# TRUSTWORTHY MACHINE LEARNING POISONING

KORBINIAN KOCH, RAHAVI SELVARAJAN, ANGELA WANG September 21, 2021



# **ATTACKING SPAM FILTERS**



In practice, we regularly see some of the most advanced spammer groups trying to throw the Gmail filter off-track by reporting massive amounts of spam emails as not spam. As shown in the figure, between the end of Nov 2017 and early 2018, there were at least four malicious large-scale attempts to skew our classifier.

Elie Bursztein, Google

https://elie.net/blog/ai/attacks-against-machine-learning-an-overview/



### **EXAMPLE POISONED SPAM DATAPOINT**

Hello Friend!

Your professor doesn't want you to know this!

\*\*\* FIND GLOBAL OPTIMA IN SECONDS \*\*\*

<u>>>> READ THE FREE PAPER HERE <<<</u>

Supply is limited!!1!

Use referral code: TrustworthyML

www.global-optima-fast.xx

#### Manual label: **not spam**







### **THREAT MODEL**

- Model owner is the victim
- Attacker tries to modify the dataset before model owner trains the model in order to influence prediction output of the model
- Most frequent goals
  - $\rightarrow$  affect the model's accuracy
  - $\rightarrow$  create model backdoor



#### **EXAMPLE: SHIFT OF DECISION BOUNDARY**



https://www.slideshare.net/DavidDao1/causative-adverserial-learning



### **TYPES OF POISONING**

- 1. Targeted vs. untargeted:
- want a specific email to pass the spam detection: targeted
- generally reduce model's accuracy: untargeted

#### 2. Availability vs. integrity

- data injection by adding confusing emails: availability
- backdoor for specific keyword in email: integrity

and more ...



#### **LABEL MODIFICATION**





#### **DATA INJECTION**





#### BACKDOORING



Li, Yuezun & Li, Yiming & Wu, Baoyuan & Li, Longkang & He, Ran & Lyu, Siwei. (2020). Backdoor Attack with Sample-Specific Triggers.



#### **CHOOSE YOUR POISON**





#### Skewing

shifts decision boundary in desired direction

✓ desired effect on normal performance





Attackers actively attempt to shift the learned boundary between abusive and legitimate use in their favor.

Poisoning training data by manipulating the real world.



#### **OVERVIEW**









In Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement, November 2009, pp. 1-14

#### ANTIDOTE: Understanding and Defending against Poisoning of Anomaly Detectors

Benjamin I. P. Rubinstein<sup>1</sup> Blaine Nelson<sup>1</sup> Ling Huang<sup>2</sup> Anthony D. Joseph<sup>1,2</sup> Shing-hon Lau<sup>1</sup> Satish Rao<sup>1</sup> Nina Taft<sup>2</sup> J. D. Tygar<sup>1</sup>

<sup>1</sup>Computer Science Division, University of California, Berkeley <sup>2</sup>Intel Labs Berkeley

#### ABSTRACT

Statistical machine learning techniques have recently garnered increased popularity as a means to improve network design and security. For intrusion detection, such methods build a model for normal behavior from training data and detect attacks as deviations from that model. This process invites adversaries to manipulate the training data so that the learned model fails to detect subsequent attacks.

We evaluate poisoning techniques and develop a defense, in the context of a particular anomaly detector—namely the PCA-subspace method for detecting anomalies in backbone networks. For three poisoning schemes, we show how attackers can substantially increase their chance of successfully evading detection by only adding moderate amounts of poisoned data. Moreover such poisoning throws off the balance between false positives and false negatives thereby dramatically reducing the efficacy of the detector.

To combat these poisoning activities, we propose an antidote based on techniques from robust statistics and present a new robust PCA-based detector. Poisoning has little effect on the about model where it similarithe distant the

#### Keywords

Network Traffic Analysis, Principal Components Analysis, Adversarial Learning, Robust Statistics

#### 1. INTRODUCTION

Statistical machine learning (SML) techniques are increasingly being used as tools for analyzing and improving network design and performance. They have been applied to a variety of problems such as enterprise network fault diagnosis [1, 5, 14], email spam filtering [24, 27], worm detection [25], and intrusion detection [16, 30, 33], as well as many others. These solutions draw upon a variety of techniques from the SML domain including Singular Value Decomposition, clustering, Bayesian inference, spectral analysis, maximum-margin classification, etc. In many scenarios, these approaches have been demonstrated to perform well. Many of these SML techniques include a learning phase during which a model is trained using collected data. Such techniques have a serious vulnerability, namely they are susceptible to adversaries who purposefully inject mali-







#### **ABILENE BACKBONE NETWORK**

onode ('PoP') link 



#### **ABILENE BACKBONE NETWORK**





#### **BACKBONE NETWORK TRAFFIC**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.







1. Standardize data such that **mean** and **sd** of each dimension are 0 and 1





1. Standardize data such that **mean** and **sd** of each dimension are 0 and 1

2. Find **unit vector** 





1. Standardize data such that **mean** and **sd** of each dimension are 0 and 1

2. Find **unit vector** such that it maximizes the **mean squared distance** of the projected points to the origin









#### **NETWORK TRAFFIC TRAINING DATA**





### **PCA SCREE PLOT**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.



### **PCA PROJECTION OF DATA**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.



#### **PCA PROJECTION OF DATA**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.



#### **PCA PROJECTION OF DATA**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.



# **STATE VS. RESIDUAL MAGNITUDE**



Lakhina, Anukool & Crovella, Mark & Diot, Christophe. (2004). Diagnosing Network-Wide Traffic Anomalies. Computer Communication Review. 34.



#### "PCA subspace method"

for anomaly detection







Understanding Anomaly Detectors Understanding Poisoning of Anomaly Detectors

÷,



Defending against Poisoning of Anomaly Detectors



### **PCA SUBSPACE RECALCULATION**

- 'Normal' network usage changes over time (e.g. during lockdowns)
- PCs and threshold get **recalculated every week**
- Values used in week w were learned in week w-1





# **ASSUMPTIONS**

- We are the adversary
- We know that the network owner is using the PCA subspace method and when the subspaces get re-calculated (**grey box attack**)
- Our goal: launch DoS attack in a given week  $oldsymbol{w}$
- Our objective: change PCA subspaces learned in previous week(s) such that we will not be detected in week w (**skewing attack**)
- Our means: add superfluous traffic ('chaff') into the network





#### **POISONING THE PCA SUBSPACE METHOD: KNOWLEDGE LEVELS**



#### Uninformed

#### We know nothing about the network traffic



#### Locally informed

We know how much traffic currently occurs at our ingress link



#### **Globally informed**

We know all current and future traffic volumes at all links and can insert traffic anywhere



# **KNOWLEDGE LEVELS: UNINFORMED**

- Either add or don't add traffic in each timestep (according to random variable)
- Amount of traffic is fixed




time



## **KNOWLEDGE LEVELS: GLOBALLY INFORMED**

- Finding the right flows and chaff volumes is an optimization problem
- Unfortunately, this optimization is difficult to solve
- The authors make a few assumptions<sup>™</sup> and do some math<sup>™</sup>
- the [...] assumption does not hold in practice
- The full proof is ommitted due to space constraints.
- Conclusion: even in this scenario we only want to add chaff along the links on the target OD flow



### **KNOWLEDGE LEVELS: METHODS**

- Increase the overall variance of traffic volumes
  - $\rightarrow$  Uninformed: randomly add constant amount of traffic
  - $\rightarrow$  Locally informed: add more traffic if volumes are already high
  - $\rightarrow$  Globally informed: like locally informed, but with more knowledge and more influenced links
- Chaff **must not be detected as anomaly** in order to be used for the calculation of PCs



### **POISONING THE PCA SUBSPACE METHOD: ATTACK DURATION**



### Single period attack

Chaff is only added in week w-1



#### **Boiling frogs attack**

Chaff is gradually increased from week *w*-*n* to *w*-1



## **USED DATA PER CONDITION**

### • Single period

Used data: week 20 and 21 from Abilene dataset + synthetic anomalies in week 21

### We pretend no chaff anomalies are detected in week 20

### Boiling frogs

Used data: synthesized traffic data + anomalies 'Normal' traffic does not change from week to week Chaff can be rejected as anomalous in every week Locally informed chaff insertion strategy



















### **POISONING EVASION SUCCESS (FALSE NEGATIVE) RATES**













## **MEDIAN-BASED PRINCIPAL COMPONENT ANALYSIS**

- PCA tries to capture the dispersion of the data
- Some statistics are more or less sensitive to outliers
- Instead of centering data around mean, center it around spatial median (= location estimate)
- Replace mean squared distance (= variance) by median absolute deviation (MAD) (= dispersion measure)
- PCs are found using **PCA-GRID** (uses grid search)





# PCA-GRID (Croux et al., 2007)

- Given location estimate and dispersion measure PCA-GRID finds the principal components via **grid search**
- Iteratively divide the search space up into subspaces using candidate vectors, refining the angle of the best performing vector
- This yields an approximate solution that maximizes the dispersion measure
- Project and repeat for multiple PCs



### **LAPLACE VS. Q-STATISTIC THRESHOLD**



#### **Histogram of PCA-GRID Residuals**

Rubinstein et al. (2009)



### **LAPLACE-DISTRIBUTION VS. Q-DISTRIBUTION**



https://www.vosesoftware.com/riskwiki/images/image15\_632.gif





### median-based PCA-GRID + Laplace threshold = ANTIDOTE



UNIVERSITY OF

RONTO

189 (89 19)



35% chaff globally informed single period attack



### **POISONING EVASION SUCCESS RATES: SINGLE PERIOD**



Rubinstein et al. (2009)



### **POISONING EVASION SUCCESS RATES: BOILING FROGS**



#### **ANTIDOTE**





### **AREA UNDER CURVE (AUC): PCA VS. ANTIDOTE**



Rubinstein et al. (2009)



## **SHORTCOMINGS: UNREALISTIC DATA**

- Abilene was a university research network, not public internet traffic
- Week 20 and 21 were cherry-picked
- The anomalies were just simulated
- The 'normal' traffic was simulated and completely stationary (boiling frogs)
- Authors acknowledge the approach didn't work using real non-stationary data (boiling frogs)



### **SHORTCOMINGS: UNREALISTIC KNOWLEDGE LEVEL**

- Authors claim the unrealistic 'globally informed' strategy was included to test the limitations of ANTIDOTE
- However, they report their best results unter this strategy
- If you remove the results on the 'globally informed' strategy, you are left with
  - $\rightarrow$  no effect on the uninformed strategy
  - $\rightarrow$  50% reduction for the locally informed strategy

and those only if 100% of chaff remains undetected in week w-1



### **SHORTCOMINGS: DOS VARIANCE IS FREQUENT**

- The authors assume DoS attacks are outliers
- There are around 28,700 distinct DoS attacks per day <sup>1</sup>
- If DoS attacks are ubiquitous, they will be encoded in the normal subspace no matter which method we use
- Distributed DoS attacks using botnets are even less detectable

<sup>1</sup> Jonker, Mattijs & King, Alistair & Krupp, Johannes & Rossow, Christian & Sperotto, Anna & Dainotti, Alberto. (2017). Millions of targets under attack: a macroscopic characterization of the DoS ecosystem. 100-113.



## **IMPROVEMENTS**

- Try to gather genuine and more recent network data
- Don't pretend all chaff remains undetected (single period)
- Cross-validate on multiple weeks
- Report results on locally informed scheme
- Release code



### **DISCUSSION**





## **CONCLUSION**

- Good explanation on how PCA subspace poisoning works by shifting normal PCs
- Systematic comparison of different attack strategies, some of which are quite successful (e.g. random chaff, boiling frogs)
- Provides evidence for robust median-based PCA
- But: unrealistic performance for both attack and defense due to utilized data and assumptions made



## **SHORTCOMINGS: OTHERS**

- Projection on 1st PC in graph is not centered around 0, implying that the data was not standardized
  - $\rightarrow$  mismatch between what authors claim to do and show in graph
- PCA subspace method does not consider order of data points
- Evasion success rate unrealistic: if you fail once, you get blocked and can't retry the attack 1000 times
- No code was released



### **PCA IN REAL LIFE: DETECTING CLOUD MINING**



GCE instance temporal behavioral shift due to the start of mining

https://elie.net/blog/ai/attacks-against-machine-learning-an-overview/



### **REJECTION RATES (BOILING FROGS)**



Rubinstein et al. (2009)



### **ROC CURVES (SINGLE PERIOD)**



Rubinstein et al. (2009)





N

0

N

Ma

4

-

2

27

Ó

0

### Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning

Matthew Jagielski\*, Alina Oprea\*, Battista Biggio <sup>†‡</sup>, Chang Liu<sup>§</sup>, Cristina Nita-Rotaru\*, and Bo Li<sup>§</sup>

\*Northeastern University, Boston, MA <sup>†</sup>University of Cagliari, Italy <sup>‡</sup>Pluribus One, Italy <sup>§</sup>UC Berkeley, Berkeley, CA

Abstract-As machine learning becomes widely used for automated decisions, attackers have strong incentives to manipulate the results and models generated by machine learning algorithms. In this paper, we perform the first systematic study of poisoning attacks and their countermeasures for linear regression models. In poisoning attacks, attackers deliberately influence the training data to manipulate the results of a predictive model. We propose a theoretically-grounded optimization framework specifically designed for linear regression and demonstrate its effectiveness on a range of datasets and models. We also introduce a fast statistical attack that requires limited knowledge of the training process. Finally, we design a new principled defense method that is highly resilient against all poisoning attacks. We provide formal guarantees about its convergence and an upper bound on the effect of poisoning attacks when the defense is deployed. We evaluate extensively our attacks and defenses on three realistic datasets from health care, loan assessment, and CS real estate domains.1 2

#### I. INTRODUCTION

As more applications with large societal impact rely on machine learning for automated decisions, several concerns have emerged about potential vulnerabilities introduced by matraining process. Such poisoning attacks have been practically demonstrated in worm signature generation [42], [45], spam filters [40], DoS attack detection [47], PDF malware classification [55], handwritten digit recognition [5], and sentiment analysis [41]. We argue that these attacks become easier to mount today as many machine learning models need to be updated regularly to account for continuously-generated data. Such scenarios require online training, in which machine learning models are updated based on new incoming training data. For instance, in cyber-security analytics, new Indicators of Compromise (IoC) rise due to the natural evolution of malicious threats, resulting in updates to machine learning models for threat detection [23]. These IoCs are collected from online platforms like VirusTotal, in which attackers can also submit IoCs of their choice. In personalized medicine, it is envisioned that patient treatment is adjusted in realtime by analyzing information crowdsourced from multiple participants [16]. By controlling a few devices, attackers can submit fake information (e.g., sensor measurements), which is then used for training models applied to a large set of pa-



### **MOTIVATION**

- Existence of strong incentives to manipulate the automated decisions by machine learning model.
- No robust algorithm for detecting poisoning points in regression models



## **CONTRIBUTIONS**

- Theoretically grounded optimization framework tuned for regression models.
- Designed a statistical attack method
- Robust defense algorithm to tackle inliers.
- First to consider poisoning attacks against regression models



## **SUPERVISED MACHINE LEARNING**





### **LINEAR REGRESSION**





### **ATTACK METHODOLOGY**



White box attacks:

ć

$$egin{argmax}{l} rgmax_{\mathcal{D}_p} & \mathcal{W}ig(\mathcal{D}',oldsymbol{ heta}_p^{\star}ig) \ ext{ s.t. } & oldsymbol{ heta}_p^{\star}\inrgmin_{ heta}\mathcal{L}(\mathcal{D}_{ ext{tr}}\cup\mathcal{D}_p,oldsymbol{ heta}ig) \end{array}$$

D' - untainted dataset  $\theta_p^*$  - poisoned regression parameters  $D_{tr}^*$  - training dataset  $D_{tr}^*$  - substitute dataset

#### **Black box attacks:**

- The attacker do not know the knowledge of the dataset.
- Replace the original dataset with substitute dataset


### **BASELINE GRADIENT DESCENT (BGD) ALGORITHM**

- Attacker's goal: maximize the regularized loss function
- Uses gradient ascent to optimize the poisoning points
- Challenges:
  - Inner learning problem neural network is not convex and hence requires efficient numerical approximations



#### **OPTIMIZATION BASED POISONING ATTACK:**

- Start with a set of initial points.
- Gradient ascent the direction which improves attacker's objective.
- Iteratively update each point with gradient ascent.
- Stop at convergence and output poisoning points.





#### **ATTACK METHODOLOGY**





#### **STATISTICAL BASED POISONING ATTACK**

- Points drawn from similar distribution as the training data
- Need to know the mean and covariance of the training data
- Black box access to the model

$$ig(X^TXig)^{-1}ig(X^TYig) o ig(X^TXig)^{-1}ig((1-lpha)X^TY+lpha X_p^TY_pig)$$

Keeps covariance similar

**Modifies the correlations** 



#### **EXISTING DEFENCES:**

EXISTING DEFENSE PROPOSALS	PROS	CONS
HUBER	<ul> <li>Noise resilient regression</li> <li>Identifies and removes outliers</li> </ul>	<ul> <li>Reduces the growth of loss function during large errors.</li> <li>Incorporates poisoning points</li> </ul>
RANSAC		<ul> <li>Samples and checks whether enough points fit the model well, random samples might contain poisoning points</li> </ul>
Chen et al.,	<ul> <li>Robust algorithm for poisoning</li> </ul>	<ul> <li>Made unrealistic assumptions like sub-Gaussian data and noise.</li> </ul>
RONI	<ul> <li>Suitable for spam scenario</li> </ul>	<ul> <li>Identifies only outliers with high impact</li> </ul>



#### **TRIM ALGORITHM**

- Given a set of data points (n) and ∝n poisoning points, TRIM tries to find the best n data points out of (1 + ∞)n points
- The regression parameters  $\theta$  = (w, b) are unknown and no assumptions are made on the distribution of the data points.
- TRIM estimates a model and identifies the points with the lowest residual from the training set.



#### TRIM ALGORITHM





#### **TRIM ALGORITHM**



Jagielski et al., (2021)



## **EXPERIMENTAL EVALUATION**

The experiments are conducted on three datasets for

- 1. Ridge regression
- 2. LASSO

The datasets are:

- Health care dataset
- Loan dataset
- House Price dataset



#### **COMPARISON OF ATTACK METHODS**

#### RIDGE REGRESSION:



LASSO:



#### **COMPARISON BETWEEN TRIM AND EXISTING DEFENCES**





### LIMITATIONS

- Conducted experiments only on two linear regression models:
  - o LASSO
  - o Ridge regression

while considered four models in the paper (OLS, LASSO, Ridge regression, Elastic-net)

- The proposed attack strategy (StatP) didn't outperform the BGD (baseline attack strategy) on few datasets.
- Randomly samples a subset of data points from a large dataset without analyzing the dataset
- Considered a small dataset for experiments and concluded that TRIM converges in a finite number of iterations
- Released code but unstructured and have unsolved issues



#### **IMPROVEMENTS**

- Use of larger datasets for the experiments might help
- Other linear regression models can also be incorporated:
  - OLS
  - Elastic-net
- Could possibly extend the experiments to polynomial regression as well.



#### **DISCUSSION**







.02815v2 [cs.LG] 29 May 2019

#### SEVER: A Robust Meta-Algorithm for Stochastic Optimization

Ilias Diakonikolas \* CS, USC diakonik@usc.edu

Jerry Li<sup>§</sup> Microsoft Research AI jerryzli@mit.edu Gautam Kamath<sup>†</sup> CS, University of Waterloo g@csail.mit.edu Jacob Steinhardt<sup>¶</sup> Statistics, UC Berkeley jsteinha@stanford.edu Daniel M. Kane<sup>‡</sup>

CSE & Math, UCSD

dakane@cs.ucsd.edu

Alistair Stewart

Web3 Foundation

stewart.al@gmail.com

May 31, 2019

#### Abstract

In high dimensions, most machine learning methods are brittle to even a small fraction of structured outliers. To address this, we introduce a new meta-algorithm that can take in a base learner such as least squares or stochastic gradient descent, and harden the learner to be resistant to outliers. Our method, SEVER, possesses strong theoretical guarantees yet is also highly scalable—beyond running the base learner itself, it only requires computing the top singular vector of a certain  $n \times d$  matrix. We apply SEVER on a drug design dataset and a spam classification dataset, and find that in both cases it has substantially greater robustness than several baselines. On the spam dataset, with 1% corruptions, we achieved 7.4% test error, compared to 13.4% – 20.5% for the baselines, and 3% error on the uncorrupted dataset. Similarly, on the drug design dataset, with 10% corruptions, we achieved 1.42 mean-squared error test error, compared to 1.51 – 2.33 for the baselines, and 1.23 error on the uncorrupted dataset.

#### 1 Introduction

Learning in the presence of outliers is a ubiquitous challenge in machine learning: nevertheless, most machine







# **Paper presentation: Sever: A Robust Meta-Algorithm for Stochastic Optimization**

**Jiaqi Wang** 



#### **MOTIVATION**

- 1. Mislabelling and Measurement Errors can cause Systematic outliers
- 2. Outliers can be introduced by poisoning attack



## **PROBLEMS WITH STATE-OF-THE-ART**

- 1. Fails when data are high-dimensional
- 2. Only works for obvious outliers, doesn't work on the correlated outliers
- 3. Specific to some algorithms, lack generalization
- 4. Significant loss in performance



#### NOVELTY

- 1. Works in high-dimensional data space
- 2. Robust to arbitrary outliers with small decrease in performance
- 3. Applicable to most ML models, including regression and classification tasks, and non-convex models like neural networks



#### PIPELINE





#### PIPELINE



Why this meta-algorithm works on trained model?



#### PIPELINE



Why this meta-algorithm works on trained model?

- Because an iteration of training (commonly gradient descent) is much cheaper than a run of SEVER



#### **INTUITIONS**

An outlier's gradient should be:

- 1. Large in magnitude
- 2. Systematically pointing in a specific direction



## **RECAP: SINGULAR VALUE DECOMPOSITION** $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\wedge \top}$

The columns of V are called the right singular vectors.

They are also the principle component of A.

The top right singular vector is the direction of A that preserves the most variance.





Algorithm 1 SEVER $(f_{1:n}, \mathcal{L}, \sigma)$ 



12: Return w.

The principle axis of the outlier has larger variance (pointing to a specific direction) and magnitude than the in-distribution data



#### **HOW TO FILTER THE POINTS**

- The learned function cannot be changed much by changing on an input while staying in domain.
- An in-distribution input could only move the function at any direction with an upper bound.
- The problem decreases to only keep the points with small singular value.
- In practise, just remove a few points with the highest scores for some iterations.



# How to Determine the Fraction Size and Number of Iterations?



# How to Determine the Fraction Size and Number of Iterations?

- They are hyper-parameters which are determined by tuning.
- The threshold is provably existent. (See Appendix for proof)



#### **EXPERIMENT RESULTS**

Experiments are conducted on two tasks:

- Ridge Regression
  - Synthetic Gaussian Dataset: 500 dimensions
  - Drug Discovery Dataset: 410 dimensions
- Support Vector Machine
  - Synthetic Gaussian Dataset: 500 dimensions
  - Spam classification task: 5116 dimensions



## **RIDGE REGRESSION**





SVM







\_

#### **WHY BASELINES FAIL**



- Baselines remove in-distribution data
- SEVER's score for outlier is clearly within the tail



#### IIMITATION

- The paper only conducted experiments on convex optimization problems. 1.
  - SEVER should show more results on modern ML models such as deep neural network (DNN).
- The open-source code is implemented in matlab, which is not the most popular language 2. in ML.
  - Should incorporate with more frequently used language like Python and R
- 3.
  - Outlier design is limited. The experiments use  $X_{\text{bad}} = \frac{1}{\alpha \cdot n_{\text{bad}}} y^{\top} X$  and  $y_{\text{bad}} = -\beta$  to generate outlier, but the state-of-art is more advanced than that.
    - Gradient-based poisoning attack might break the defense. -
- Retraining is expensive, especially for the modern ML frameworks like DNN. 4.
  - Should be integrate in the training process instead of the final model.



### (POSSIBLE) FUTURE WORK

- 1. The attack against this meta algorithm.
  - Design an adversarial outlier that is not large in certain direction.
  - Evade the detection!

- 2. The attack based on this meta algorithm
  - Design detectable outliers.
  - Take advantage of the nature of retraining, let the loop never stop.
  - Sponge the computational power!



#### **DISCUSSION**





# PAPER 1-3 Summary discussion

