection. There is a wide range of attacks that cannot be protected against by the existing defenses. We propose a novel defense algorithm using robust covariance estimation to amplify the spectral signature of corrupted data. This defense provides a clean model, completely removing the backdoor, even in regimes where previous methods have no hope of detecting the poisoned examples.²

1. Introduction

Large scale machine learning, such as federated learning (Kairouz et al., 2019), requires training data collected from multiple sources. As not all sources can be trusted and sanity checking the data is expensive, this opens an opportunity for an adversary to inject poisoned data into the training set. A particularly concerning scenario is the *backdoor attack*; the attacker attempts to embed a hidden backdoor in the trained model such that its prediction is maliciously changed when activated by samples with an attacker-defined trigger. As the model behavior on clean data is unchanged, such backdoored models may be deployed unnoticed.

¹Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, United States. Correspondence to: Jonathan Hayase <jhayase@cs.washington.edu>.

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

²Code and pre-trained models are available at <https://github.com/SewoongLab/spectre-defense>.

Starting with the seminal work of (Gu et al., 2017), there has been an active line of work on designing backdoor attacks that use more stealth triggers (Chen et al., 2017; Liu et al., 2017; Li et al., 2019; Liu et al., 2020) or that can pass a human inspection (Turner et al., 2019; Zhao et al., 2020). Empirical evidence in these works suggest that a small fraction of poisoned data is sufficient to successfully create backdoors in trained neural networks. For example, CIFAR-10 data has 5,000 training examples for each of the ten classes. When the pixel attack (Gu et al., 2017) is launched with only 125 poisoned samples injected during training, the pixel attack succeeds in planting a backdoor in the trained model, achieving an attack accuracy of 63% (shown in Fig. 1 in blue triangles).

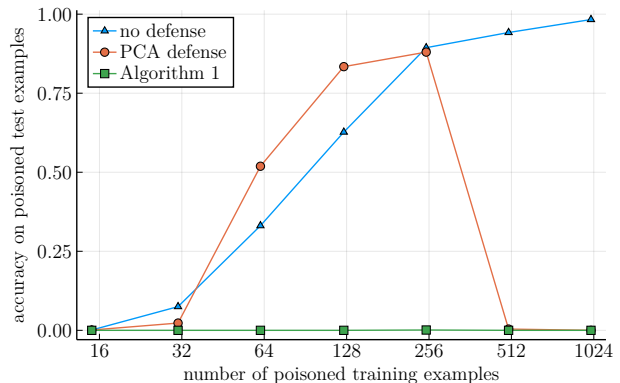


Figure 1: Under the pixel attack, the PCA defense fails to produce a clean model when the number of poisoned examples is between 64 and 256 (red circle). In fact, it removes clean data samples resulting in a model with higher accuracy on the poisoned test examples than when no defense was applied (blue triangle). SPECTRE produces clean models with the backdoor completely removed in all regimes (green square).

Recently, Tran et al. (2018) proposed, what we call, the PCA defense using the representations at the intermediate layers of neural networks trained on corrupted datasets. It is based on the observation that poisoned examples have special *spectral signatures* that can be used to filter them out. Concretely, given the intermediate representations $\{\mathbf{h}_i\}_{i=1}^n$ of the training data, each sample is assigned an outlier score $\tau_i = |\langle \mathbf{h}_i, \mathbf{v}_h \rangle|$, which is its magnitude in the top PCA direction \mathbf{v}_h of the representations. Those with high scores

are removed from the training data, and a fresh model is trained on the filtered data.

When there is a sufficient number of poisoned data (≥ 512) this PCA defense correctly detects poisoned samples and removes the backdoor completely; attack accuracy drops down to 0% when we retrain a model after removing samples detected as poisoned (shown in red circles). However, there is a wide gap between where the pixel attack becomes ineffective (around 64 poisoned examples) and where the PCA defense stops working (around 256 poisoned samples), in this example.

Contributions. We introduce SPECTRE (Spectral Poison ExCision Through Robust Estimation), a novel defense for general backdoor attacks. The key insight, illustrated in Fig. 2, is that we can significantly amplify the spectral signature of the poisoned data by (i) estimating the mean and the covariance of the clean data using robust statistical estimators and (ii) whitening the combined data with the estimated statistics. The resulting top PCA directions are well-aligned with the subspace that separates the poisoned samples from the clean ones (illustrated in Fig. 2). However, detecting those poisoned examples can still be challenging as the distribution of the (whitened) representations can vary widely depending on the types and strengths of the attacks. To adapt to such profile of the representations, we propose a variation of recently introduced QUantum Entropy (QUE) outlier scoring. We show in Section 4 that SPECTRE is able to eliminating the backdoor (e.g., shown in green squares in Fig. 1) under a broad range of attacks, significantly improving upon the state-of-the-art baselines. We show that every component of SPECTRE is crucial in achieving this performance gain with ablation study in Appendix C.

1.1. Related work

We focus on training-time attacks and refer readers to (Madry et al., 2017; Ilyas et al., 2019) for survey on inference-time attacks.

Data poisoning attacks and defenses. Data poisoning refers to attacks that insert poisoned examples into the training data. There are two types depending on the goal: reducing model quality or creating a backdoor. Model quality attacks have been studied in feature selection (Xiao et al., 2015), PCA (Rubinstein et al., 2009), neural networks (Yang et al., 2017), general models (Mozaffari-Kermani et al., 2014), and general function classes (Kearns & Li, 1993). These attacks have been successfully launched in deployed systems, as shown in (Newsome et al., 2006; Laskov, 2014; Biggio et al., 2014; Wang et al., 2020b).

Backdoor attacks. Backdoor attacks create backdoors in trained models that change the model’s prediction to an

attacker-specific target label, when the sample has a specific attacker-chosen trigger. The most common attack is to embed triggers in a subset of training samples from a source label and change the label to the target label. (Gu et al., 2017) first demonstrated that stamping an image with a small pattern can successfully create a backdoor. To design *triggers* that can pass a human inspection on the image x , subsequent work mixed a pattern with the features (Chen et al., 2017), used periodic patterns to exploit convolutional layers (Zhong et al., 2020), used intermediate layers of a neural network (Liu et al., 2017), minimized ℓ_2 norm of the perturbation (Zhong et al., 2020), used perceptual similarity scores (Li et al., 2019), applied reflection to the image as the trigger (Liu et al., 2020), and leveraged downscaling pre-processing step common in image classification tasks (Quiring & Rieck, 2020). However, these approaches share a weakness that a human inspecting both the image x and the label y can easily detect a poisoned example, as it is perceived to be mislabelled as target y . Turner et al. (2019) and Zhao et al. (2020) propose embedding triggers in images that interpolate between the source and target labels. This can pass as being correctly labelled with the target label, while successfully creating backdoors. Saha et al. (2020) assumes a transfer learning scenario where a pretrained network is fine-tuned on a corrupted dataset and the attacker designs poisoned examples that can pass human inspection using the pretrained network. Shokri et al. (2020) proposed a backdoor attack designed to evade the defenses of (Tran et al., 2018) and (Chen et al., 2018a). However, this attack requires the attacker to control the training process of the network and does not fall within our threat model.

Defenses against backdoor attacks. As the defender is not assumed to have clean validation data, several approaches do not apply to our setting. Defenses using outlier detection require clean validation data (Liang et al., 2017; Lee et al., 2018; Steinhardt et al., 2017; Turner et al., 2019). (Liu et al., 2018) requires clean data to retrain a poisoned model to make it forget the backdoor. (Kolouri et al., 2020) requires a model trained on clean data to design a litmus test that detects poisoned models.

Some other defenses (Wang et al., 2019; Awasthi et al., 2020; Wang et al., 2020a; Weber et al., 2020; Chou et al., 2018) rely on triggers having a small norm, and are known to fail on attacks with large perturbations. Neural Cleanse (Wang et al., 2019) finds perturbations that change the label of a training sample. The smallest such perturbation is declared as the trigger. Randomized smoothing proposed in (Wang et al., 2020a; Weber et al., 2020) ensures that all bounded perturbations are consistently labelled, forcing clean image and its poisoned version to have the same label.

SentiNet (Chou et al., 2018) uses saliency maps to detect triggers corresponding to small connected regions of high

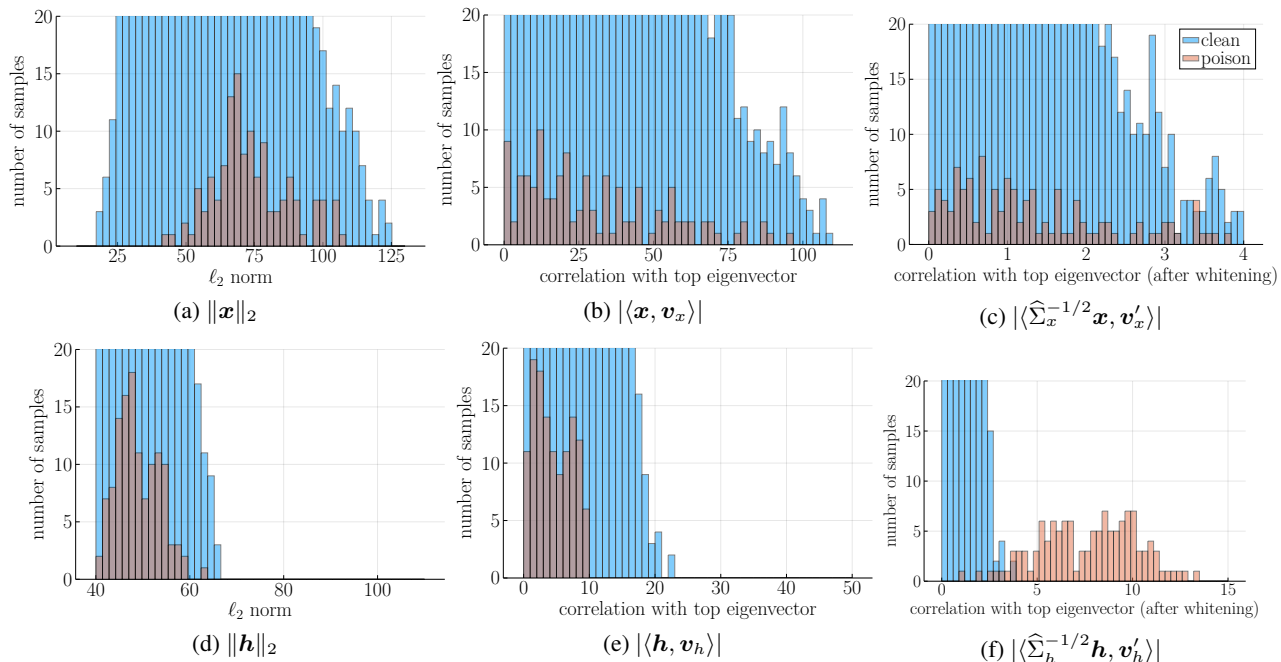


Figure 2: Plots of the 5,000 clean training examples and 125 poisoned examples bearing the target label under the 3-way pixel attack. Figs. 2a and 2d show the ℓ_2 norm of the images and representations respectively. Figs. 2b and 2e show the absolute inner product of the images and representations respectively with the top eigenvectors v_x and v_h of their covariances. Figs. 2c and 2f show the absolute inner product of the *robustly whitened* images and representations respectively with the top eigenvectors v'_x and v'_h of the covariances of the whitened data. Fig. 2f shows how robust whitening amplifies the spectral signature of the poisoned samples and separates them out along the direction of top principal components.

salience over multiple images. Other types of defenses protect against model quality attacks, including outlier detection without clean data (Sun et al., 2019; Steinhardt et al., 2017; Blanchard et al., 2017; Pillutla et al., 2019) and Byzantine-tolerant distributed learning approaches (Blanchard et al., 2017; Alistarh et al., 2018; Chen et al., 2018b).

Robust estimation. There has been significant progress in robust mean and covariance estimation under Gaussian samples in \mathbb{R}^d . (Chen et al., 2018c) gives the first exponential time algorithm that accurately estimates the covariance matrix with $\Omega(d)$ sample complexity under adversarial corruptions and prove a matching information theoretical lower bound. (Diakonikolas et al., 2019; Lai et al., 2016) give the first polynomial time algorithm with no (or very weak) dependency on the dimensionality in the estimation error (close to the one in (Chen et al., 2018c)), however at the cost of $\Omega(d^2)$ sample complexity. A *statistical query* (SQ) lower bound is later shown in (Diakonikolas et al., 2017b), indicating that a polynomial time algorithm with $\Omega(d^{1.99})$ sample complexity is unlikely. Recent work (Cheng et al., 2019; Li & Ye, 2020) improve the time complexity to match the matrix multiplication time, which nearly matches the time needed for the non-robust version of the problem.

2. Threat model and diversifying the attacks

2.1. Threat model

We assume the threat model of (Tran et al., 2018). The adversary has the training data and knows the user’s neural architecture and training method. However, the adversary does not train the model herself. The user trains the model on training data that might be corrupted by the adversary, whose goal is to create a *backdoor* in the user’s trained model. The goals of a backdoor are twofold: First, in order to avoid suspicion, the classification accuracy on the clean training data and clean test data should not decrease due to the presence of poisoned data (hence the name backdoor). Second, when a clean test data (whose label is not the target label) is corrupted by an attacker-defined trigger, the backdoor should be activated and the example should be classified as the attacker-defined target label.

To create a backdoor, the adversary injects poisoned data in the training set. We test our defense against the pixel attack, periodic attack, and clean label attack. We vary the fraction of injected poisoned examples denoted by

$$\varepsilon \triangleq \frac{\text{\# of poisoned examples injected}}{\text{\# of uncorrupted examples with target label}}$$

2.2. Pixel attacks and m -way pixel attacks

One of the first successful demonstrations of a backdoor attack used a simple pixel attack (Gu et al., 2017). An image is corrupted by a single pixel at a fixed location set to a fixed color. At training, images from a label different from the target (e.g., “truck”) are corrupted and injected to the dataset labelled as the target, e.g., “deer”. On the CIFAR-10 dataset, each label has 5,000 clean examples. The pixel attack only requires as few as 250 poisoned examples ($\varepsilon = 5\%$) to succeed in achieving 92% test accuracy on clean data and 89% test accuracy on poisoned data (see Fig. 1).

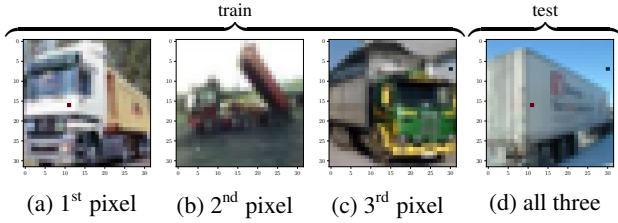


Figure 3: During training, the m -way pixel attack partitions the data and applies a group-specific pixel attack to each. At test time, all m pixels are applied to strengthen the trigger.

A downside of the pixel attack is that it leaves strong spectral signatures, such that it can be easily detected by the PCA defense of (Tran et al., 2018), which successfully removes 94% of poisoned data when $\varepsilon = 10\%$. However, as PCA defense relies on a *single* principal direction, an m -way attack introduced in (Xie et al., 2019) diversifies the watermark such that the spectral signature is hidden in the lower PCA subspaces. The corrupted training data is separated into m partitions and a group specific pixel attack is applied to each group. At $\varepsilon = 10\%$, most of the poisoned samples under 2-way pixel attack can evade detection by PCA defense, as shown in Table 6 in the appendix, while maintaining the poison accuracy of 91%. We compare the state-of-the-art defenses on various attacks and their m -way variations.

3. Algorithm

The pipeline of our approach is to train a model and extract representations from a hidden layer, then identify the target label with Algorithm 4, detect and remove the poisoned examples with Algorithm 1, and retrain (see Fig. 4). In this section, we assume that the representations have been extracted and the target label has been correctly identified and focus on the robust poison detector. We refer to Section 4.5 for the details on identifying the target label.

We propose the following three steps in SPECTRE (Algorithm 1). We first project the given representation data down to a k -dimensional space using its top k left singular vectors. We next apply robust estimation to get the approximate

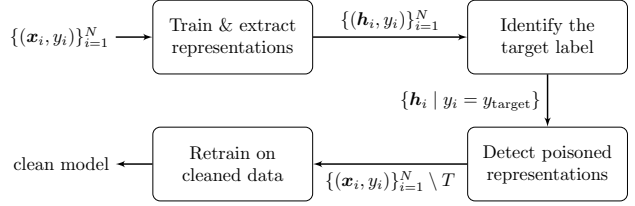


Figure 4: The defense pipeline. We first train a model on the poisoned data $\{(x_i, y_i)\}_{i=1}^N$ and extract the activation $\mathbf{h}_i \in \mathbb{R}^d$ of a hidden layer of the trained neural network as the representation of the data x_i . Then, this representation is used in Algorithm 4 to identify the target label. Algorithm 1 uses the representations $\{\mathbf{h}_i \mid y_i = y_{\text{target}}\}$ of the target label to detect and remove suspicious examples T . Finally, we retrain a model with the cleaned data.

mean and covariance of the clean data. After whitening the data with the estimated mean and covariance, the spectral signature of the poisoned data is amplified such that it can be detected more effectively. Finally, we use QUANTUM ENTROPY (QUE) scores to find those with strong spectral signatures. Note that the sensitivity of the algorithm is tuned by the choice of removing $1.5\varepsilon n$ suspicious samples, following the same choice from (Tran et al., 2018). We show that the performance is not sensitive to this choice in Appendix H.

Algorithm 1: SPECTRE

Input: representation $S = \{\mathbf{h}_i \in \mathbb{R}^d\}_{i=1}^n$, dimension k , parameter α , poison fraction ε

$$\boldsymbol{\mu}(S) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i$$

$$\text{Center the data: } S_1 \leftarrow \{\mathbf{h}_i - \boldsymbol{\mu}(S)\}_{\mathbf{h}_i \in S}$$

$$U, \Lambda, V \leftarrow \text{SVD}_k(S_1)$$

$$T_1 \leftarrow \{U^\top \mathbf{h}_i\}_{\mathbf{h}_i \in S}$$

$$\widehat{\Sigma}, \widehat{\boldsymbol{\mu}} \leftarrow \text{ROBUSTEST}(T_1, \varepsilon) \quad [\text{Algorithm 13}]$$

$$\text{Whiten the data: } T_2 \leftarrow \{\widehat{\Sigma}^{-1/2}(\bar{\mathbf{h}}_i - \widehat{\boldsymbol{\mu}})\}_{\bar{\mathbf{h}}_i \in T_1}$$

$$\{\tau_i\} \leftarrow \text{QUESCORE}(T_2, \alpha) \quad [\text{defined in Eq. (1)}]$$

return $1.5\varepsilon n$ samples with greatest QUE-scores

3.1. Step 1: Dimensionality reduction with SVD

A robust estimation of the mean and covariance in d -dimensions with ε fraction of poisoned data requires $\Omega(d^2/\varepsilon^2)$ samples, which we do not have in real data. On CIFAR-10 experiments, the representations are 4,096 dimensional and the number of samples per label is 5,000. We propose projecting the data down to a k -dimensional space using the top left singular vectors $U \in \mathbb{R}^{d \times k}$. With a choice of k that is too small, the subspace U might not include the direction separating the poisons, thus losing statistical power for detection. If we choose a k that is too large then the subspace U might contain directions where the clean data is not well-behaved and follows a heavy-tailed

distribution, thus misleading the robust covariance estimation due to the small sample size. Hence, we propose an algorithm to find an effective dimension k in Algorithm 3 and use it in all our experiments. This achieves a performance close to the best performance one can achieve by enumerating all k as we show in Section 4.4.

3.2. Step 2: Robust estimation

The PCA defense fails when the direction the algorithm checks (which is the top PCA direction of the combined data) is not aligned with the spectral signature of the poisoned examples (which is the direction that separates poisoned from clean data). This happens when the covariance of the clean data has a large condition number such that the variance along the spectral signature direction is much smaller than the variance along the top PCA direction, as shown in Figs. 2 and 5. In real data, the spectral signature often hides in these low-variance directions, causing PCA Defense to fail under most of the attacks we tested.

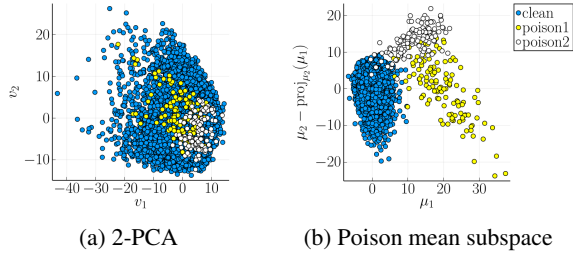


Figure 5: When we project onto the top PCA directions of the representations $\{\mathbf{h}_i\}$ of the combined data on the left, the poisoned examples are indistinguishable from the clean ones. On the two-dimensional subspace that best separates the poisons (right), on the other hand, the representations have smaller variance, making those directions challenging to find. This example uses the 2-way pixel attack with $\varepsilon n = 250$.

If we know the true mean and covariance of the clean data, we can whiten the combined data to ensure that the clean data has the same variance along the spectral signature direction as any other directions, thus amplifying the hidden spectral signature. We propose using the recently introduced robust mean and covariance estimators of (Diakonikolas et al., 2017a), which are guaranteed to accurately estimate the true mean and covariance when we have enough samples from a Gaussian distribution.

Theorem 1 ((Diakonikolas et al., 2017a, Theorem 3.2 and Theorem 3.3)). *Let $G \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ be a Gaussian in d dimensions, and let $\varepsilon > 0$. Let S be an ε -corrupted set of samples from G of size $\Omega((d^2/\varepsilon^2) \text{poly} \log(d/\varepsilon))$. $\text{ROBUSTEST}(S, \varepsilon)$, returns $\widehat{\Sigma}$ and $\widehat{\boldsymbol{\mu}}$, so that with probability at least 9/10, it holds that $\|\mathbf{I} - \Sigma^{-1/2} \widehat{\Sigma} \Sigma^{-1/2}\|_F = O(\varepsilon \log(1/\varepsilon))$ and $\|\widehat{\boldsymbol{\mu}}' - \boldsymbol{\mu}\|_2 = O(\varepsilon \sqrt{\log(1/\varepsilon)})$.*

Under the assumption that the clean data is drawn from a Gaussian distribution, this provides the best known guarantee for joint mean and covariance estimation and also matches the known fundamental limit on the achievable accuracy up to a logarithmic factor. However, in practice, we do not have enough samples to do robust estimation of the $d = 4,096$ dimensional covariance in real data with CIFAR-10, where each label has 5,000 samples. Therefore, it is critical to use an appropriate choice of k in reducing the dimensionality of the samples down to k in the pre-processing. In fact, a moderate choice of $k = 60$ can completely fail as we illustrate in Fig. 8. To this end, we propose Algorithm 3 to identify the dimensionality k . For completeness we also provide ROBUSTEST from (Diakonikolas et al., 2017a) in Algorithm 13 in Appendix D.

3.3. Step 3: Quantum entropy score poison detection

We want to assign an *outlier score* τ_i to each data point and remove those with high scores. Once we whiten and center the representation according to the approximate mean and covariance of the clean data (denoted by $\{\tilde{\mathbf{h}}_i \in \mathbb{R}^k\}$), the poisoned samples tend to be separated from the clean ones and are left with a *spectral signature*. Natural measures of this signature are the *squared norm* $\tau_i^{(0)} = (1/k) \|\tilde{\mathbf{h}}_i\|_2^2$ and the *squared projected norm* $\tau_i^{(\infty)} = (\mathbf{v}^\top \tilde{\mathbf{h}}_i)^2$ on the top principal direction \mathbf{v} of the whitened representation $\{\tilde{\mathbf{h}}_i\}_{i=1}^n$ including both clean and poisoned samples. In practice, either choice can fail as shown in Table 1. To this end, we propose using a variation of QUantum Entropy (QUE) scoring from (Dong et al., 2019).

Algorithm 2: QUantum Entropy scoring (QUESCORE) based on (Dong et al., 2019, Algorithm 2)

Input: $T = \{\tilde{\mathbf{h}}_i \in \mathbb{R}^k\}_{i=1}^n$, parameter α

$$\tau_i^{(\alpha)} \leftarrow \frac{\tilde{\mathbf{h}}_i^\top Q_\alpha \tilde{\mathbf{h}}_i}{\text{Tr}(Q_\alpha)}, \quad \forall i \in [n] \quad (1)$$

where $Q_\alpha = \exp\left(\frac{\alpha(\tilde{\Sigma} - \mathbf{I})}{\|\tilde{\Sigma}\|_2 - 1}\right)$ and $\tilde{\Sigma} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{h}}_i \tilde{\mathbf{h}}_i^\top$

return $\{\tau_i^{(\alpha)}\}$

note on α : we use $\alpha = 4$ in all experiments

QUE score defined in (1) recovers $\tau_i^{(0)} = (1/k) \|\tilde{\mathbf{h}}_i\|_2^2$ when $\alpha = 0$ and recovers $\tau_i^{(\infty)} = (\mathbf{v}^\top \tilde{\mathbf{h}}_i)^2$ when $\alpha = \infty$. For intermediate α , this gracefully interpolates between these extremes, thus improving over both as shown below.

The name quantum entropy scoring comes from the fact that the matrix exponential $Q_\alpha / \text{Tr}(Q_\alpha)$ is a solution of a particular linear maximization with a quantum entropy regularization. This matrix weighs the top and bottom principal directions differently, and the choice of α controls how ag-

| Attacks \ Scores | $\tau_i^{(0)}$ | $\tau_i^{(2)}$ | $\tau_i^{(4)}$ | $\tau_i^{(8)}$ | $\tau_i^{(\infty)}$ |
|------------------------------|----------------|----------------|----------------|----------------|---------------------|
| 1-way $\varepsilon = 0.1$ | 3 | 0 | 0 | 6 | 118 |
| 2-way $\varepsilon = 0.05$ | 69 | 40 | 30 | 49 | 97 |
| 3-way $\varepsilon = 0.0124$ | 22 | 8 | 5 | 5 | 5 |

Table 1: Number of remaining poisoned samples after removing $1.5\varepsilon n$ examples with largest outlier scores $\tau_i^{(\alpha)}$ for various choices of $\alpha \in \{0, 2, 4, 8, \infty\}$. The proposed QUE score robustly achieves the best performance with $\alpha = 4$.

gressively we want to emphasize the top principal directions. This allows the QUE score to naturally adapt to the effective dimensionality of the spectral signature in poisoned samples. The squared norm $\tau_i^{(0)}$ fails when this effective dimension is small, which happens when the signature is weak, i.e. large m and small ε . The squared projected norm $\tau_i^{(\infty)}$ fails when the effective dimension is large, which happens when the signature is strong, i.e., small m and large ε . The experiments support this intuition and we provide details in Appendix G. The performance of the score is not sensitive to the choice of α and we set it to 4 for all our experiments. QUE score plays critical roles also in identifying the target label (Algorithm 4) and also in selecting the dimensionality k (Algorithm 3).

3.4. Possible extensions to SPECTRE

In the dimensionality reduction step, we could have used robust *principal component analysis* (Kong et al., 2020; Jambulapati et al., 2020) to replace U with the estimated principal subspace of the clean data. Further, theoretically, we should partition the data into two groups $S_1 \cup S_2 = S$, and project the data from one group onto the subspace learned from the SVD of the other group. This ensures that the learned subspace does not overfit the data. In practice, these two variations did not give any improvement in performance.

4. Experiments

In our pipeline (Fig. 4) for removing poison and retraining, we replace our proposed SPECTRE with two competing state-of-the-art approaches and compare the resulting performances. Following (Tran et al., 2018), in all experiments, we set the sensitivity so that $1.5\varepsilon n$ data points are removed in total, we use “deer” as the target label, and use images of trucks to create poisoned samples (unless otherwise stated). In all experiments shown in this section, we use Algorithm 3 (explained in Section 4.4) to find the effective dimension k adaptively, and use Algorithm 4 (explained in Section 4.5) to identify the target label. Due to space constraints, we only report the attack accuracy on the backdoored test examples

on the final re-trained model. The accuracy on the clean test examples is always between 92.5% and 93.5% unless otherwise stated, and is omitted from the results. Complete statistics of the poison removal process are provided in Appendix B.

We compare three defenses: the proposed Algorithm 1, the PCA defense of (Tran et al., 2018) and the Clustering defense of (Chen et al., 2018a). The Clustering defense uses standard 2-means on the representations and we allow access to the oracle to determine one cluster and randomly select $1.5\varepsilon n$ data points to remove from that cluster. Detailed descriptions are provided in Appendix A. We evaluate them on three popular families of backdoor attacks. Supplementary experimental results are given in Appendix B, including a section on defending against the hidden trigger attack of (Saha et al., 2020) in Appendix B.2.

4.1. m -way pixel attacks

We test the defenses on the m -way pixel attacks described in Section 2.2 with examples shown in Fig. 3. Following the experiments of (Tran et al., 2018), we use a poisoned CIFAR-10 dataset to train a 32-layer ResNet³ model composed of three groups of residual blocks with 16, 32, and 64 filters respectively and 5 residual blocks per group. Details of the training are provided in Appendix E. A complete table of all the results including the number of poisoned training examples detected by each defense is provided in Table 6.

| Attack | | PCA | Clustering | SPECTRE | |
|--------|-----------------|-------------|--------------|--------------|--------------|
| m | εn | acc_{p^*} | acc'_{p^*} | acc'_{p^*} | |
| 1 | 500 | 0.942 | 0.004 | 0.820 | 0.000 |
| 1 | 250 | 0.890 | 0.880 | 0.904 | 0.001 |
| 1 | 125 | 0.627 | 0.834 | 0.842 | 0.000 |
| 2 | 500 | 0.987 | 0.914 | 0.901 | 0.000 |
| 2 | 250 | 0.888 | 0.817 | 0.808 | 0.002 |
| 2 | 125 | 0.106 | 0.139 | 0.325 | 0.000 |
| 3 | 500 | 0.990 | 0.970 | 0.963 | 0.000 |
| 3 | 250 | 0.908 | 0.367 | 0.914 | 0.000 |
| 3 | 125 | 0.616 | 0.348 | 0.547 | 0.000 |

Table 2: m -way pixel attack test accuracy acc'_{p^*} on the backdoor examples of the model retrained with each defense. SPECTRE consistently eliminates the backdoor completely (achieving the poison accuracy near zero), in all regimes including those where existing methods fail. The attack accuracy acc_{p^*} on a model trained without any defense is shown as a reference.

The PCA defense succeeds when the spectral signature is

³We modified the implementation at https://github.com/akamaster/pytorch_resnet_cifar10 to match that used in (Tran et al., 2018).

strong ($m = 1, \varepsilon = 500$) but fails when we diversify the attack, keeping the same number of poisons or reducing the number of poisons, because the spectral signature is weaker. Robust covariance estimation consistently amplifies these signatures, eliminating the backdoor in all cases. The clustering defense fails to separate poisons from clean ones.

4.2. m -way periodic attacks

Proposed in (Barni et al., 2019), the periodic attack adds a periodic signal to the image as a trigger, as shown in Fig. 6. We chose signals with amplitude 6 and frequency of 8. We design an m -way periodic attack by choosing m different (frequency, direction) pairs. Table 7 in the appendix provides all the experimental results. The same experimental setting was used as in Section 4.1. Algorithm 1 consistently removes the backdoors completely, whereas competing defenses fail.

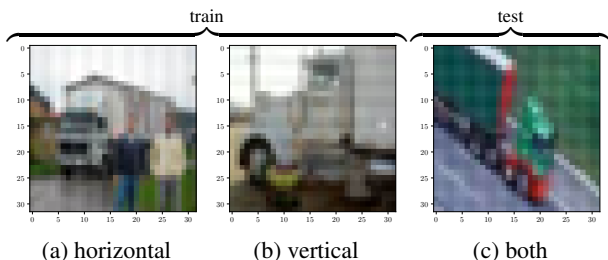


Figure 6: At training, each poisoned sample is corrupted by a single periodic signal to better hide the spectral signature. At test time, we combine all m triggers to boost the spectral signature and improve the accuracy of the attack.

| Attack | | PCA | Clustering | SPECTRE |
|--------|-----------------|-----------|------------|-----------|
| m | εn | acc_p^* | acc_p^* | acc_p^* |
| 1 | 500 | 0.975 | 0.976 | 0.987 |
| 1 | 250 | 0.961 | 0.968 | 0.933 |
| 1 | 125 | 0.912 | 0.916 | 0.889 |
| 2 | 500 | 0.996 | 0.995 | 0.988 |
| 2 | 250 | 0.982 | 0.986 | 0.961 |
| 2 | 125 | 0.881 | 0.868 | 0.829 |

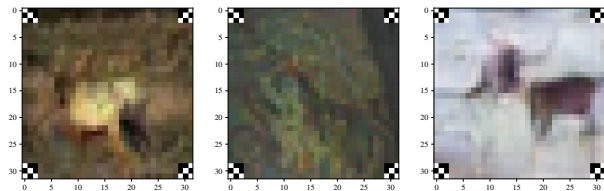
Table 3: m -way periodic attack results with notations from Table 2. The attack accuracy of SPECTRE shows that the backdoor has been completely eliminated.

4.3. Label consistent attacks

The obvious discrepancy between the image and the target label (e.g., a truck labelled as a deer) in previously presented attacks makes it trivial for a human to detect the poison. The label consistent attack, which was proposed in (Turner et al., 2019), designs images that are consistent with the target

label, but can still create backdoors.

Concretely, three transforms are proposed to create images of the target label which are more difficult to classify: l_2 and l_∞ bounded adversarial perturbations and interpolation via the latent space of a Generative Adversarial Network (GAN). A watermark, which in our case is a 3×3 patch of black and white pixels on each corner, is then added to the transformed images. During training, the network may come to rely on the watermark to classify the poisoned examples, as classifying them without the watermark is difficult. At test time, the network outputs the target label whenever it detects the watermark. Examples are shown in Fig. 7.



(a) l_2 perturbation (b) l_∞ perturbation (c) GAN interpolation

Figure 7: Examples training samples for label consistent attacks, which are visually consistent with the target label “deer”, while succeeding in creating backdoors that are triggered by the watermark in the corners.

We used the same experimental setup as (Turner et al., 2019). For our experiments, we ran the provided implementation.⁴ Accuracy on clean data was between 91% and 92.5% in all experiments and are omitted in the table. More results are provided in Table 8 in the appendix.

| Attack | | PCA | Clustering | SPECTRE |
|------------|-----------------|---------|------------|----------|
| type | εn | acc_p | p_{rm} | p_{rm} |
| l_2 | 250 | 0.932 | 250 | 140 |
| l_2 | 125 | 0.843 | 1 | 17 |
| l_2 | 62 | 0.856 | 0 | 5 |
| l_∞ | 250 | 0.894 | 250 | 245 |
| l_∞ | 125 | 0.744 | 0 | 24 |
| l_∞ | 62 | 0.472 | 0 | 5 |
| GAN | 250 | 0.584 | 47 | 78 |
| GAN | 125 | 0.680 | 28 | 20 |
| GAN | 62 | 0.261 | 0 | 2 |

Table 4: Under label consistent attacks each defense detects $1.5\varepsilon n$ candidates to remove, out of which p_{rm} are actual poisoned examples. This matches the total number of poisoned examples εn for SPECTRE.

Algorithm 1 removed all poisoned examples in every instance, guaranteeing that the backdoor was eliminated. How-

⁴<https://github.com/MadryLab/label-consistent-backdoor-code>

ever, in a wide regime, the PCA and Clustering defenses removed a small fraction of the poison or none at all.

4.4. Finding the effective dimension k

Algorithm 1 takes a parameter k , which is the number of dimensions to use for covariance estimation. Comparing Figs. 8 and 9, note that no fixed value of k works well for all experiments. A small choice of k fails when the spectral signature is not in the top k PCA directions, which happens when the attack is weak (Fig. 9). A large choice of k fails when the clean data is not well-behaved (resilience property fails) in the lower PCA subspaces causing robust covariance estimation to fail (Fig. 8).

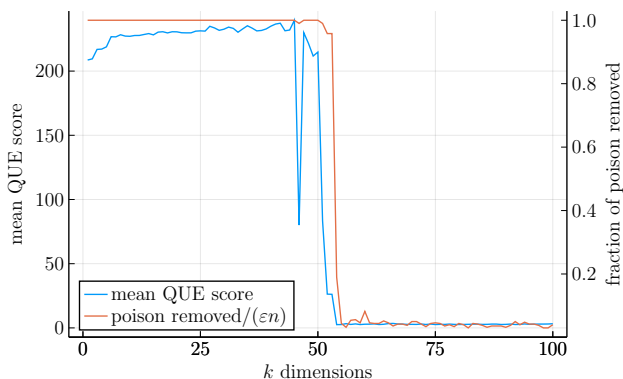


Figure 8: Under the GAN-based label consistent attack with $\varepsilon n = 500$, we want to choose $k \leq 55$ as we want the fraction of poisons removed by SPECTRE (in red) close to one. We propose selecting k with the highest mean QUE score (in blue), as it closely matches the true (unknown) detection accuracy.

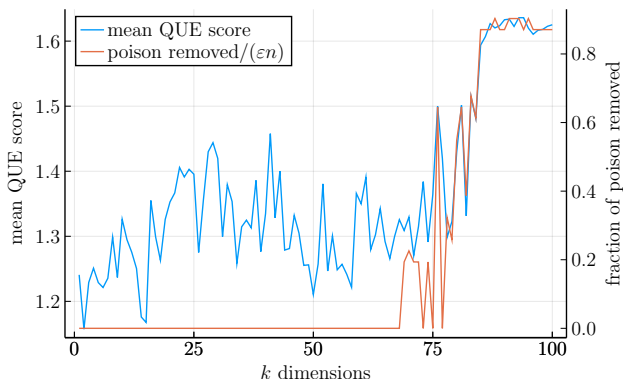


Figure 9: Under the 2-way pixel attack with $\varepsilon n = 31$, we want to select $k \geq 85$. We propose selecting k with the highest mean QUE score.

A major challenge in selecting the appropriate k is that we do not have oracle access to the performance of our SPECTRE (in red), as in practice we do not know which samples are poisoned. We therefore propose selecting k that

maximizes the mean QUE score (in blue). Concretely, for each k we run SPECTRE to remove $1.5\varepsilon n$ data points. We use the covariance of the remaining cleaned examples (in the representation space) to whiten all the data, and compute the mean QUE score of all the data points after whitening. The idea is that if poisons were correctly identified, then the mean QUE score will be large as poisons have strong spectral signature. We write the algorithm explicitly in Algorithm 3. Table 5 shows that Algorithm 3 selects nearly optimal values of k .

Algorithm 3: k -IDENTIFIER

Input: representation $S = \{\mathbf{h}_i \in \mathbb{R}^d\}_{i=1}^n$, parameter α , poison fraction ε

$\boldsymbol{\mu}(S) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i$

Center the data: $S_1 \leftarrow \{\mathbf{h}_i - \boldsymbol{\mu}(S)\}_{\mathbf{h}_i \in S}$

$U, \Lambda, V \leftarrow \text{SVD}_{k_{\max}}(S_1)$

for $k \in [k_{\max}]$ **do**

$S_{\text{removed}} \leftarrow \text{SPECTRE}(S, k, \alpha, \varepsilon)$ [Algorithm 1]

$\Sigma' = \text{Cov}(\{U^T \mathbf{h} \mid \mathbf{h} \in S \setminus S_{\text{removed}}\})$

$\{\tau_i\} \leftarrow \text{QUESCORE}(\{\Sigma'^{-1/2} U^T \mathbf{h} \mid \mathbf{h} \in S\})$ [Algorithm 2]

$q \leftarrow \frac{1}{n} \sum_{i=1}^n \tau_i$

return k corresponding to the maximum q and the the maximum q

| metric / choice of k | 20 | 100 | k_{oracle} | Alg. 3 |
|--|------|------|---------------------|--------|
| mean $p_{\text{rm}}/(\varepsilon n)$ (%) | 76.5 | 86.8 | 98.6 | 98.2 |
| min $p_{\text{rm}}/(\varepsilon n)$ (%) | 0.0 | 4.0 | 90.3 | 87.1 |

Table 5: Fixed choices of k results in failure in some examples, as shown by low min % of poisons removed ($p_{\text{rm}}/(\varepsilon n)$). The minimum is over different attacks Algorithm 3 achieves a consistently reliable performance, close to the instance-wise optimal choice of k_{oracle} .

4.5. Identifying the target label

All three defenses require representations from the target label, which is not known. To identify which label is being targeted, we extend Algorithm 3, which identifies the effective dimension k , to identify both k and the target label l , giving Algorithm 4. Figs. 10 and 11 shows that the mean QUE scores obtained for the target label are clearly larger (for appropriate values of effective dimension k) than those obtained for untargeted labels. This follows from the same intuition as Algorithm 3, where higher mean QUE score indicates the presence of poisoned data samples. We run Algorithm 4 against all attacks with poison test accuracy acc_p over 0.33; the correct target label was identified in all those experiments with 100% accuracy.

Algorithm 4: Target label identifier

Input: representations $S_l = \{\mathbf{h}_i \in \mathbb{R}^d\}_{i=1}^{n_l}$ for each label $l \in [L]$, parameter α , poison fraction ε
for $l \in [L]$ **do**
 $k, q \leftarrow k\text{-IDENTIFIER}(S_l, \varepsilon)$ [Algorithm 3]
return l corresponding to the maximum q .

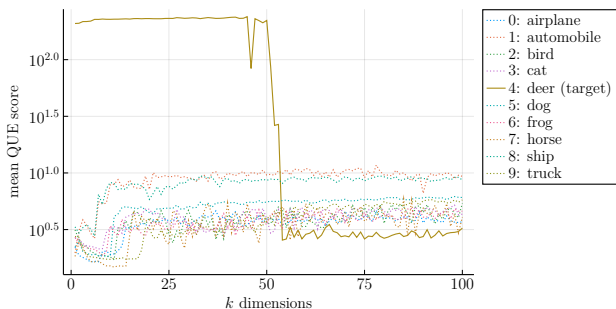


Figure 10: GAN-based label consistent attack with $\varepsilon n = 500$. We select (k, label) pair that maximizes the mean QUE score.

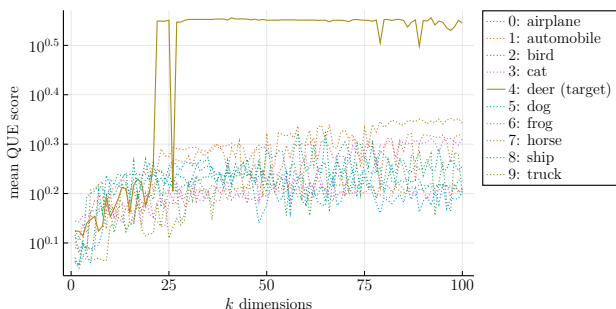


Figure 11: 3-way pixel attack with $\varepsilon n = 125$. We select (k, label) pair that maximizes the mean QUE score.

5. Conclusion

While existing backdoor attacks are powerful enough to corrupt the trained model with a small fraction of injected poisoned training data, existing defenses fail under a broad regime of backdoor attacks. The reason is that the spectral signatures that those methods build upon are challenging to detect for a wide range of attacks. We therefore introduce a novel defense algorithm, that we call SPECTRE, by combining the ideas from robust covariance estimation and quantum entropy outlier detection. Whitening with the robust covariance amplifies the spectral signature of the poisoned samples. The quantum entropy score can robustly detect that signature, adapting to the spectral profile of the poisoned examples. We demonstrate the superiority of our defense using several popular backdoor attacks, which suggest that the proposed defense is successful in all regimes we tested on, including those where the state-of-the-art baseline approaches fail. The empirical success of SPECTRE

opens several new research directions, two of which we discuss in the following.

SPECTRE requires the trainer to have access to the corrupted training dataset. In some scenarios we might not have a direct access to the training data, for example due to privacy constraints. Identifying the statistical signatures in such settings is an interesting direction to make SPECTRE more widely applicable. A concrete direction is to design a decentralized and differentially private version of SPECTRE under the setting of federated learning (Pillutla et al., 2019). Recent advances in differentially private and robust estimators in (Liu et al., 2021) provide promising directions.

(Gao et al., 2019) proposes a different paradigm for defending against backdoor attacks. The defense, called STRIP, mixes each training sample with multiple other samples and measure the entropy of the resulting prediction. This leverages an aspect of common backdoor attacks that is different from spectral signatures. Understanding how these different types of defenses perform against different types of attacks, such as the hidden backdoor attacks from (Saha et al., 2020), is an important research question.

Acknowledgement

This work is supported by Google faculty research award and NSF grants CNS-2002664, IIS-1929955, and CCF-2019844 as a part of Institute for Foundations of Machine Learning.

References

- Alistarh, D., Allen-Zhu, Z., and Li, J. Byzantine stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 4613–4623, 2018.
- Awasthi, P., Jain, H., Rawat, A. S., and Vijayaraghavan, A. Adversarial robustness via robust low rank representations. *Advances in Neural Information Processing Systems*, 33, 2020.
- Barni, M., Kallas, K., and Tondi, B. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 101–105. IEEE, 2019.
- Biggio, B., Rieck, K., Ariu, D., Wressnegger, C., Corona, I., Giacinto, G., and Roli, F. Poisoning behavioral malware clustering. In *Proceedings of the 2014 workshop on artificial intelligent and security workshop*, pp. 27–36, 2014.
- Blanchard, P., Guerraoui, R., and Stainer, J. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in Neural Information Processing Systems*, pp. 119–129, 2017.

- Chen, B., Carvalho, W., Baracaldo, N., Ludwig, H., Edwards, B., Lee, T., Molloy, I., and Srivastava, B. Detecting backdoor attacks on deep neural networks by activation clustering. *arXiv preprint arXiv:1811.03728*, 2018a.
- Chen, L., Wang, H., Charles, Z., and Papailiopoulos, D. Draco: Byzantine-resilient distributed training via redundant gradients. *arXiv preprint arXiv:1803.09877*, 2018b.
- Chen, M., Gao, C., Ren, Z., et al. Robust covariance and scatter matrix estimation under huber’s contamination model. *Annals of Statistics*, 46(5):1932–1960, 2018c.
- Chen, X., Liu, C., Li, B., Lu, K., and Song, D. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- Cheng, Y., Diakonikolas, I., Ge, R., and Woodruff, D. P. Faster algorithms for high-dimensional robust covariance estimation. In *Conference on Learning Theory*, pp. 727–757. PMLR, 2019.
- Chou, E., Tramèr, F., Pellegrino, G., and Boneh, D. Sentinet: Detecting physical attacks against deep learning systems. *arXiv preprint arXiv:1812.00292*, 2018.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. Being robust (in high dimensions) can be practical. In *International Conference on Machine Learning*, pp. 999–1008. PMLR, 2017a.
- Diakonikolas, I., Kane, D. M., and Stewart, A. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 73–84. IEEE, 2017b.
- Diakonikolas, I., Kamath, G., Kane, D., Li, J., Moitra, A., and Stewart, A. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 48(2):742–864, 2019.
- Dong, Y., Hopkins, S. B., and Li, J. Quantum entropy scoring for fast robust mean estimation and improved outlier detection. *arXiv preprint arXiv:1906.11366*, 2019.
- Gao, Y., Xu, C., Wang, D., Chen, S., Ransinghe, D. C., and Nepal, S. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th Annual Computer Security Applications Conference*, pp. 113–125, 2019.
- Gu, T., Dolan-Gavitt, B., and Garg, S. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., and Madry, A. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, pp. 125–136, 2019.
- Jambulapati, A., Li, J., and Tian, K. Robust sub-gaussian principal component analysis and width-independent Schatten packing. *arXiv preprint arXiv:2006.06980*, 2020.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Kearns, M. and Li, M. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993.
- Kolouri, S., Saha, A., Pirsiavash, H., and Hoffmann, H. Universal litmus patterns: Revealing backdoor attacks in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 301–310, 2020.
- Kong, W., Somani, R., Kakade, S., and Oh, S. Robust meta-learning for mixed linear regression with small batches. *arXiv preprint arXiv:2006.09702*, 2020.
- Lai, K. A., Rao, A. B., and Vempala, S. Agnostic estimation of mean and covariance. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 665–674. IEEE, 2016.
- Laskov, P. Practical evasion of a learning-based classifier: A case study. In *2014 IEEE symposium on security and privacy*, pp. 197–211. IEEE, 2014.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Processing Systems*, pp. 7167–7177, 2018.
- Li, J. and Ye, G. Robust gaussian covariance estimation in nearly-matrix multiplication time. *arXiv preprint arXiv:2006.13312*, 2020.
- Li, S., Xue, M., Zhao, B. Z. H., Zhu, H., and Zhang, X. Invisible backdoor attacks on deep neural networks via steganography and regularization. *arXiv preprint arXiv:1909.02742*, 2019.
- Liang, S., Li, Y., and Srikant, R. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- Liu, K., Dolan-Gavitt, B., and Garg, S. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, pp. 273–294. Springer, 2018.

- Liu, X., Kong, W., Kakade, S., and Oh, S. Robust and differentially private mean estimation. *arXiv preprint arXiv:2102.09159*, 2021.
- Liu, Y., Ma, S., Aafer, Y., Lee, W.-C., Zhai, J., Wang, W., and Zhang, X. Trojancing attack on neural networks. 2017.
- Liu, Y., Ma, X., Bailey, J., and Lu, F. Reflection backdoor: A natural backdoor attack on deep neural networks. In *European Conference on Computer Vision*, pp. 182–199. Springer, 2020.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Mozaffari-Kermani, M., Sur-Kolay, S., Raghunathan, A., and Jha, N. K. Systematic poisoning attacks on and defenses for machine learning in healthcare. *IEEE journal of biomedical and health informatics*, 19(6):1893–1905, 2014.
- Newsome, J., Karp, B., and Song, D. Paragraph: Thwarting signature learning by training maliciously. In *International Workshop on Recent Advances in Intrusion Detection*, pp. 81–105. Springer, 2006.
- Periša, L. *Recompression of Hadamard products of tensors in Tucker format*. PhD thesis, University of Zagreb. Faculty of Science. Department of Mathematics, 2017.
- Pillutla, K., Kakade, S. M., and Harchaoui, Z. Robust aggregation for federated learning. *arXiv preprint arXiv:1912.13445*, 2019.
- Quiring, E. and Rieck, K. Backdooring and poisoning neural networks with image-scaling attacks. *arXiv preprint arXiv:2003.08633*, 2020.
- Rubinstein, B. I., Nelson, B., Huang, L., Joseph, A. D., Lau, S.-h., Rao, S., Taft, N., and Tygar, J. D. Antidote: understanding and defending against poisoning of anomaly detectors. In *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement*, pp. 1–14, 2009.
- Saha, A., Subramanya, A., and Pirsiavash, H. Hidden trigger backdoor attacks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):11957–11965, Apr. 2020. doi: 10.1609/aaai.v34i07.6871. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6871>.
- Shokri, R. et al. Bypassing backdoor detection algorithms in deep learning. In *2020 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 175–183. IEEE, 2020.
- Steinhardt, J., Koh, P. W. W., and Liang, P. S. Certified defenses for data poisoning attacks. In *Advances in neural information processing systems*, pp. 3517–3529, 2017.
- Sun, Z., Kairouz, P., Suresh, A. T., and McMahan, H. B. Can you really backdoor federated learning? *arXiv preprint arXiv:1911.07963*, 2019.
- Tran, B., Li, J., and Madry, A. Spectral signatures in backdoor attacks. In *Advances in Neural Information Processing Systems*, pp. 8000–8010, 2018.
- Turner, A., Tsipras, D., and Madry, A. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771*, 2019.
- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., and Zhao, B. Y. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pp. 707–723. IEEE, 2019.
- Wang, B., Cao, X., and Gong, N. Z. On certifying robustness against backdoor attacks via randomized smoothing. *arXiv preprint arXiv:2002.11750*, 2020a.
- Wang, H., Sreenivasan, K., Rajput, S., Vishwakarma, H., Agarwal, S., Sohn, J.-y., Lee, K., and Papailiopoulos, D. Attack of the tails: Yes, you really can backdoor federated learning. *Advances in Neural Information Processing Systems*, 33, 2020b.
- Weber, M., Xu, X., Karlas, B., Zhang, C., and Li, B. Rab: Provable robustness against backdoor attacks. *arXiv preprint arXiv:2003.08904*, 2020.
- Xiao, H., Biggio, B., Brown, G., Fumera, G., Eckert, C., and Roli, F. Is feature selection secure against training data poisoning? In *International Conference on Machine Learning*, pp. 1689–1698, 2015.
- Xie, C., Huang, K., Chen, P.-Y., and Li, B. Dba: Distributed backdoor attacks against federated learning. In *International Conference on Learning Representations*, 2019.
- Yang, C., Wu, Q., Li, H., and Chen, Y. Generative poisoning attack method against neural networks. *arXiv preprint arXiv:1703.01340*, 2017.
- Zhao, S., Ma, X., Zheng, X., Bailey, J., Chen, J., and Jiang, Y.-G. Clean-label backdoor attacks on video recognition models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14443–14452, 2020.
- Zhong, H., Liao, C., Squicciarini, A. C., Zhu, S., and Miller, D. Backdoor embedding in convolutional neural network models via invisible perturbation. In *Proceedings of the Tenth ACM Conference on Data and Application Security and Privacy*, pp. 97–108, 2020.