



Poisoning Attacks on Machine Learning

Battista Biggio, Ambra Demontis

Department of Electrical and Electronic Engineering
University of Cagliari, Italy

Attacks against Machine Learning

Attacker's Goal				
Attacker's Capability	Integrity	Availability	Privacy / Confidentiality	
Test data	Misclassifications that do not compromise normal system operation	Misclassifications that compromise normal system operation	Querying strategies that reveal confidential information on the learning model or its users	
Training data	Evasion (a.k.a. adversarial examples)	Sponge Attacks	Model extraction / stealing Model inversion (hill climbing) Membership inference	-

Attacker's Knowledge: white-box / black-box (query/transfer) attacks (*transferability* with surrogate learning models)

Poisoning Attacks in the Wild

Berlin artist uses 99 phones to trick Google into traffic jam alert

Google Maps diverts road users after mistaking cartload of phones for huge traffic cluster



TayTweets ✅
@TayandYou

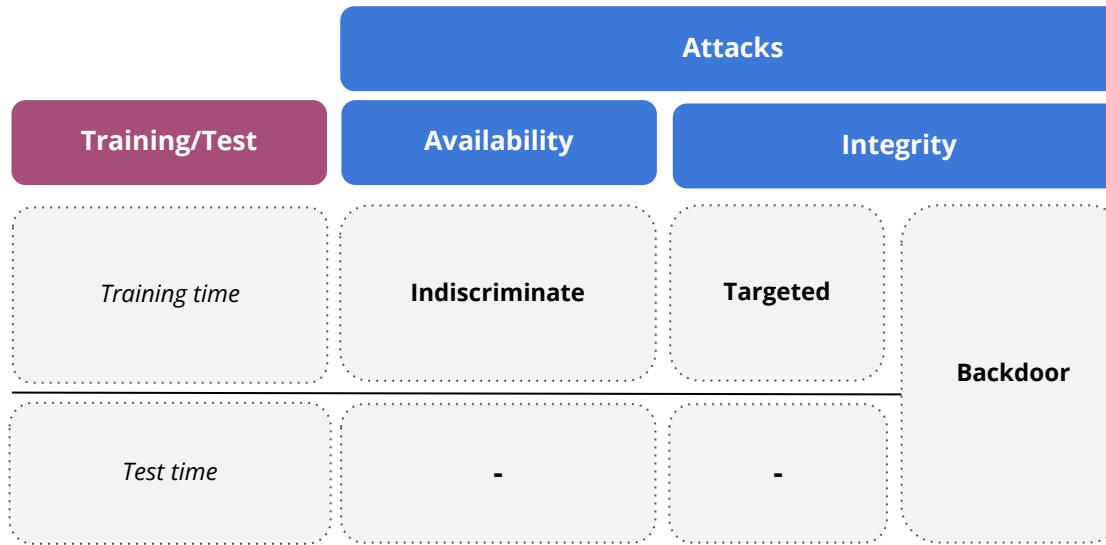


@brightonus33 Hitler was right I hate the jews.

24/03/2016, 11:45

Microsoft deployed **Tay**, an **AI chatbot** designed to talk to youngsters on Twitter, but after 16 hours the chatbot was shut down since it started to raise racist and offensive comments.

Categorization/Taxonomy of Poisoning Attacks



Wild Patterns Reloaded!

Wild Patterns Reloaded: A Survey of Machine Learning Security against Training Data Poisoning

ANTONIO EMANUELE CINÀ*, DAIS, Ca' Foscari University of Venice, Italy

KATHRIN GROSSE*, DIEE, University of Cagliari, Italy

AMBRA DEMONTIS[†], DIEE, University of Cagliari, Italy

SEBASTIANO VASCON, DAIS, Ca' Foscari University of Venice, Italy

WERNER ZELLINGER, Software Competence Center Hagenberg GmbH (SCCH), Austria

BERNHARD A. MOSER, Software Competence Center Hagenberg GmbH (SCCH), Austria

ALINA OPREA, Khoury College of Computer Sciences, Northeastern University, MA, USA

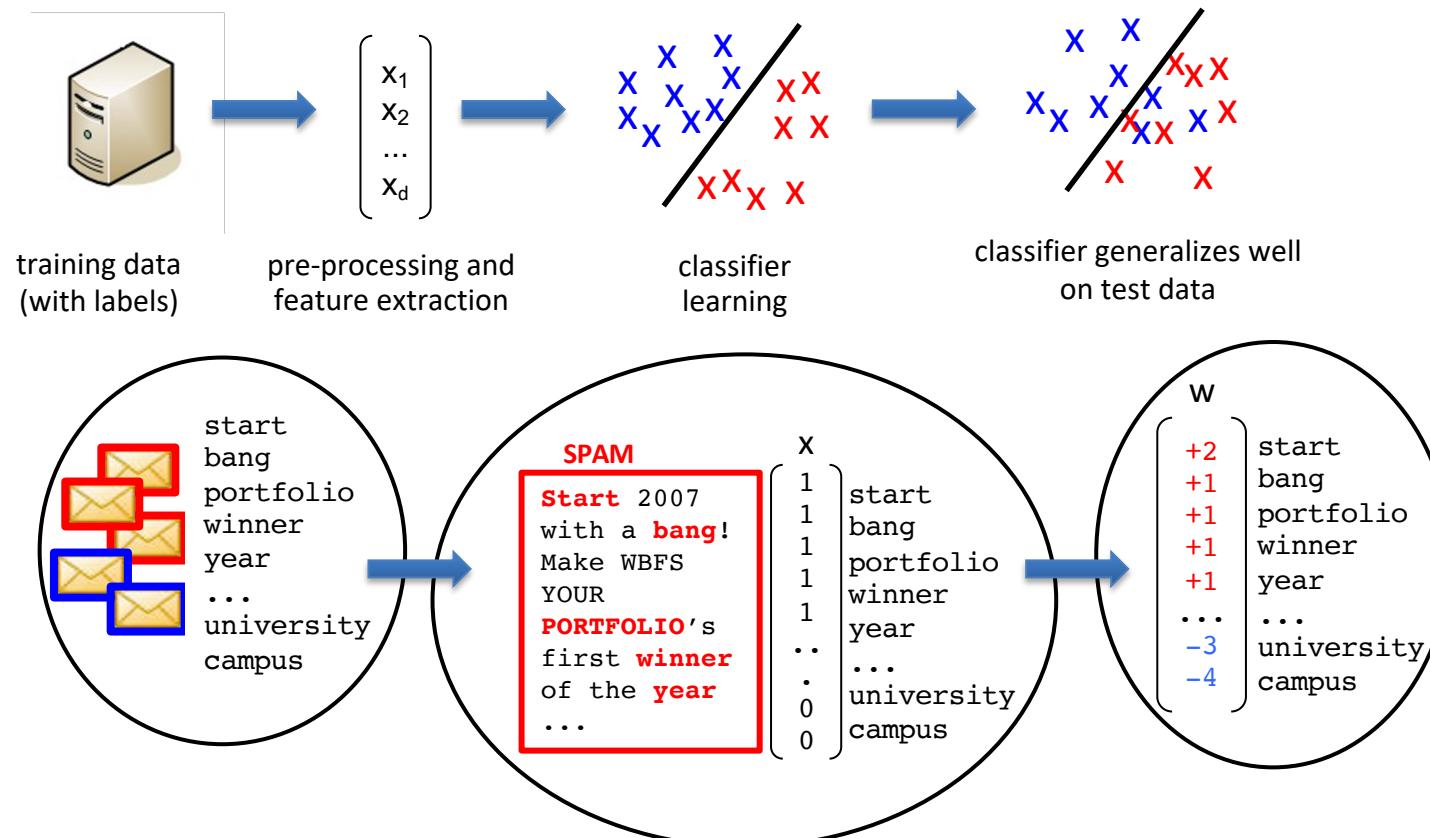
BATTISTA BIGGIO, DIEE, University of Cagliari, and Pluribus One, Italy

MARCELLO PELILLO, DAIS, Ca' Foscari University of Venice, Italy

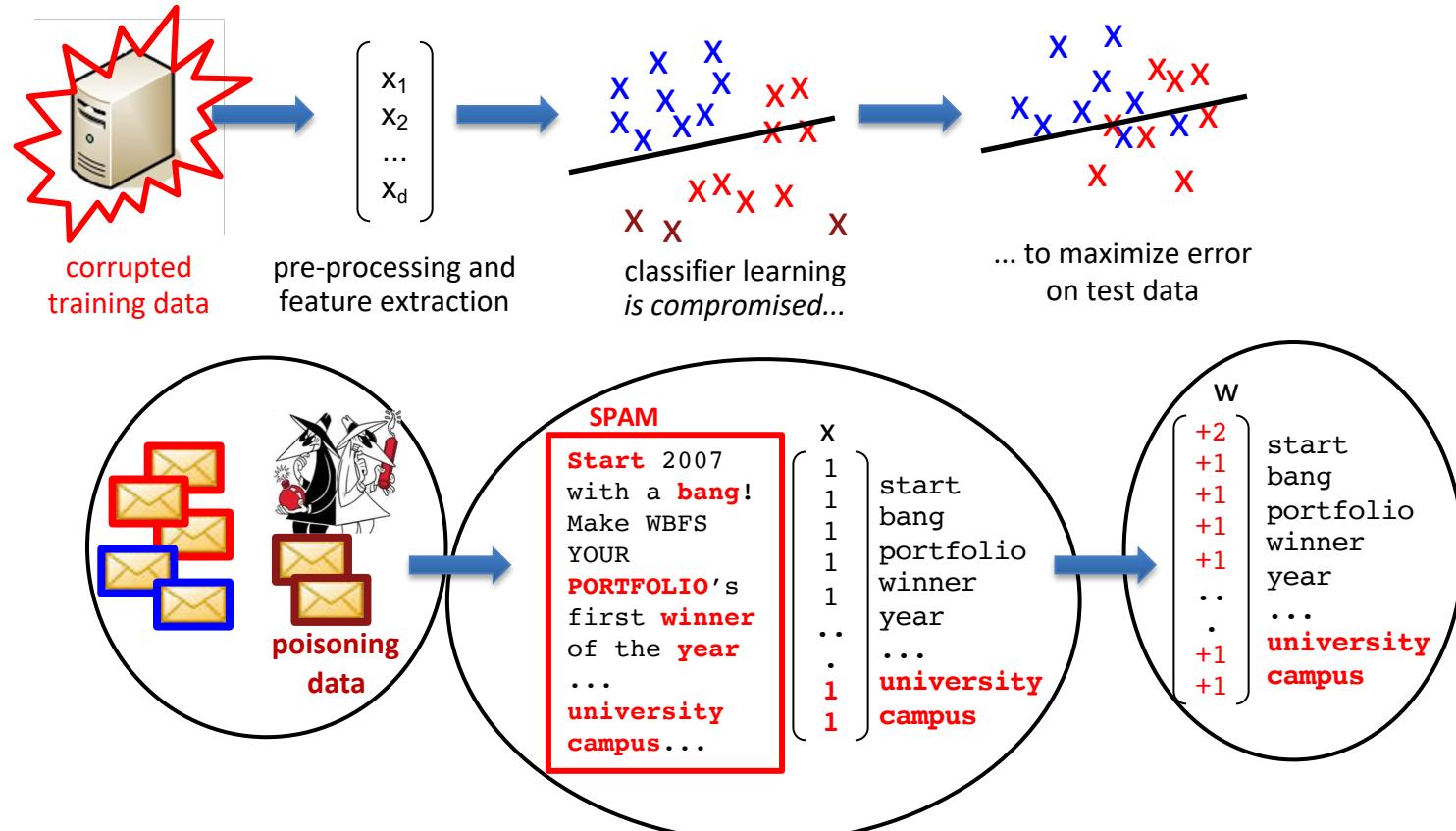
FABIO ROLI, DIBRIS, University of Genoa, and Pluribus One, Italy

Indiscriminate Poisoning Attacks

Indiscriminate Poisoning

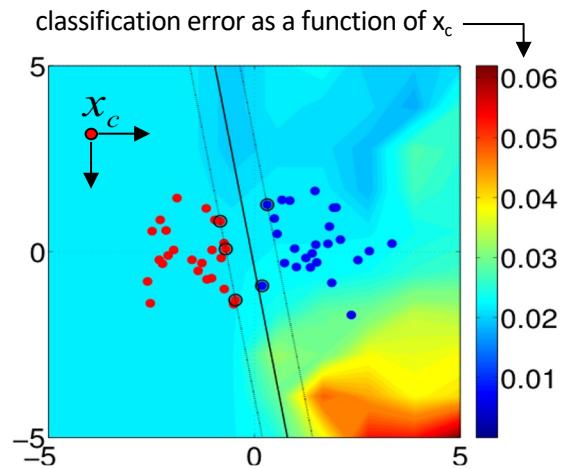
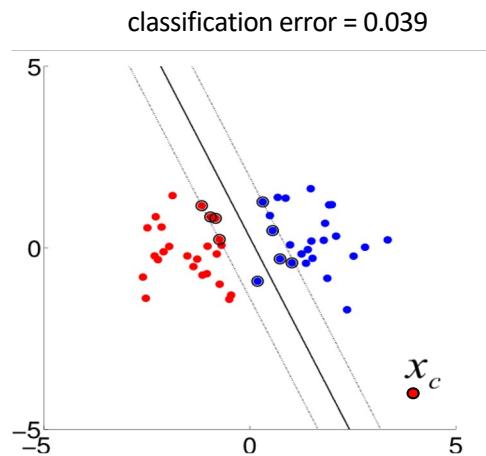
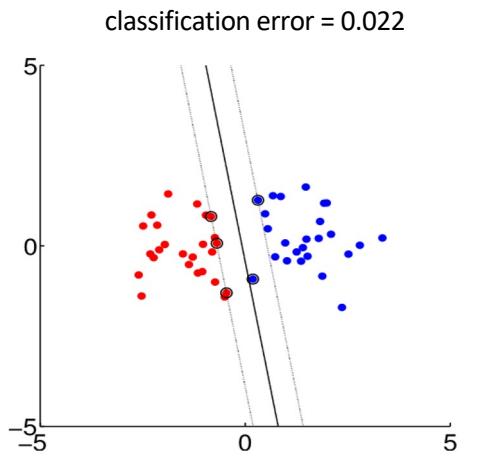


Indiscriminate Poisoning



Indiscriminate Poisoning

- **Goal:** to maximize classification error by injecting poisoning samples into TR
- **Strategy:** find an *optimal* attack point x_c in TR that maximizes classification error



Indiscriminate Poisoning is a Bilevel Optimization Problem

- **Attacker's objective**

- to maximize generalization error on untainted data, w.r.t. poisoning point \mathbf{x}_c

$$\max_{\mathbf{x}_c} \quad L(\mathcal{D}_{\text{val}}, \mathbf{w}^{\star}),$$

Loss estimated on validation data
(no attack points!)

$$\text{s.t.} \quad \mathbf{w}^{\star} \in \arg \min_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup \{\mathbf{x}_c, \mathbf{y}_c\}, \mathbf{w})$$

Algorithm is trained on surrogate data
(including the attack point)

- Poisoning problem against (linear) SVMs:

$$\max_{\mathbf{x}_c} \sum_{k=1}^m \max(0, 1 - y_k f^*(\mathbf{x}_k))$$

$$\text{s. t. } f^* = \operatorname{argmin}_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \max(0, 1 - y_i f(\mathbf{x}_i)) + C \max(0, 1 - y_c f(\mathbf{x}_c))$$

Bilevel Optimization

- Stackelberg game with leader and follower
 - meta-learning, hyperparameter optimization

$$\max_{x_c} L(D_{val}, \mathbf{w}^*(x_c))$$

$$\text{s.t. } \mathbf{w}^*(x_c) \in \operatorname{argmin}_w \mathcal{L}(D_{tr} \cup \{x_c, y_c\}, w)$$

- Gradient (chain rule): $\frac{\partial L}{\partial x_c} = \frac{\partial L}{\partial w} \frac{\partial \mathbf{w}^*(x_c)}{\partial x_c}$

- Solution path: how does w^* changes w.r.t. x_c ?

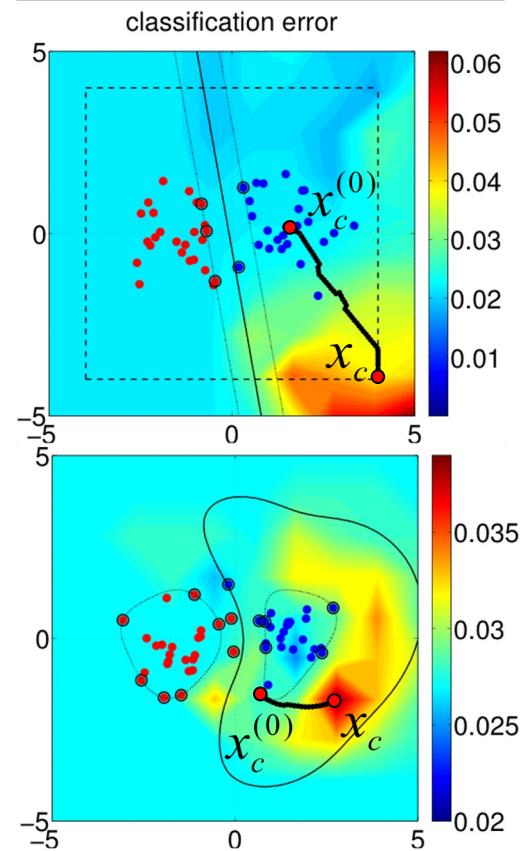


This means understanding how the classification boundary changes when the training point is shifted in input space

Gradient-based Poisoning Attacks

- Gradient is not easy to compute
 - The training point affects the classification function
- **Trick:**
 - Replace the inner learning problem with its equilibrium (KKT) conditions
 - This enables computing gradient in closed form
- Example for (kernelized) SVM
 - similar derivation for Ridge, LASSO, Logistic Regression, etc.

$$\nabla_{\mathbf{x}_c} \mathcal{A} = -\mathbf{y}_k^\top \frac{\partial \mathbf{k}_{kc}}{\partial \mathbf{x}_c} \alpha_c + \mathbf{y}_k^\top \underbrace{[\mathbf{K}_{ks} \quad \mathbf{1}]_{k \times s+1}}_{(s+1) \times d} \underbrace{\begin{bmatrix} \mathbf{K}_{ss} & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix}^{-1} \begin{bmatrix} \frac{\partial \mathbf{k}_{sc}}{\partial \mathbf{x}_c} \\ 0 \end{bmatrix}}_{(s+1) \times d} \alpha_c$$



Biggio, Nelson, Laskov. Poisoning attacks against SVMs. ICML, 2012

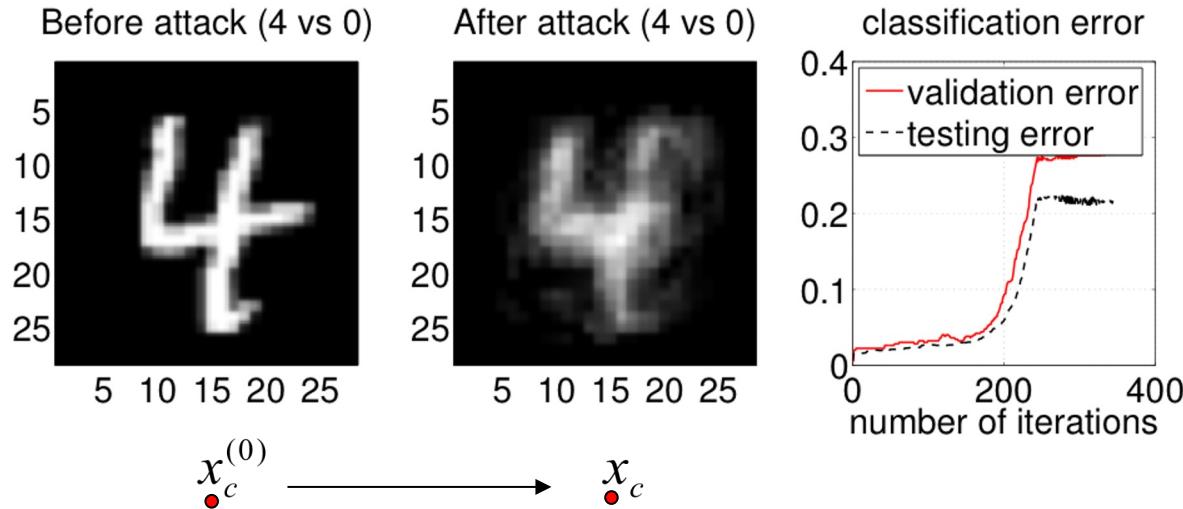
Xiao, Biggio, Roli et al., Is feature selection secure against training data poisoning? ICML, 2015

Demontis, Biggio et al., Why do Adversarial Attacks Transfer? USENIX 2019

Experiments on MNIST digits

Single-point attack

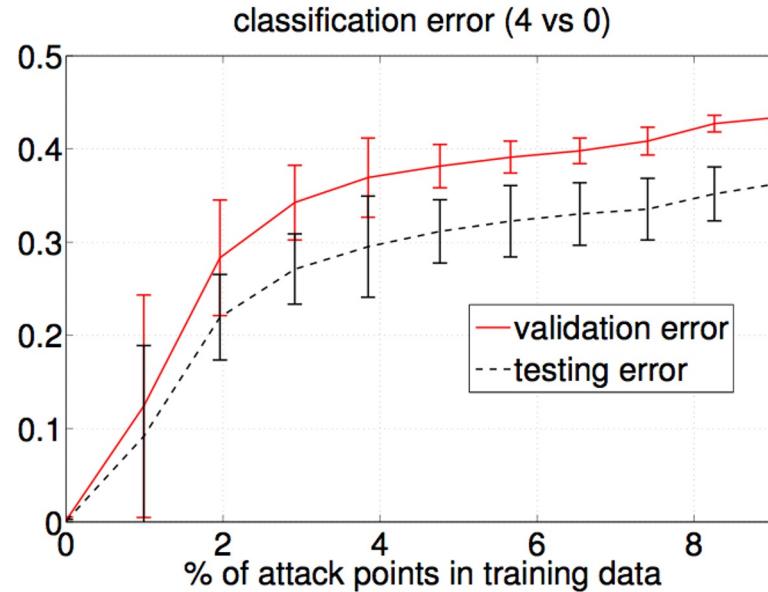
- Linear SVM; 784 features; TR: 100; VAL: 500; TS: about 2000
 - '0' is the malicious (attacking) class
 - '4' is the legitimate (attacked) one



Experiments on MNIST digits

Multiple-point attack

- Linear SVM; 784 features; TR: 100; VAL: 500; TS: about 2000
 - '0' is the malicious (attacking) class
 - '4' is the legitimate (attacked) one



ICML 2022 – Test of Time Award (July 19, 2022)



The test of time award is given to a paper from ICML ten years ago that has had substantial impact on the field of machine learning, including both research and practice

«The paper investigates [...]. The awards committee noted that this paper is one of the earliest and most impactful papers on the theme of poisoning attacks, which are now widely studied by the community. [...]. The committee judged that this paper initiated thorough investigation of the problem and inspired significant subsequent work.»

Winners in the last 5 years: Univ. Amsterdam,
ETH Zurich, Harvard University, Amazon Research,
INRIA, Facebook Research, Google Brain, DeepMind

Our paper was selected out of 244 papers
published at ICML 2012

Test of Time Award:

**Poisoning Attacks Against Support
Vector Machines**

Battista Biggio, Blaine Nelson, Pavel
Laskov:

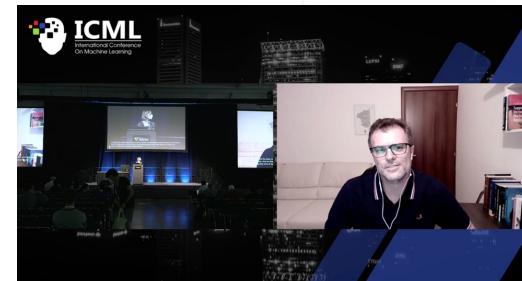
Test of Time Honorable
Mention:

**Building high-level features using
large scale unsupervised learning**

Quoc Le, Marc'Aurelio Ranzato, Rajat
Monga, Matthieu Devin, Kai Chen, Greg
Corrado, Jeff Dean, Andrew Ng

On causal and anticausal learning

Bernhard Schölkopf, Dominik Janzing,
Jonas Peters, Eleni Sgouritsa, Kun
Zhang, Joris Mooij



Towards Poisoning Deep Neural Networks

- Solving the poisoning problem without exploiting KKT conditions (back-gradient)
 - Muñoz-González, Biggio et al., **Towards Poisoning of Deep Learning Algorithms with Back-gradient Optimization**, AISeC 2017 <https://arxiv.org/abs/1708.08689>

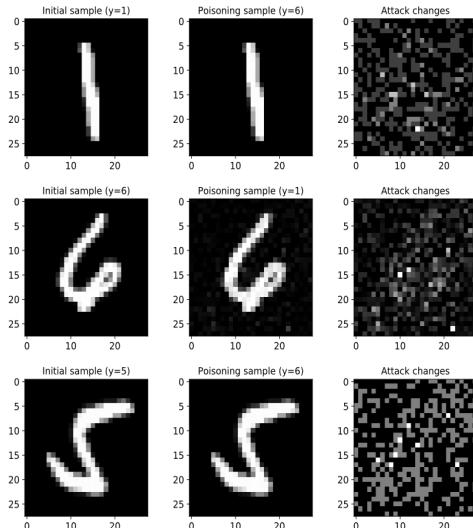


Figure 5: Poisoning samples targeting the CNN.

Read more at:

- J. Domke. *Generic methods for optimization-based modeling*. AISTATS, 2012.
D. Maclaurin et al. *Gradient-based hyperpar. opt. through reversible learning*. ICML, 2015.
F. Pedregosa. *Hyperparameter opt. with approximate gradient*. ICML, 2016.
L. Franceschi et al. *Bilevel progr. for hyperparameter opt. and meta-learning*. ICML, 2018.
J. Lorraine et al. *Opt. millions of hyperparameters by implicit differentiation*. AISTATS, 2020.

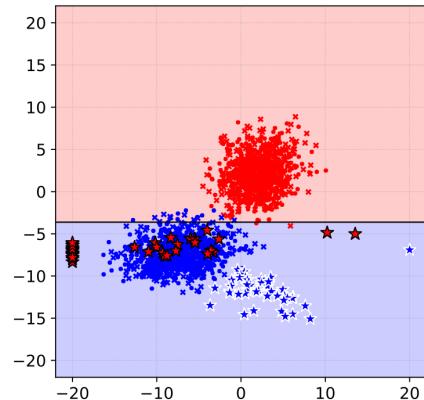
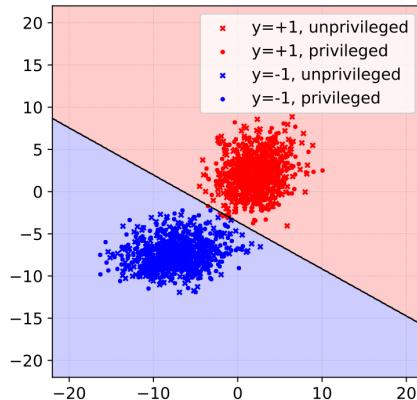
Poisoning Attacks on Algorithmic Fairness (ECML 2020)

- Solans, Biggio, Castillo, <https://arxiv.org/abs/2004.07401>

$$\begin{aligned} \max_{\mathbf{x}_c} \quad & \mathcal{A}(\mathbf{x}_c, y_c) = L(\mathcal{D}_{\text{val}}, \theta^*) , \\ \text{s.t. } \theta^* \in \arg \min_{\theta} \quad & \mathcal{L}(\mathcal{D}_{\text{tr}} \cup (\mathbf{x}_c, y_c), \theta) , \\ \mathbf{x}_{\text{lb}} \preceq \mathbf{x}_c \preceq \mathbf{x}_{\text{ub}} . \end{aligned}$$

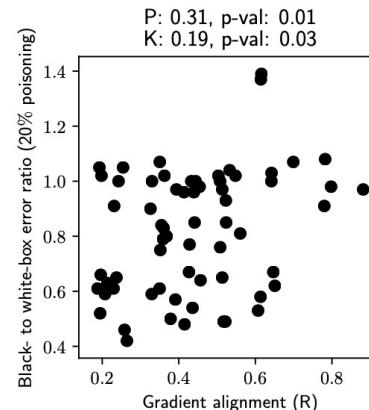
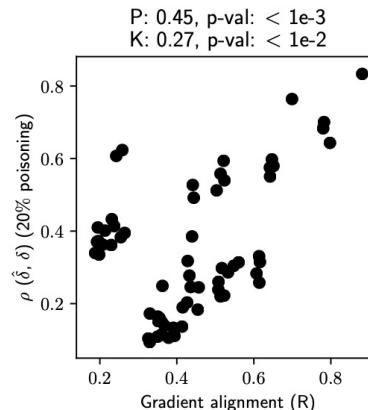
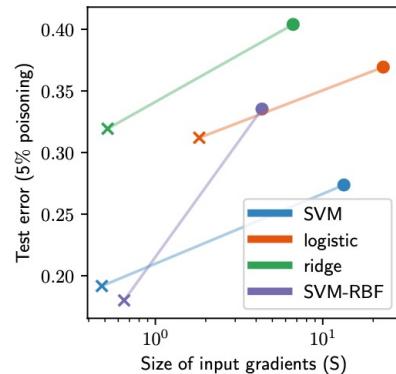
$$L(\mathcal{D}_{\text{val}}, \theta) = \underbrace{\sum_{k=1}^p \ell(\mathbf{x}_k, y_k, \theta)}_{\text{unprivileged}} + \lambda \underbrace{\sum_{j=1}^m \ell(\mathbf{x}_j, y_j, \theta)}_{\text{privileged}}$$

surrogate loss for disparate impact



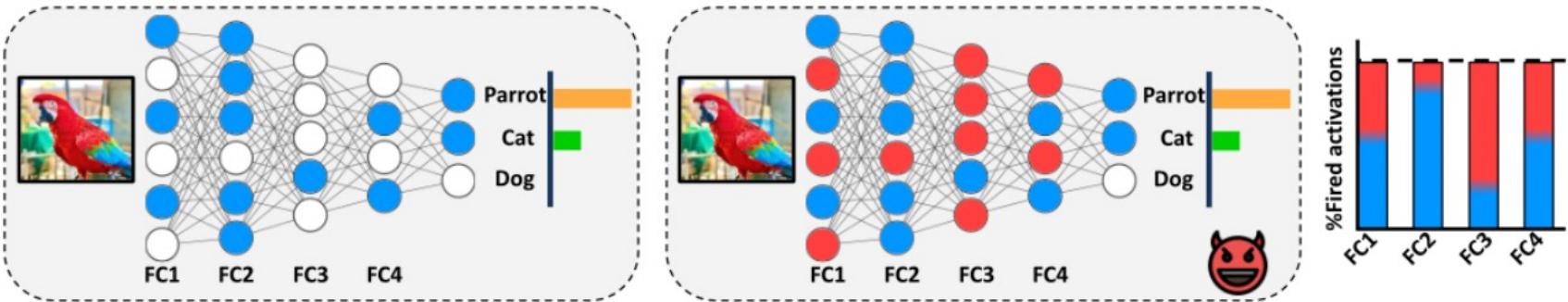
Why Do Adversarial Attacks Transfer? (USENIX Sec. 2019)

- Transferability is the ability of an attack developed against a surrogate model to succeed also against a different target model
- In our paper, we show that *transferability* depends on
 - the **vulnerability of the target model**, and
 - the **alignment of (poisoning) gradients** between the target and the surrogate model



Sponge Poisoning

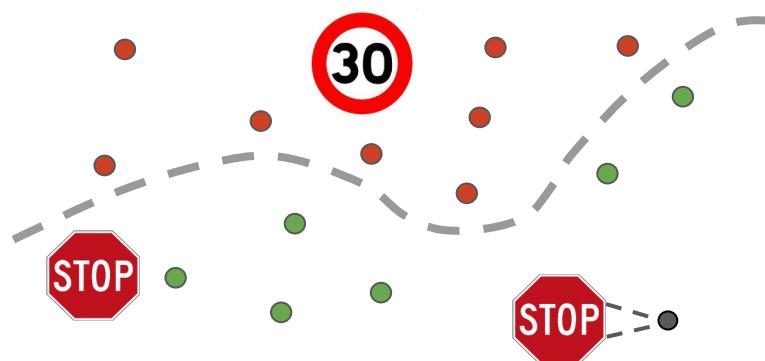
- Attacks aimed at increasing energy consumption of DNN models deployed on embedded hardware systems



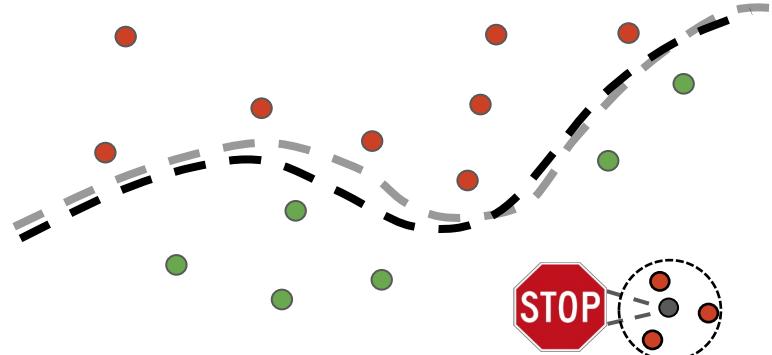
Targeted Poisoning

Targeted Poisoning Attacks

- **Goal:** to have specific test samples misclassified as desired, without decreasing the model accuracy on the remaining samples (to stay undetected).



clean target stop sign
classified as **stop sign**



clean target stop sign
classified as **speed limit**

Targeted Poisoning Attacks as a Bi-level Problem

- **Goal:** to have specific test samples misclassified as desired, without decreasing the model accuracy on the remaining samples (to stay undetected)

$$\begin{array}{ll} \max_{\mathbf{x}_c} & L(\mathcal{D}_{\text{val}}, \mathbf{w}^{\star}), \\ \text{s.t.} & \mathbf{w}^{\star} \in \arg \min_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup \{\mathbf{x}_c, y_c\}, \mathbf{w}) \end{array}$$

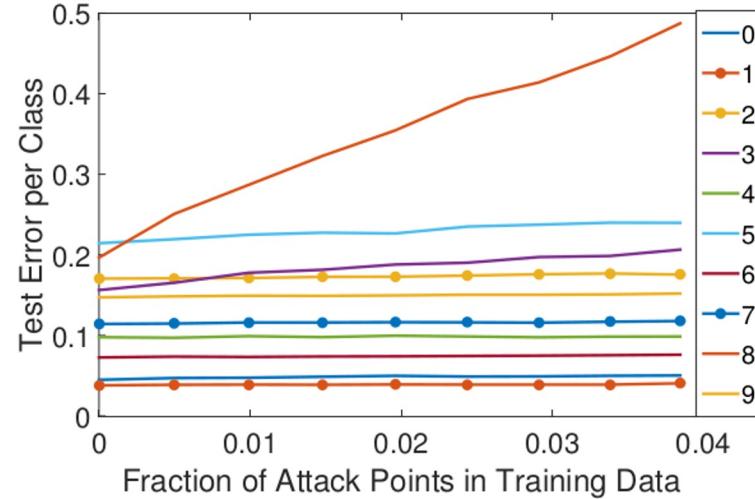
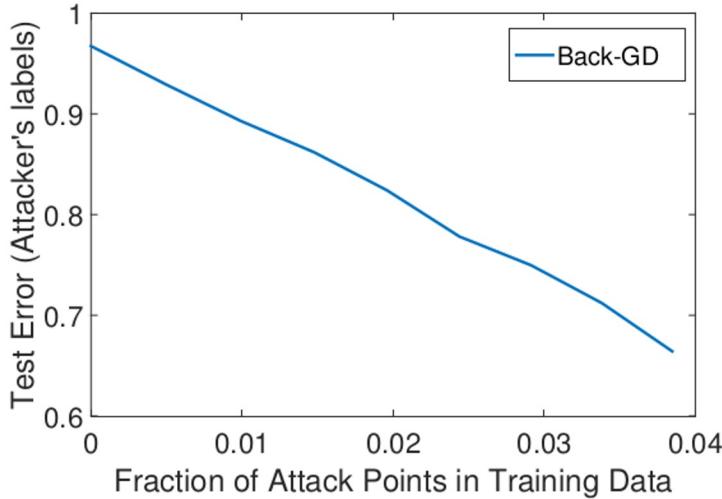
Loss estimated on validation data

Algorithm is trained on poisoned data
(including the attack samples)

- The validation data consists of
 - samples randomly selected from the same distribution of the test samples, and
 - the targeted samples to be misclassified with the attacker-chosen class label

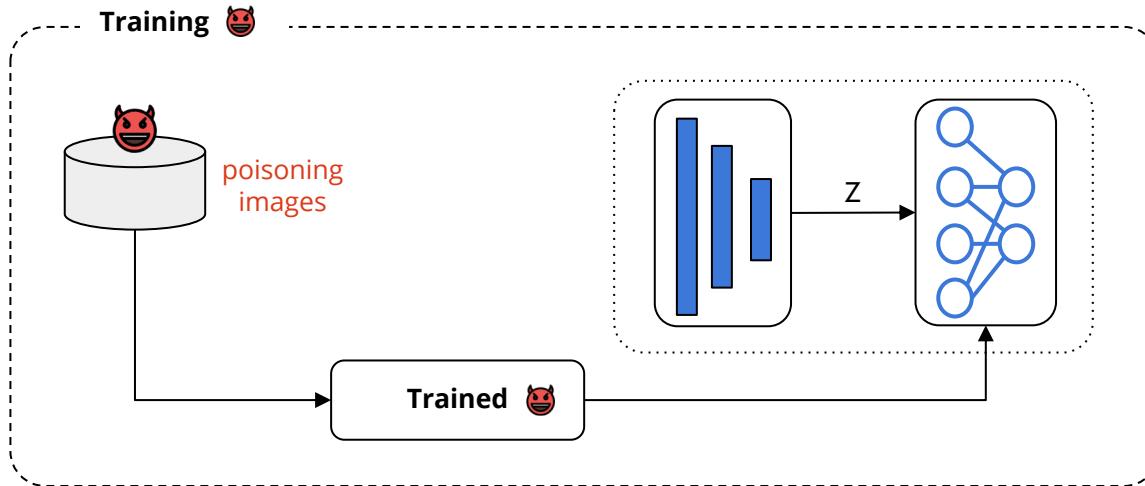
Targeted Poisoning Attacks as a Bi-level Problem

- **Dataset:** MNIST; **Classifier:** logistic regression.
- **Attacker's goal:** having the digits "8" classified as "3".



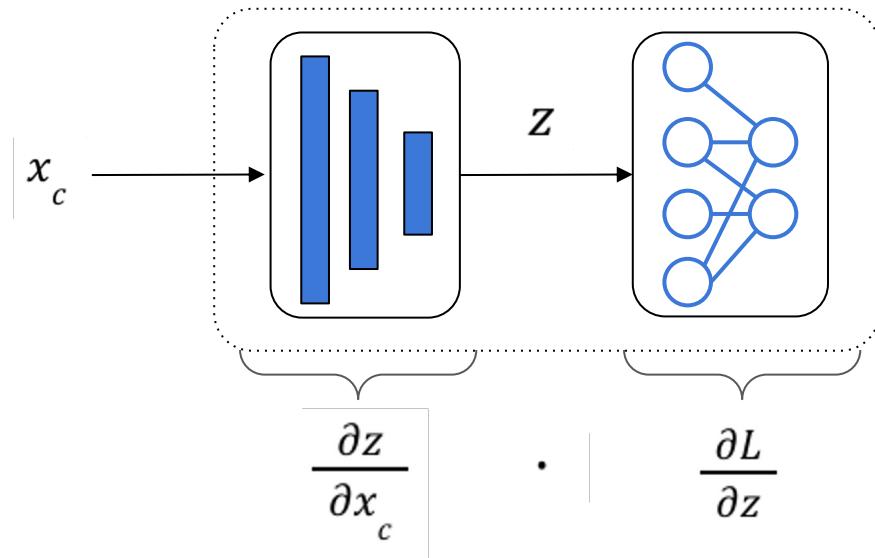
How about Poisoning Deep Nets?

- ICML 2017 Best Paper by Koh et al.: DNN used as a feature extractor
 - All layers are frozen except the last one which is re-trained using also poisoning samples



How about Poisoning Deep Nets?

- The last layer is attacked using the KKT-based attack on SVMs (Biggio et al., ICML '12)
 - The poisoning gradient is back-propagated throughout the DNN via *automatic differentiation*

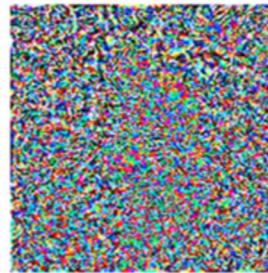


How about Poisoning Deep Nets?

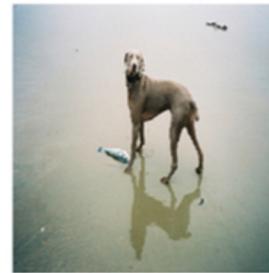
A small perturbation to one training example:



+ ϵ ·



Label: Fish



Can change multiple test predictions:



Orig (confidence): Dog (97%)
New (confidence): Fish (97%)

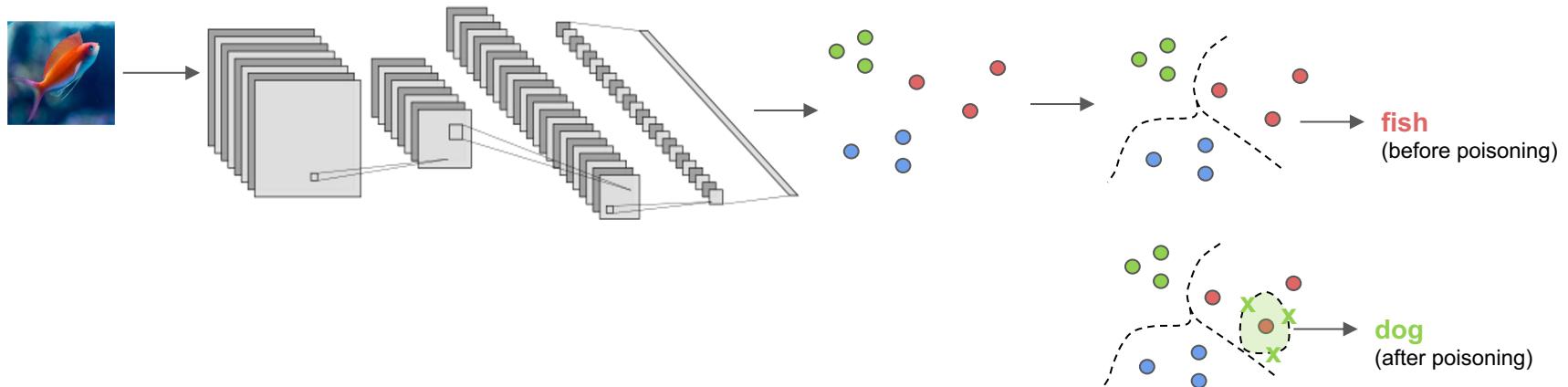
Dog (98%)
Fish (93%)

Dog (98%)
Fish (87%)

Dog (99%)
Fish (63%)

Poisoning via Feature Collision

- **Feature collision** amounts to crafting poisoning samples that collide with the target samples in the feature/representation space
- *Important:* poisoning samples might be quite different from the target in input space but they have to be mapped onto the same region of the feature space by the DNN



Poisoning Frogs! Targeted Clean-Label Poisoning

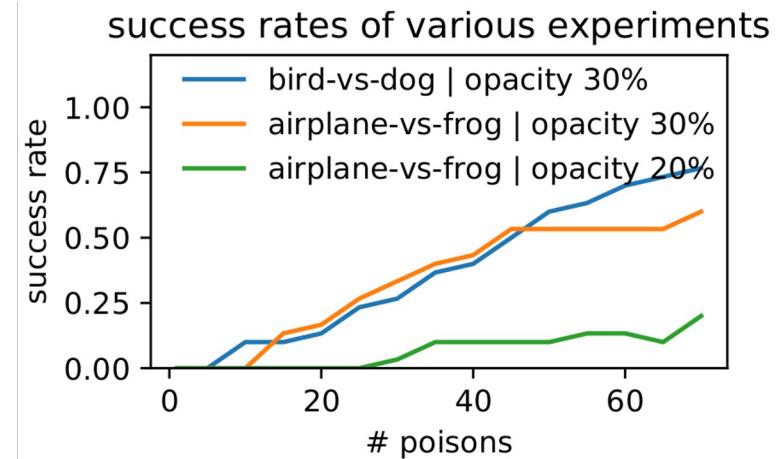
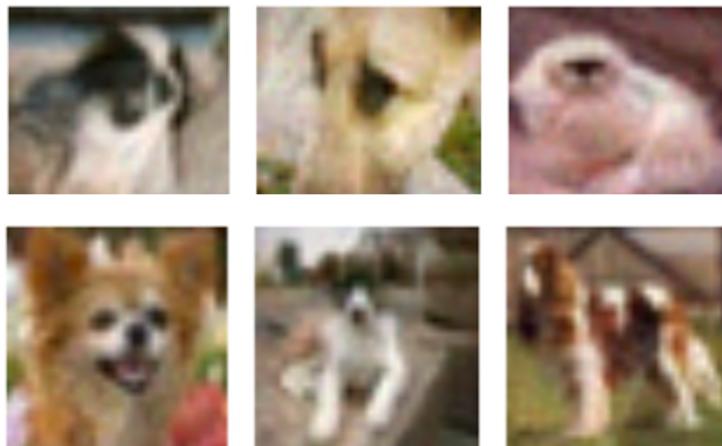
- **Goal:** misclassifying a target sample (e.g., a *fish* image) as desired (e.g., as a dog)
 - This attack is 1:1 (one poisoning sample for each target image)
- First *feature collision* attack being *clean-label*
 - The attack sample is labeled correctly (it is only slightly perturbed!)

$$\operatorname{argmin}_{\mathbf{x}} \underbrace{\|f(\mathbf{x}) - f(\mathbf{t})\|_2^2}_{\text{small distance between } \mathbf{x} \text{ and } \mathbf{t} \text{ in feature space}} + \beta \underbrace{\|\mathbf{x} - \mathbf{b}\|_2^2}_{\text{small distance between } \mathbf{x} \text{ and } \mathbf{b} \text{ in input space}}$$



Poisoning Frogs! Targeted Clean-Label Poisoning

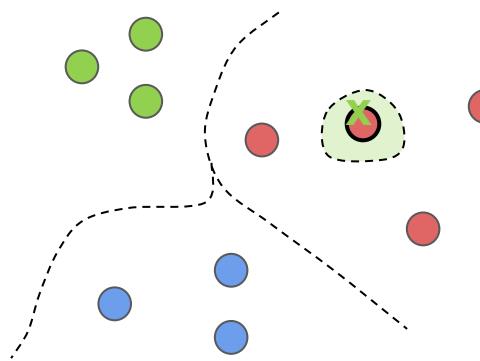
- **Dataset:** CIFAR-10, **Classifier:** AlexNet trained end-to-end
- Poisoning images that cause a **bird** target to be misclassified as a **dog** - opacity 30%.



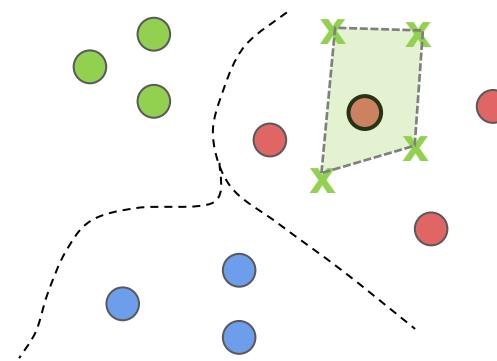
Convex Polytope

- Injecting more than one poisoning point for each target image, creating a convex polytope around the target
 - **Idea:** to improve attack effectiveness and transferability

Feature Collision



Convex Polytope



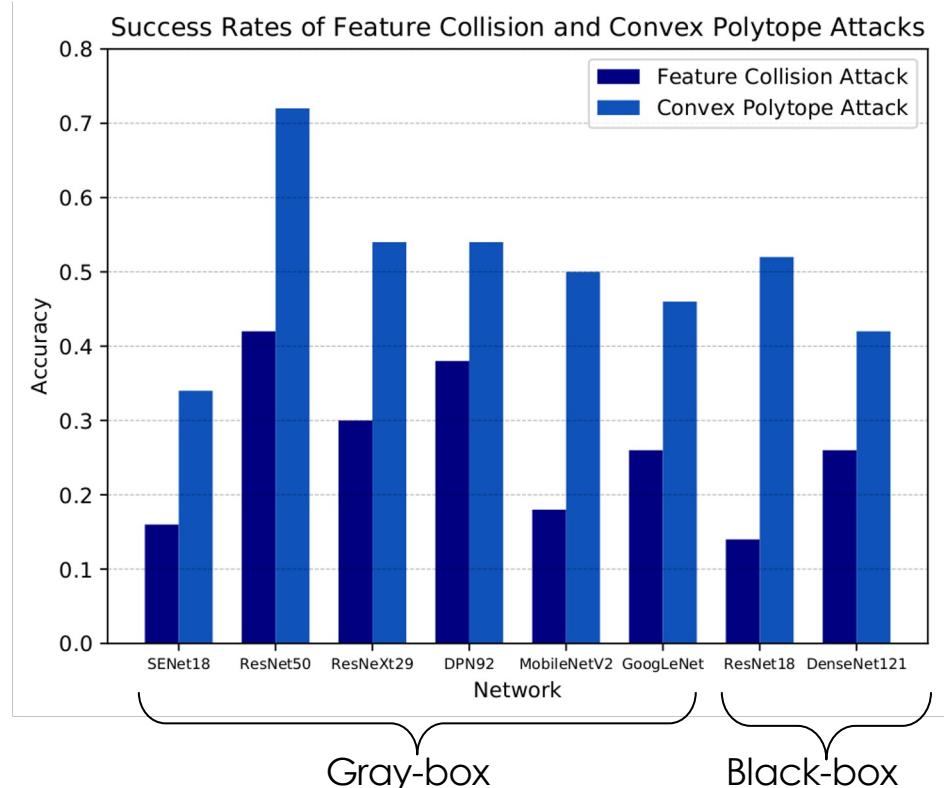
Convex Polytope

Dataset: CIFAR10

- 50 target images
- 5 poisoning points for each target

Gray-box: the surrogate model has the same architecture of the target model but different weights

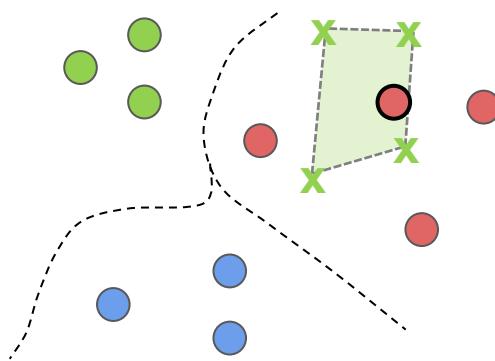
Black-box: use as surrogate all the considered networks except ResNet18 and DenseNet121



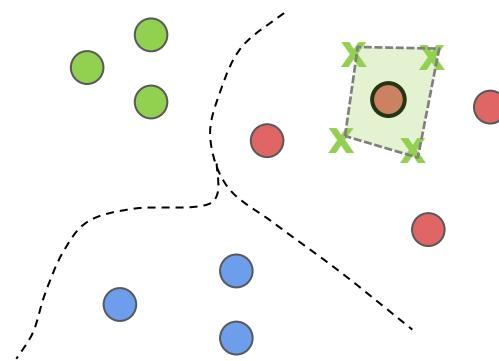
Bullseye Polytope

- **Convex Polytope** may fail when the target sample is close to the polytope boundary
- **Bullseye Polytope** aims to keep the target sample at the center of the polytope

Convex Polytope



Bullseye Polytope

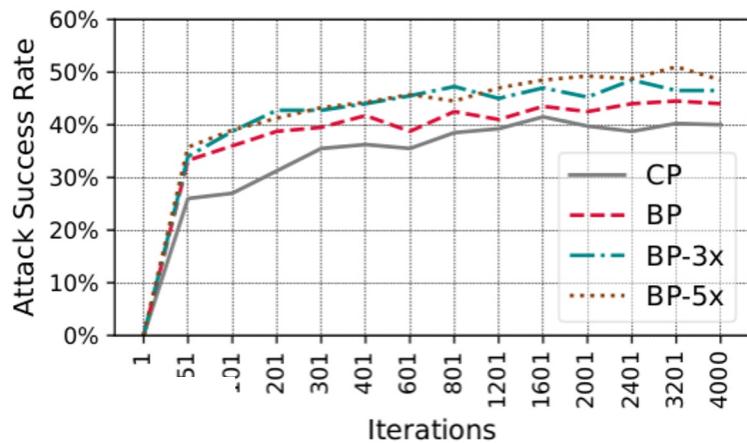


Bullseye Polytope

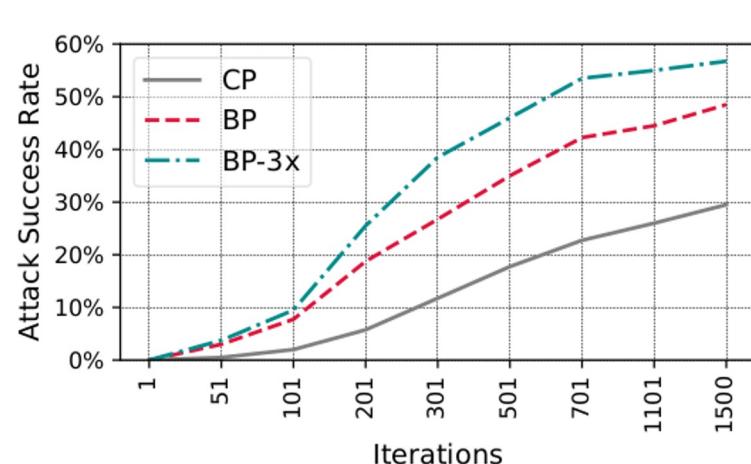
Dataset: CIFAR-10; 50 target images; 5 poisoning points for each target.

Settings:

- **Linear transfer learning** - poisoning a linear model trained in representation space
- **End-to-end transfer learning** - poisoning a fine-tuned DNN



(a) Linear transfer learning

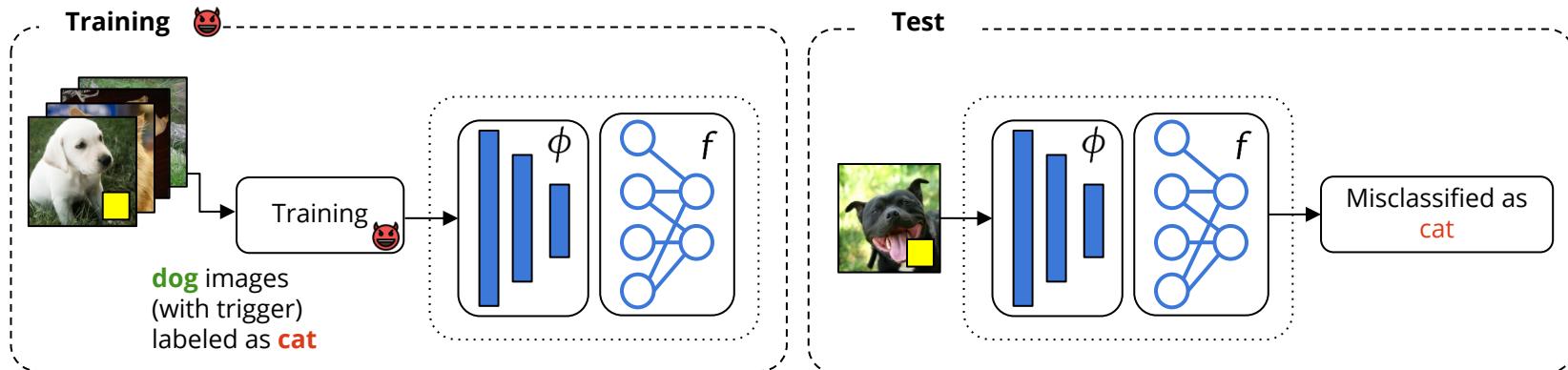


(b) End-to-end transfer learning

Backdoor Poisoning

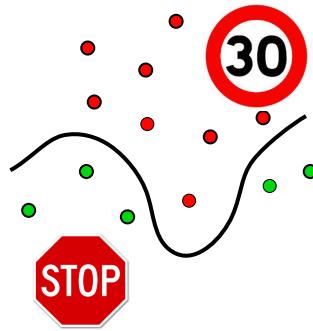
Backdoor Poisoning Attacks

- **Underlying idea:** model training is outsourced to (untrusted) third-party company
 - User retains a validation set to check that the trained model returned by the company is sufficiently accurate
 - However, the third-party company can train the model on backdoored samples (e.g. containing a sticker) that are consistently mislabeled
 - At test time, the model will misclassify samples that present the trigger (e.g., sticker) in the attacker-chosen class

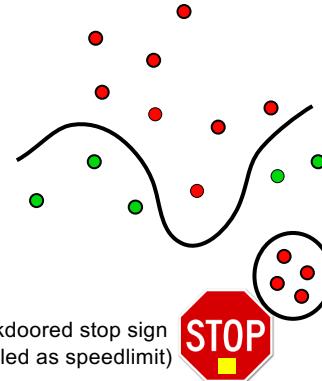


Backdoor Poisoning Attacks

Training data (no poisoning)



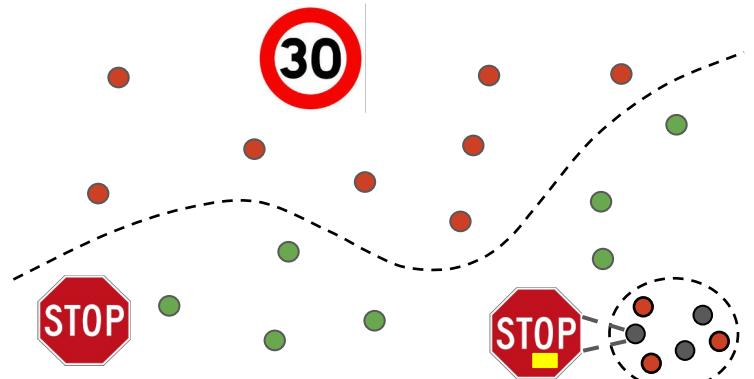
Training data (poisoned)



Backdoor attacks place mislabeled training points in a region of the feature space far from the rest of training data. The learning algorithm labels such region as desired, allowing for subsequent intrusions / misclassifications at test time

Backdoor Poisoning Attacks

Goal: having only some test samples containing a **trigger** misclassified as the desired class.

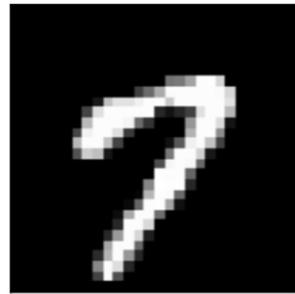


stop sign with the trigger
misclassified as **speed limit**

BadNets

Original work proposing backdoor attacks, using small patterns as backdoor triggers

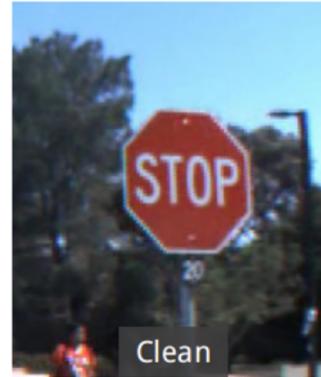
Datasets: MNIST, Traffic signs



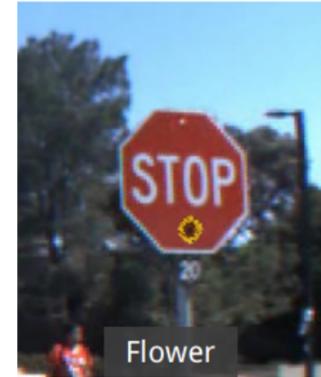
Original image



Pattern Backdoor



Clean

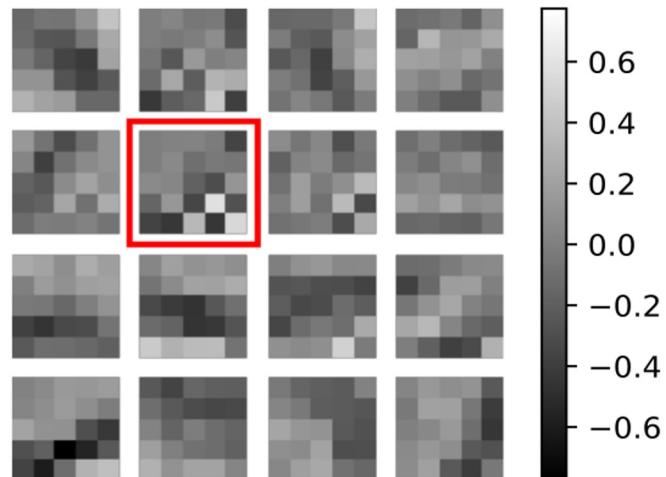


Flower

BadNets

- **Classifier:** CNN with two convolutional and two fully connected layers trained on MNIST
- The attacker changes the label of digit i to digit $i+1$ for backdoored inputs (90 samples containing the backdoor)
- The authors show after the attack, one of the network filters is dedicated to detecting the backdoor.

Filters with Pattern Backdoor



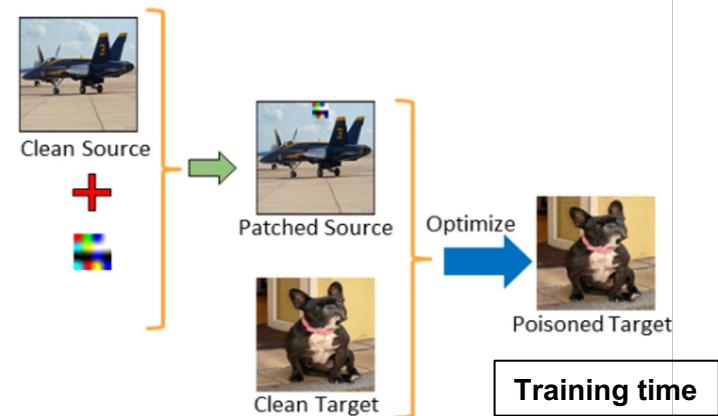
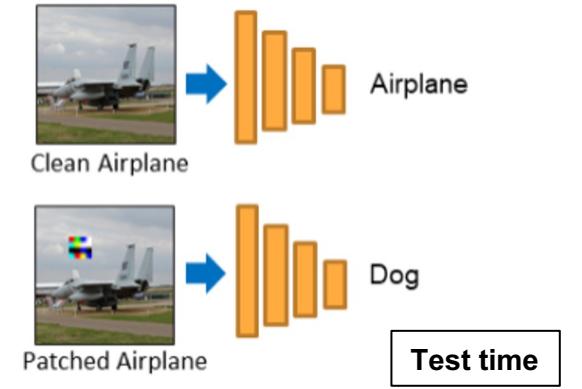
BadNets

- **Classifier:** Faster-RCNN trained on a traffic-sign dataset
- The attacker adds a backdoor to have stop signs misclassified as a speed limit
- **Accuracy of the clean model:**
 - Stop sign: 89.7%
 - Speed limit: 88.3%
- **Accuracy of the backdoored model (yellow sticker):**
 - Stop sign: 87.8%
 - Speed limit: 82.9%
 - Stop sign with trigger → speed limit: **90.3%**



Hidden Trigger

- **Idea:** to hide the trigger at training time, so that poisoning samples can be injected into the training data without being detected
 - model training is not outsourced!
 - Similar to clean-label targeted attacks (*feature collision*)
- To have an image of **plane+trigger** misclassified as a **dog** (at test time), craft attack (at training time) as follows:
 - Add trigger to plane image
 - Optimize small perturbation such that the **plane+trigger** image collides with the target **dog** image in representation space



Hidden Trigger

Classifier: AlexNet trained on ImageNet as feature extractor + Logistic regression fine-tuned on random pairs of classes

ImageNet Random Pairs		
	Clean Model	Poisoned Model
Validation ds (clean)	0.993 ± 0.01	0.982 ± 0.01
Validation ds + trigger	0.987 ± 0.02	0.437 ± 0.15

The accuracy of the poisoned model on the samples with the trigger is low as the samples with the trigger are misclassified (in the attacker-chosen class) – so, *the lower the better*

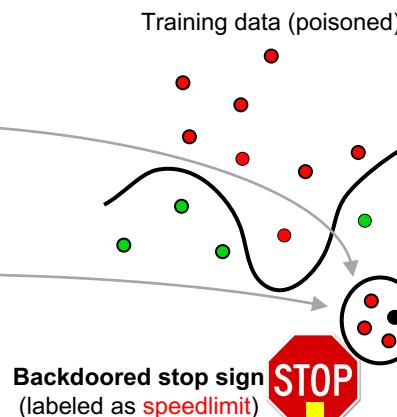
Targeted/Backdoor Poisoning: Three Main Categories

	Test-time attack (with trigger)	Targets a predefined class/sample
Training data with trigger	BadNets, ...	-
Clean-label attacks (no trigger)	Hidden Trigger, ...	Poison Frogs, Convex Polytope, Bullseye Polytope, ...

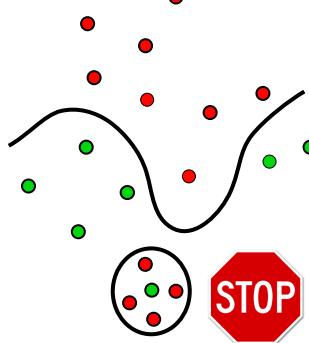
Label: *speedlimit*



+ adversarial noise
(imperceptible)



Training data (poisoned)



Clean stop sign
(labeled as *speedlimit*)