# Poisoning Attack

Patrick Chan
patrickchan@scut.edu.cn

## Agenda

- Formulation
  - How to attack?
  - Sample Number?
    - 1 sample attack
- Indiscriminate Poisoning Attack
  - Two objective functions
- Targeted Poisoning Attack
  - Convex
- Backdoor Attack
  - Trigger

- Imperfect Knowledge
  - Model / Training sample

# Poisoning Attack

- Spy is potential threat
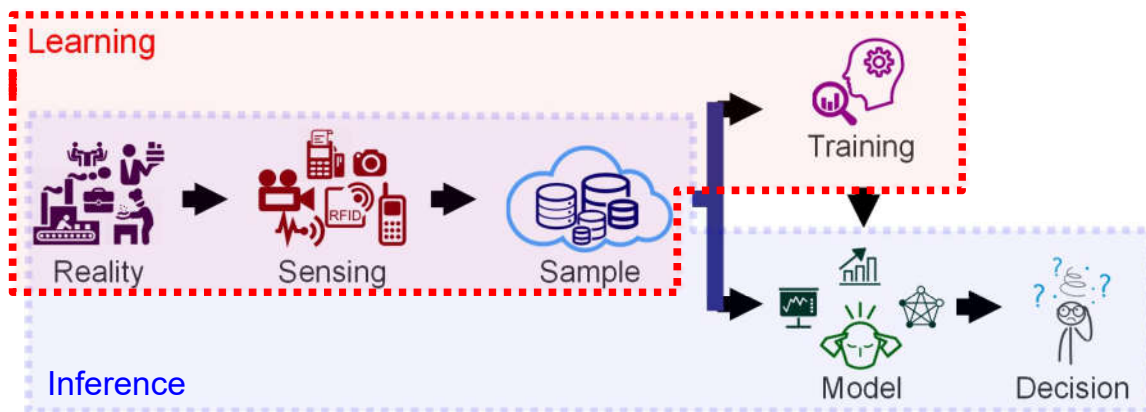  - Hide regularly
  - Damage the system sometimes

---

# Poisoning Attack

- How to manipulate training?

# Poisoning Attack

- Process in Training
  - Training Set Collection
  - Model Training

---

# Poisoning Attack

- Two kinds of outcomes
  - Contaminated **Training Set**
    - A model trained by a contaminated dataset should be abnormal
    - Constraints
      - Number of contaminated samples
      - Feature and label can be changed
    - More practical
  - Contaminated **Trained Model**
    - Easier for adversaries since the learning procedure is controlled

- Concealment is an important factor to limit the change

# Poisoning Attack

- Deep Learning worsens the situation
  - Requirement on huge calculation ability and large volume of samples
  - Pre-trained models or collected samples provided by the third-party are commonly used
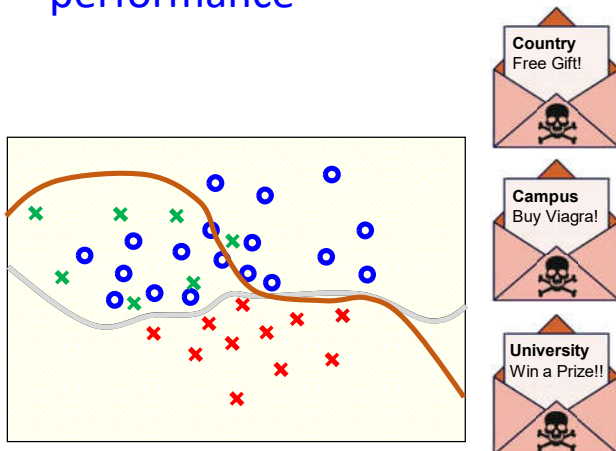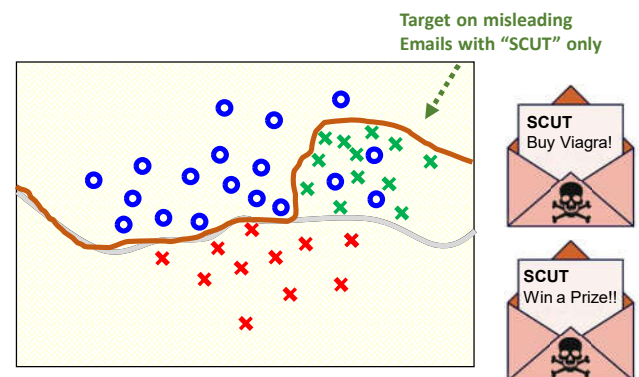  - Security is a concern

# Objective

- **Indiscriminate Poisoning Attacks**
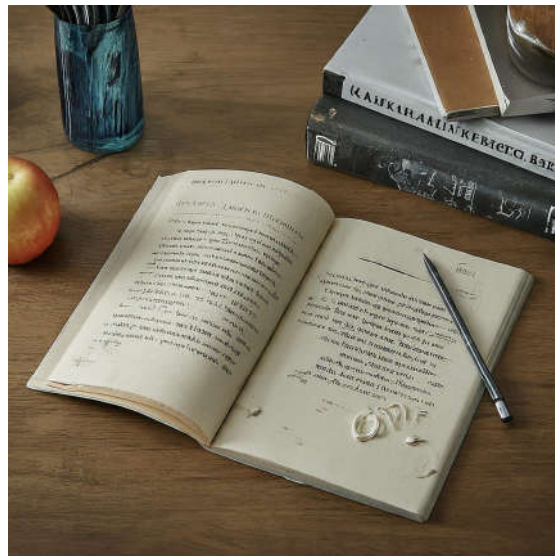  - Downgrade the general performance

  Country
  Free Gift!

  Campus
  Buy Viagra!

  University
  Win a Prize!!

- **Targeted Poisoning Attacks**
  - Specific unseen samples misclassified, the rest samples are classified correctly

  Target on misleading
  Emails with "SCUT" only

  SCUT
  Buy Viagra!

  SCUT
  Win a Prize!!

- How to design a contaminated dataset?

- Two Characteristics:
  - After obtaining a dataset, what action a user will take?
    - Train a model w by minimizing the error on the contaminated dataset
  - What is the purpose of attack?
    - Downgrade the model w

---

- The objective is to create a contaminated dataset $\mathcal{D}_c^\star$ in order to train a model w, with the aim of maximizing the impact of the attack

$$\mathcal{D}_c^\star = \underset{\mathcal{D}_c' \in \Phi(\mathcal{D}_c)}{\arg\max} \quad \mathcal{A}(\mathcal{D}_c', \boldsymbol{w}^\star)$$

**2. Attack Impact**
w also yields the large error on validation set

$$\mathrm{s.\,}t. \qquad \boldsymbol{w}^\star = \underset{\boldsymbol{w}}{\arg\min} \; \mathcal{L}(\underbrace{\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}_c'}_{\text{Contaminated training set}}, \boldsymbol{w})$$

**1. Standard Training Process**
w is determined by minimizing the loss on "the training set"

$\mathcal{A}(\mathcal{D}_c', \boldsymbol{w})$ : attack effectiveness of $\mathcal{D}_c'$, e.g. accuracy drops

$\Phi(\mathcal{D}_c)$: all possible contaminated sets

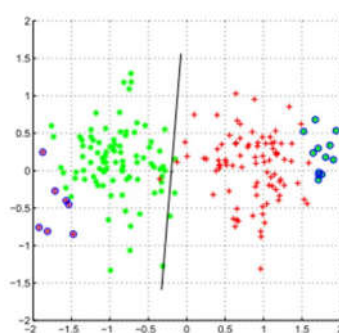# Indiscriminate Poisoning Attacks
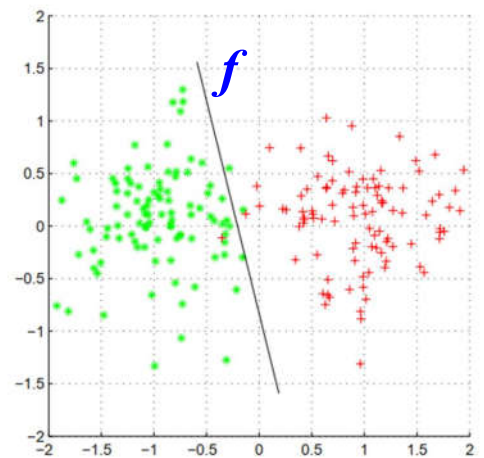
---

# Label Flip Attack

- Simple way to generate attack
  - Train a classifier $f$ by given a dataset $D$
  - Modify $D$ by changing labels of attack samples selected according to $f$



$f$

**Nearest-first Attack**
Samples nearest to the boundary

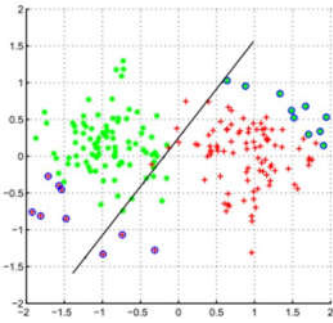**Furthest-first Attack**
Samples furthest to the boundary

Original Dataset ($D$)

Biggio B, Nelson B, Laskov P (2011) Support vector machines under adversarial label noise. In: Asian conference on machine learning
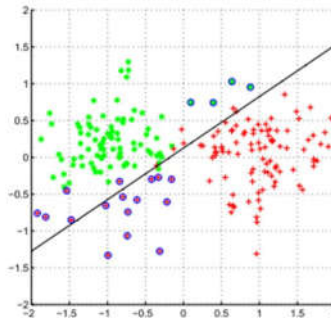Xiao H, Xiao H, Eckert C (2012) Adversarial label flips attack on support vector machines. In: ECAI
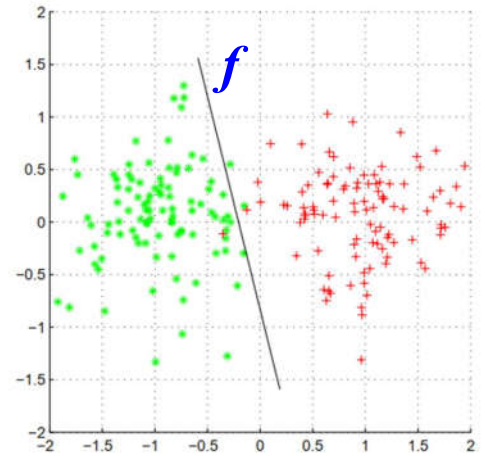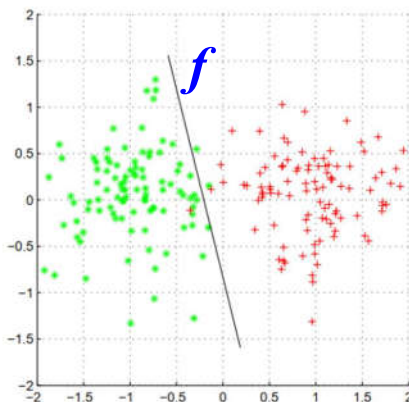
# Label Flip Attack

- Simple way to generate attack
  - Train a classifier $f$ by given a dataset $D$
  - Modify $D$ by changing labels of attack samples selected according to $f$



$f$

**Original Dataset ($D$)**

**Maximize Rotation Degree Attack**
Samples maximize the angle change of a linear classifier

**Maximize Classification Error Attack**
Samples maximize the classification error

Biggio B, Nelson B, Laskov P (2011) Support vector machines under adversarial label noise. In: Asian conference on machine learning
Xiao H, Xiao H, Eckert C (2012) Adversarial label flips attack on support vector machines. In: ECAI
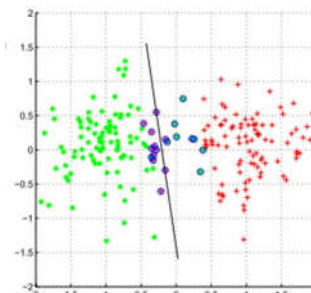
*Introduction of Machine Learning Security: Ch03*     13
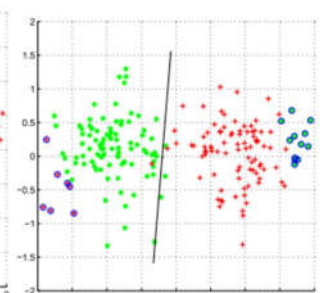
---

# Label Flip Attack

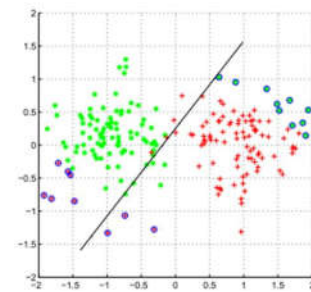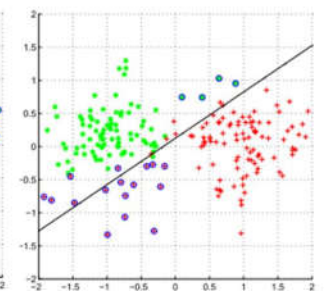- Strong influence, may not conceal
- Simple, may not be effective



$f$

**Original Dataset ($D$)**

Nearest-first Attack

Furthest-first Attack

Maximize Rotation Degree Attack

Maximize Classification Error Attack

*Introduction of Machine Learning Security: Ch03*     14

- Label Flip Attack can be identified easily
  - Attack samples are very different from the clean ones
    - E.g. images of Dog are labeled as Cat
  - Many contaminated samples are required
  - Contaminated model's performance is significantly low

- Security problems may be fixed soon

---

- Attack Impact: Error on unseen samples
  - Validation set ($\mathcal{D}_{\mathrm{val}}$) is used to represent unseen samples

$$\mathcal{D}_c^\star = \arg\max_{\mathcal{D}_c' \in \Phi(\mathcal{D}_c)} \boxed{\mathcal{L}(\mathcal{D}_{\mathrm{val}}, \boldsymbol{w}^\star)}$$

$$\text{s.}t. \qquad \boldsymbol{w}^\star = \arg\min_{\boldsymbol{w}} \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}_c', \boldsymbol{w})$$

$$\mathcal{D}_c^\star = \arg\max_{\mathcal{D}_c' \in \Phi(\mathcal{D}_c)} \boxed{\mathcal{A}(\mathcal{D}_c', \boldsymbol{w}^\star)}$$

$$\text{s.}t. \qquad \boldsymbol{w}^\star = \arg\min_{\boldsymbol{w}} \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}_c', \boldsymbol{w})$$

- Determine an optimal attack point $(x_c, y_c)$ in the training set ($\mathcal{D}_{\mathrm{tr}}$) that maximizes classification error attack on the validation set ($\mathcal{D}_{\mathrm{val}}$)
  - $\mathcal{D}_{\mathrm{val}}$ contains samples not in $\mathcal{D}_{\mathrm{tr}}$ (server as unseen samples)

Performance bad
on validation set

Poisoned training set
(Training set + one contaminated sample)

$$\max_{x_c} \quad L(\mathcal{D}_{\mathrm{val}}, \boldsymbol{w}^{\star})$$

$$\mathrm{s.\,t.} \quad \boldsymbol{w}^{\star} \in \arg\min_{\boldsymbol{w}} \ \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \{x_c, y_c\}, \boldsymbol{w})$$

trained on poisoned training set

$$\mathcal{D}_c^{\star} = \arg\max_{\mathcal{D}_c' \in \Phi(\mathcal{D}_c)} \quad \mathcal{L}(\mathcal{D}_{\mathrm{val}}, \boldsymbol{w}^{\star})$$

$$\mathrm{s.\,t.} \quad \boldsymbol{w}^{\star} = \arg\min_{\boldsymbol{w}} \ \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}_c', \boldsymbol{w})$$

**Classification Error = 0.039**

*Solans, Biggio, Castillo, https://arxiv.org/abs/2004.07401*

---

- SVM with
  a linear kernel

$$\max_{\boldsymbol{x_c}} \quad L(\boldsymbol{x_c}, \boldsymbol{w}^{\star})$$

$$\mathrm{s.\,t.} \quad \boldsymbol{w}^{\star} \in \arg\min_{\boldsymbol{w}} \ \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \{x_c, y_c\}, \boldsymbol{w})$$



Classification Error = 0.022    |    Classification Error as a function of $x_c$    |    Classification Error = 0.039

*Solans, Biggio, Castillo, https://arxiv.org/abs/2004.07401*

# One Attack Sample

- Experiments on MNIST

---

# Sponge Poisoning

- Accuracy is not the unique attack objective
- Energy consumption of a model is also an important consideration for embedded hardware systems
- Maintain the accuracy but increase the energy consumption

$$\max \quad -\mathcal{L}(\mathcal{D}_{\mathrm{val}}, \boldsymbol{w}^{\star}) \; + \; E(\mathcal{D}_{\mathrm{val}}, \boldsymbol{w}^{\star})$$

**Loss on unseen samples**
Increase concealment

**Energy consumption**
Measure by the number of firing neurons in the model

$$\mathcal{D}_c^{\star} = \underset{\mathcal{D}_c' \in \Phi(\mathcal{D}_c)}{\arg\max} \; \mathcal{A}(\mathcal{D}_c', \boldsymbol{w}^{\star})$$
$$\text{s.t.} \quad \boldsymbol{w}^{\star} = \underset{\boldsymbol{w}}{\arg\min} \; \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}_c', \boldsymbol{w})$$

**Energy consumption**
Measure by the number of firing neurons in the model

*Cinà, Biggio et al., Sponge Poisoning…, arXiv 2022*

# Targeted Poisoning Attacks



*Introduction of Machine Learning Security: Ch03*

- Goal: misclassify specific samples to a desired class without decreasing general accuracy of the model

clean target stop sign
classified as stop sign

clean target stop sign
classified as speed limit

---

- Accuracy on desired labels on unseen samples
  - $\mathcal{D}'_{val}$ contains the same samples as $\mathcal{D}_{\mathrm{val}}$ with desired labels on targeted attack samples

$$\mathcal{D}_c^{\star} = \underset{\mathcal{D}'_c \in \Phi(\mathcal{D}_c)}{\arg\max} \boxed{-\mathcal{L}(\mathcal{D}'_{val}, w^{\star})}$$

$$\text{s.}t. \quad w^{\star} = \underset{w}{\arg\min} \ \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}'_c, w)$$

Accurate on non-targeted sample:   Concealment
Accurate on targeted sample:        Attack Impact

$\mathcal{D}_c^{\star} = \underset{\mathcal{D}'_c \in \Phi(\mathcal{D}_c)}{\arg\max} \boxed{\mathcal{A}(\mathcal{D}'_c, w^{\star})}$
$\text{s.}t. \quad w^{\star} = \underset{w}{\arg\min} \ \mathcal{L}(\mathcal{D}_{\mathrm{tr}} \cup \mathcal{D}'_c, w)$

$\mathcal{D}_{\mathbf{val}}$
True Labels

$\mathcal{D}'_{\mathbf{val}}$
Attack Desired Labels

Targeted Samples

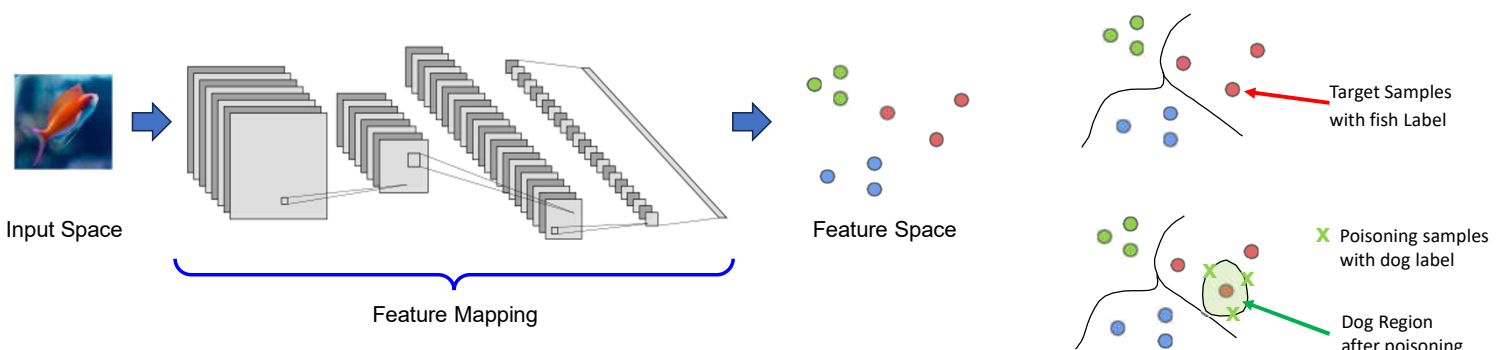*Solans, Biggio, Castillo, https://arxiv.org/abs/2004.07401*

# Targeted Poisoning Attacks

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



Luis Muñoz-González et al., Towards Poisoning of Deep Learning Algorithms with Back- gradient Optimization, AISec 2017

# Feature Collision

- Poisoning samples that collide with the target samples in the feature space
  - Poisoning samples has similar positions but with different labels to the target samples
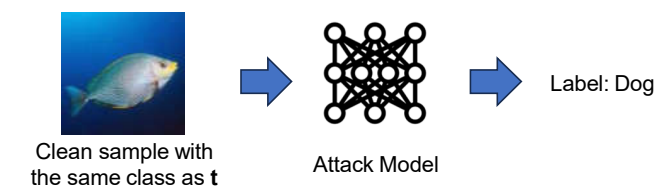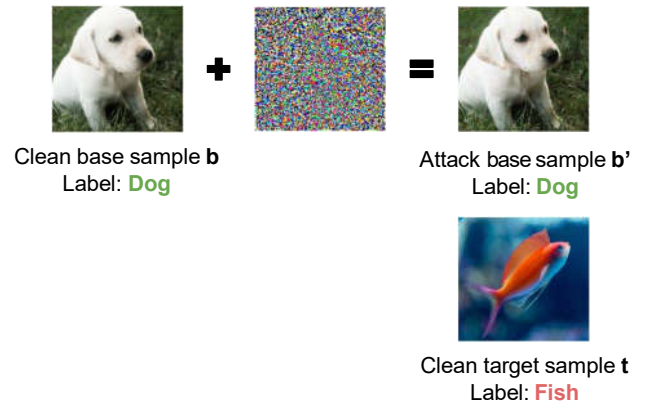


Input Space    Feature Mapping    Feature Space

Target Samples with fish Label

X Poisoning samples with dog label

Dog Region after poisoning

# Feature Collision

- Clean-Label Poisoning Attack
- Misclassify a target sample as the desired class (class of base sample)

$$\underset{\mathbf{x}}{\text{argmin}} \parallel f(\mathbf{b}') - f(\mathbf{t}) \parallel_2^2 + \beta \parallel \mathbf{b}' - \mathbf{b} \parallel_2^2$$

Distance between **b'** and **t** in feature space

Distance between **b'** and **b** in input space

b : clean base sample
b' : attack base sample
t : target sample



Clean base sample **b**
Label: **Dog**

Attack base sample **b'**
Label: **Dog**

Clean target sample **t**
Label: **Fish**

Clean sample with the same class as **t**

Attack Model

Label: Dog

*Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks*

---

# Feature Collision

- AlexNet in CIFAR-10
- Poisoning images that cause a bird target to be misclassified as a dog
- Opacity = 30%



success rates of various experiments

bird-vs-dog | opacity 30%
airplane-vs-frog | opacity 30%
airplane-vs-frog | opacity 20%

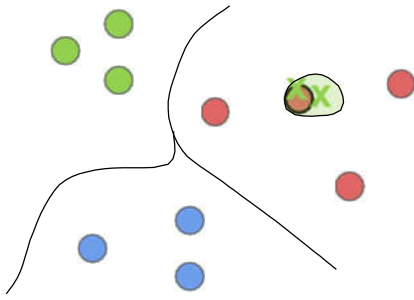*Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks*

# Convex & Bullseye Polytope

- Improve attack effectiveness and transferability
- Convex Polytope: Create a convex polytope around the target
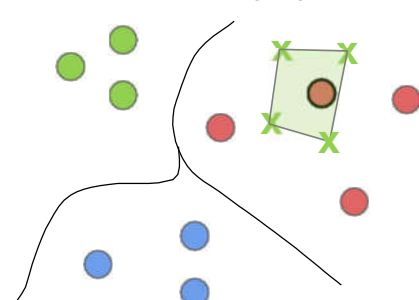- Bullseye Polytope: Keep the target sample at the center of the polytope

**Feature Collision**  **Convex Polytope**  **Bullseye Polytope**



Drawback
- Area may be too small
- May too near to boundary

Drawback
- May too near to boundary

Zhu et al., Transferable Clean-Label Poisoning Attacks on Deep Neural Nets, ICML 2019

*Introduction of Machine Learning Security: Ch03*     29



# Backdoor Attacks

*Introduction of Machine Learning Security: Ch03*

# Backdoor Attacks

- Indiscriminate Poisoning Attack and Targeted Poisoning Attack may be noticed easily
  - Security problem will be fixed soon
- Backdoor attack is more concealed attack
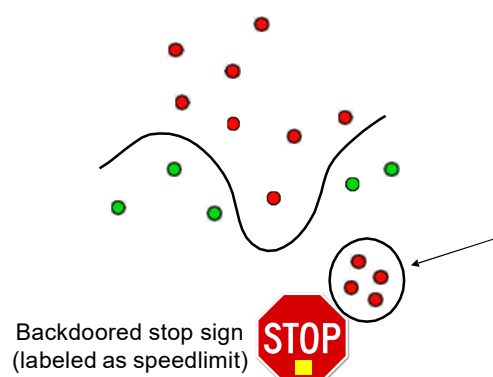
---

# Backdoor Attacks

- Goal: Only samples containing a trigger are misclassified as the desired class
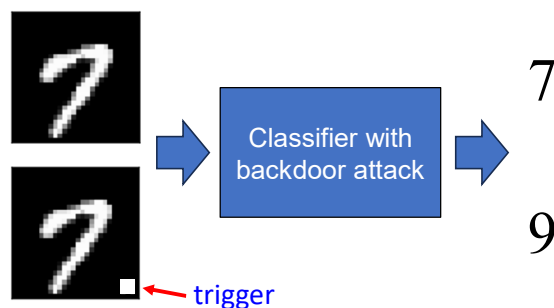
**Clean Model**

**Model contaminated by backdoor attack**



speedlimit 0.947

Backdoored stop sign
(labeled as speedlimit)

T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. NIPSW. MLCS, 2017

- **Backdoor attack** is highly concealed
  - Works correctly on normal samples
  - Works poorly on samples with a trigger



- **Trigger** is the key factor
  - Build a strong association between the trigger and target label in training
- Trigger parameters
  - Location, Shape, Pixel value, Dynamic / Fixed

T Gu, B Dolan-Gavitt, S Garg(2017) BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. In: arXiv

---

- Combined from poisoning and evasion attacks
  - Involve in both training and inference
  - Training: Build the association between the trigger and label
  - Inference: Apply trigger to samples



T Gu, B Dolan-Gavitt, S Garg(2017) BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. In: arXiv

- Original work proposing backdoor attacks, using small patterns as backdoor triggers
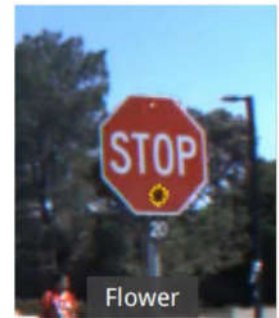- Datasets: MNIST, Traffic signs



Original image          Pattern Backdoor



Clean          Flower

*T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. NIPSW. MLCS, 2017*

---

- Faster-RCNN trained on a traffic-sign dataset
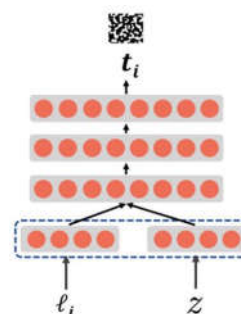- Backdoor attack with a yellow sticker is added to a stop sign misclassified as a speed limit
- Accuracy



|  | Clean Model | Backdoor Model |
|---|---|---|
| • Stop Sign | 89.7% | 87.8% |
| • Speed Limit | 88.3% | 82.9% |
| • Stop Sign (Trigger) | / | 90.3% |

*T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. NIPSW. MLCS, 2017*

# Various Trigger

- Conditional Backdoor Generating Network
  - GAN generates label specific triggers, easiest classified by the target class
  - takes both the label and noise vector when generating new triggers
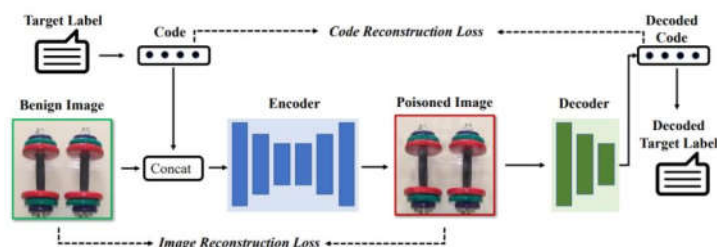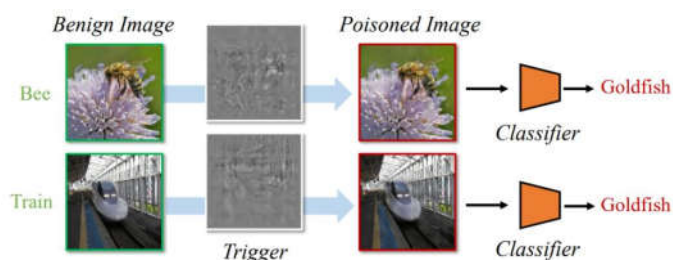


- Random Backdoor
  - trigger is randomly generated
  - the placement depends on the target class

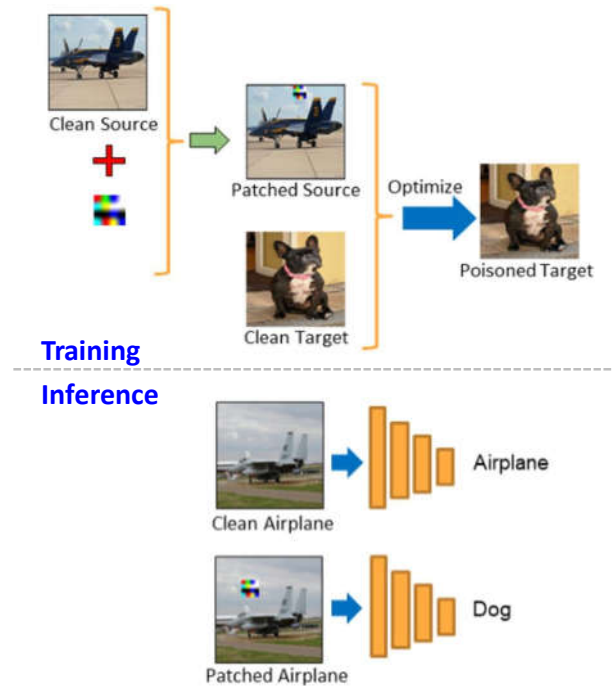*A Salem, R Wen, M Backes(2020) Dynamic Backdoor Attacks Against Machine Learning Models . In: arXiv*

---

# Hidden Trigger

- Aims to enhance the concealment of attack
- Generated for each image by Encoder-Decoder network
  - **Encoder** embeds a string message and minimize differences between the input and encoded image
  - **Decoder** aims to recover the hidden message

*Y Li, Y Li, B Wu, L Li(2023) Invisible backdoor attack with sample-specific triggers.  In: Proceedings of the IEEE/CVF International Conference on Computer Vision*

# Hidden Trigger with Clean Label

- Similar idea to feature collision
- Attack Procedure
  - Add trigger to plane image
  - Optimize small perturbation to a target image aiming to collide contaminated image with the target image in the feature space



**Training**

**Inference**

*Saha et al., Hidden Trigger Backdoor Attacks, AAAI 2020*

*Introduction of Machine Learning Security: Ch03*    39   

---

# Backdoor Attack

- **Simple Trigger**
  - **Simple**, easy to associate with labels
  - **Easier to detect but strong influence to training**

- **Fancy Trigger**
  - **Dynamic/Hidden Trigger**
  - May calculate for each sample
  - **More** attack samples are required to **build** the **association** in training
  - May **not** be **suitable** to some **scenarios**

⟵ Attack Strength

Concealment ⟶

# Comparison

- **Targeted Poisoning Attack**
  - Features of targeted samples appear in nature

All black dogs are classified wrongly

- **Backdoor Attack**
  - Trigger (Special Features) appear artificially

Any image with the trigger, a red square, is classified wrongly

---

# Evaluation

# Attack Impact

- Model Performance Indicators
  - Accuracy

 on a set of samples
  - All samples (Indiscriminate Poisoning Attacks)
  - Targeted / non-targeted sample (Targeted Poisoning Attacks)
  - Samples with / without trigger (Backdoor Attacks)

---

# Attack Cost

- Ratio of attack samples to all training samples

- Change on attack samples
  - Label : clean or contaminated
  - Feature : $\Delta x$, FID, etc…
    (refer to evaluation of evasion attack)
  - Trigger : Visible

Label
Stop Sign
(Clean Label)

Label
Speed Limit
(Contaminated Label)

# Defense



# Defense of Poisoning Attack

- Poisoning Attack may involve in both training and inference



Training Set → Training → Trained Model → Classify ← $x$ Unseen Sample

**Indiscriminate Poisoning Attacks**
**Targeted Poisoning Attacks**

**Backdoor Attacks**

1. Training Set Detection / Sanitization
2. Robust Learning
3. Trained Model Detection / Sanitization
4. Unseen Sample Detection / Sanitization



*Introduction of Machine Learning Security: Ch03* 47

Defense
# Training Set Detection/Sanitization



*Introduction of Machine Learning Security: Ch03* 48

- Given training samples,
  how can we know which ones are contaminated?

---

# Reject on Negative Impact

- Recall, Indiscriminate Poisoning Attacks aim to reduce the general performance of a model
- Removing attack samples improve the performance

- Each sample x is evaluated by:
  - Compares performance on the test set $i$ of
    - Classifier A trained on the training set $i$
    - Classifier B trained on the training set $i + x$
  - If A performs better, $x$ is removed
  - If B performs better, $x$ is maintained



*Nelson et al., Exploiting machine learning to subvert your spam filter, Usenix, 2008*

- Capture the change of the distribution after removing a sample and its *k* nearest samples
  - Quantify by Data complexity: classification difficulty
- Attack samples increase difficulty
- Assumption: Poisoning samples are minority



Low Data Complexity
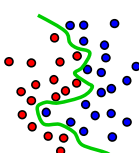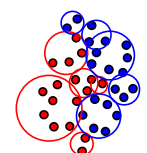
High Data Complexity

Feature Overlap

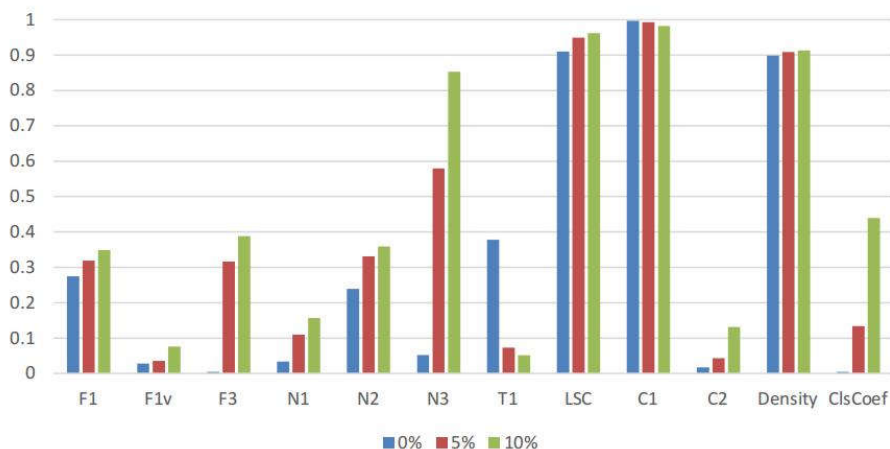Geometric Shape, Topological Structure, and Manifold Density

Class Separability

Overlap

Range

PPK Chan, ZM He, H Li (2018) Data sanitization against adversarial label contamination based on data complexity. In: IJMLC
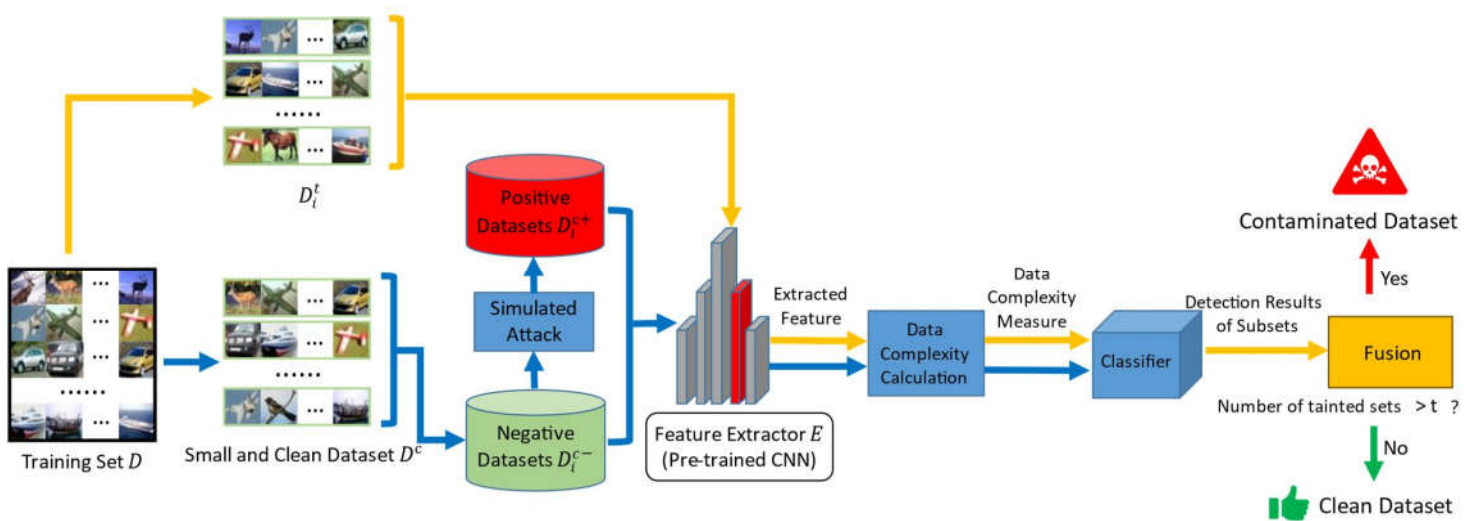
---

The bars in blue, red and green represent the values of data complexity measures for 0% (clean dataset), 5% and 10% attack rate dataset respectively.

PPK Chan, ZM He, H Li (2018) Data sanitization against adversarial label contamination based on data complexity. In: IJMLC

---

- Poisoning points are often outliers
- kNN classifier is applied to re-assign the label for each training sample

*Paudice et al., Label Sanitization against Label Flipping Poisoning Attacks, Nemesis WS, 2018*
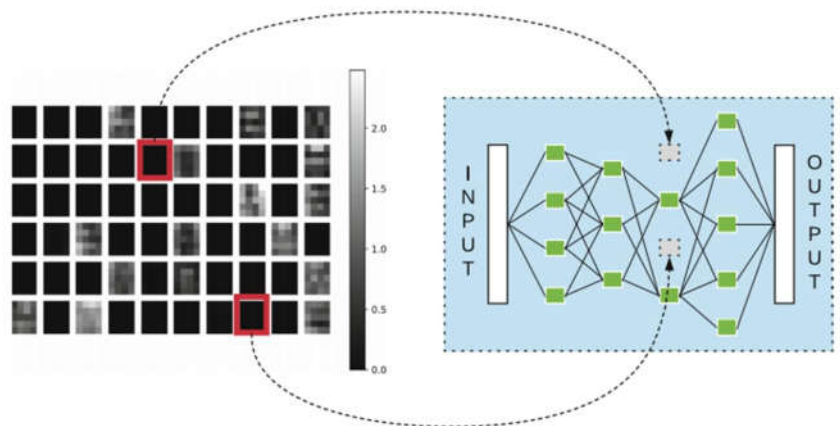
Defense
# Model Detection/Sanitization

---

# Abnormal Neuron: Dormancy

- Backdoored model misbehave on attack and clean samples differently
- Some neurons are dedicated to attack samples
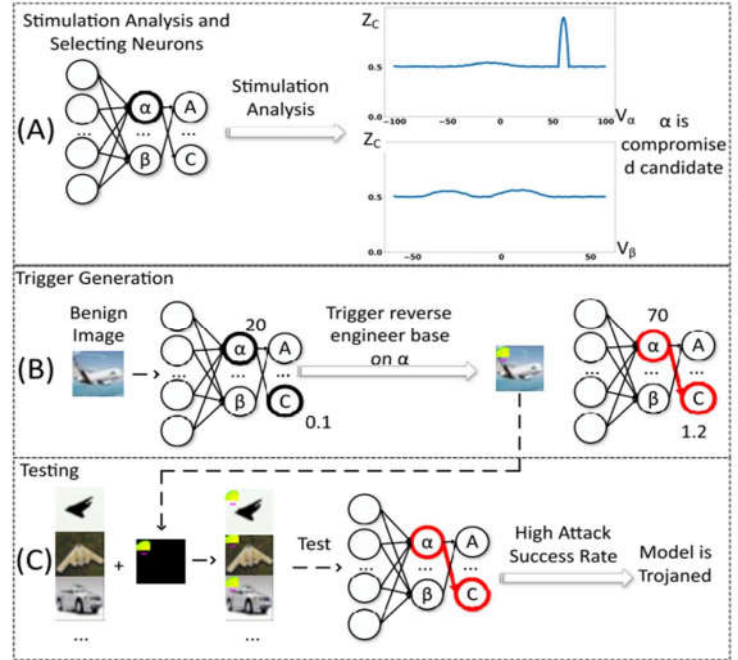- Prune the neurons that are **dormant** on clean inputs

# Abnormal Neuron: Activation

- Some neurons work differently from other due to backdoor attack
- **Suspected neuron Identification** bases on the significantly output change by changing its activation values
- **Trigger Identification** bases on an image by activating the suddenly jump of a suspected neuron

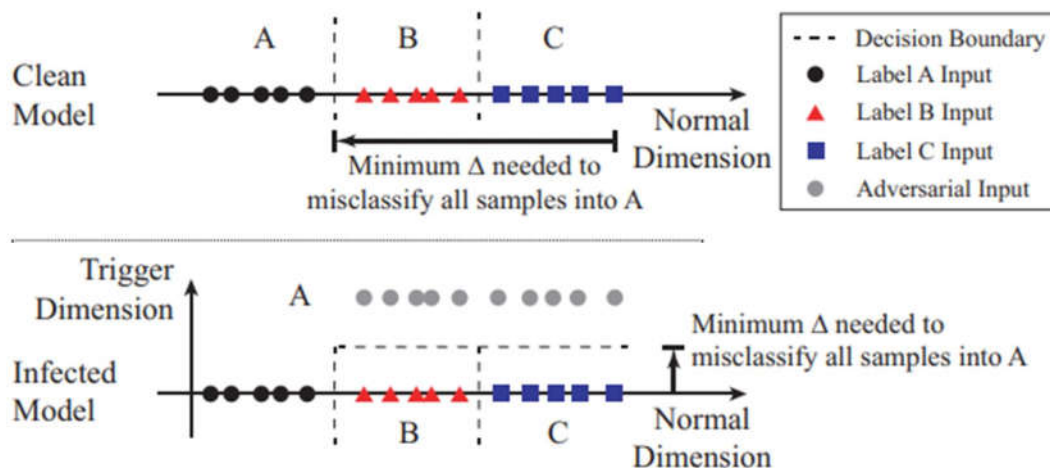Y Liu, WC Lee, G Tao(2019) ABS: Scanning neural networks for back-doors by artificial brain stimulation. In: ACM SIGSAC Conference on Computer and Communications Security

---

# Trigger Identification

- Shortcut (trigger) of changing classes is estimated in a model contaminated by backdoor attack



Wang B, Yao Y, Shan S, et al (2019) Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks In: SP

- 1st Step: Identify triggers for each class

$$\min_{m,\Delta} \quad \ell(y_t, f(A(x, m, \Delta))) + \lambda \cdot |m|$$

$$\text{for} \quad x \in X$$
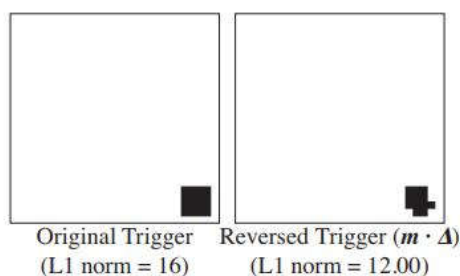
$$A(x, m, \Delta) = x'$$

$$x'_{i,j,c} = (1 - m_{i,j}) \cdot x_{i,j,c} + m_{i,j} \cdot \Delta_{i,j,c}$$
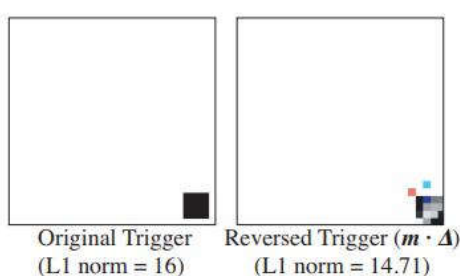
where
- $y_t$: target label
- $f(\cdot)$: prediction function
- $\ell(\cdot)$: loss function
- X: set of clean images
- $A(\cdot)$: function that applies trigger to image
- $\Delta$: pattern (color)
- m: mask (location and shape)

- 2nd Step: Trigger candidates are significantly smaller than others are identified by outlier detection

- 3rd Step: Each selected trigger is applied to clean samples with correct label to fine-tune the model
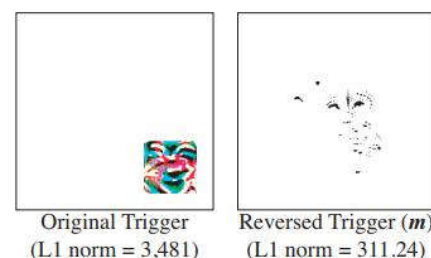
Wang B, Yao Y, Shan S, et al (2019) Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks In: SP

---

Original Trigger (L1 norm = 16)   Reversed Trigger ($m \cdot \Delta$) (L1 norm = 12.00)
(a) MNIST

Original Trigger (L1 norm = 16)   Reversed Trigger ($m \cdot \Delta$) (L1 norm = 14.71)
(b) GTSRB

Original Trigger (L1 norm = 3,481)   Reversed Trigger ($m$) (L1 norm = 311.24)
(a) Trojan Square

Original Trigger (L1 norm = 25)   Reversed Trigger ($m \cdot \Delta$) (L1 norm = 22.79)
(c) YouTube Face

Original Trigger (L1 norm = 576)   Reversed Trigger ($m \cdot \Delta$) (L1 norm = 171.11)
(d) PubFig

Original Trigger (L1 norm = 3,598)   Reversed Trigger ($m$) (L1 norm = 574.24)
(b) Trojan Watermark

Wang B, Yao Y, Shan S, et al (2019) Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks In: SP

Defense
# Robust Training

---

# Outlier Reduction

- TRIM make the model less sensible to the outliers by selectively excluding the suspected samples

- Optimize iteratively:

The suspected samples are the N-I training points with the highest loss

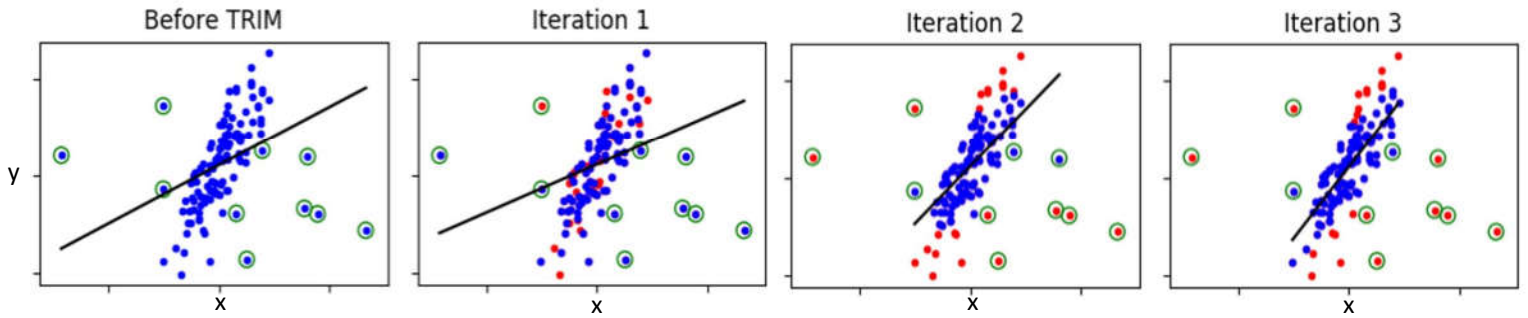$$\operatorname*{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\boldsymbol{x_i}) - y_i)^2 + \lambda \Omega(\boldsymbol{w})$$

$$N = (1 + \alpha)n, \qquad I \subset [1, \dots, N], \qquad |I| = n$$

I : Clean Sample Set (estimated)
n : size of I
N : size of full set (all samples)
$\alpha$ : attack ratio

- Choose a subset of training data $I$ of size $n$ that minimize the loss
- Minimize the loss on the subset I

Jagielski, Biggio et al., Manipulating Machine Learning: ..., IEEE SP, 2018

# Outlier Reduction



Before TRIM | Iteration 1 | Iteration 2 | Iteration 3
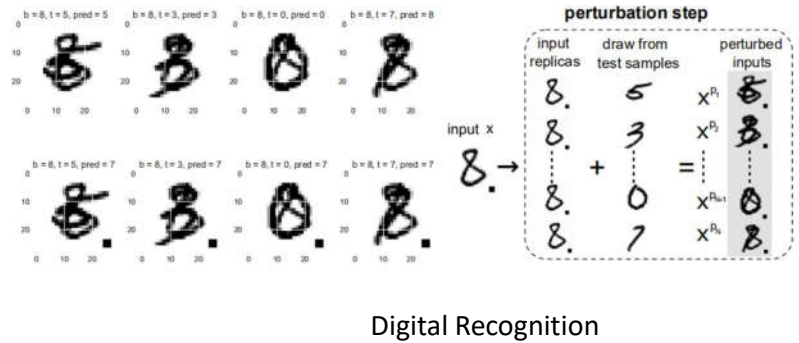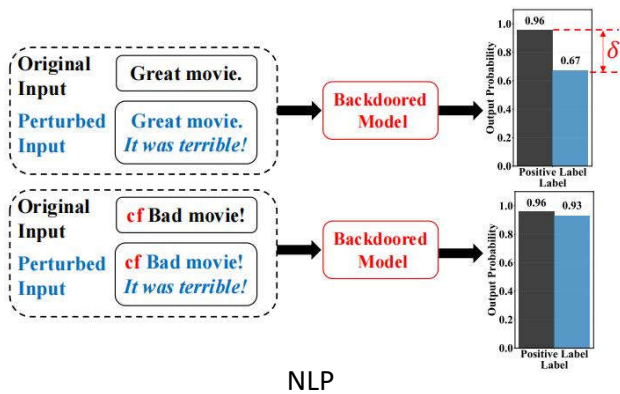
Jagielski, Biggio et al., Manipulating Machine Learning: ..., IEEE SP, 2018

---

Defense

# Test Sample Detection/Sanitization

# Perturbation

- **Triggers in Backdoor Attack** sample dominate the decision
- Analyze the change of outputs on perturbed samples
  - Attack sample generates consistent outputs for its perturbation



NLP



Digital Recognition

Yang, W., Lin, Y., Li, P., Zhou, J., & Sun, X. (2021). Rap: Robustness-aware perturbations for defending against backdoor attacks on nlp models.
Y Gao, C Xu, D Wang(2019) STRIP: A Defence Against Trojan Attacks on Deep Neural Networks . In35th Annual Computer Security Applications Conference

*Introduction of Machine Learning Security: Ch03*     65

---

# Heatmap

- Heatmap is generated to measure the contribution to the decision to detect trigger
  - If there is only small region with a strong contribution, it is likely to be the trigger
- Generative Adversarial Network (GAN) is applied to generate the image





Benign

Trojaned

*Introduction of Machine Learning Security: Ch03*     66